Forschungszentrum Jülich GmbH Institute of Complex Systems Jülich Centre for Neutron Science Peter Grünberg Institute Institute for Advanced Simulation

Lecture Notes of the 49th IFF Spring School 2018

Gerhard Gompper, Jan Dhont, Jens Elgeti, Christoph Fahlke, Dmitry Fedosov, Stephan Förster, Pavlik Lettinga, Andreas Offenhäusser (Eds.)

Physics of Life

This Spring School was organized by the Institute of Complex Systems of the Forschungszentrum Jülich on 26 February until 9 March 2018.

In collaboration with Universities and research institutions.

Schriften des Forschungszentrums Jülich Reihe Schlüsseltechnologien / Key Technologies

Band / Volume 158

ISSN 1866-1807

ISBN 978-3-95806-286-3

Bibliographic information published by the Deutsche Nationalbibliothek. The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at http://dnb.d-nb.de

Publisher:	Forschungszentrum Jülich GmbH IAS, ICS, JCNS, PGI 52425 Jülich Tel.: +49 2461 61-6048 Fax: +49 2461 61-2410
Cover Design:	Grafische Medien, Forschungszentrum Jülich GmbH
Printer:	Schloemer + Partner GmbH, Düren
Copyright:	Forschungszentrum Jülich 2018
Distributor:	Forschungszentrum Jülich GmbH Zentralbibliothek, Verlag 52425 Jülich Tel.: +49 2461 61-5368 Fax: +49 2461 61-6103 zb-publikation@fz-juelich.de www.fz-juelich.de/zb

Schriften des Forschungszentrums Jülich Reihe Schlüsseltechnologien / Key Technologies, Band / Volume 158

ISSN 1866-1807 ISBN 978-3-95806-286-3

The complete volume is freely available on the Internet on the Jülicher Open Access Server (JuSER) at www.fz-juelich.de/zb/openaccess.



This is an Open Access publication distributed under the terms of the <u>Creative Commons Attribution License 4.0</u>, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Contents

Preface

I	Introduction: Physics of Life J. Elgeti, D. A. Fedosov, G. Gompper
Α	Methods
A1	Super-resolution in Optical Microscopy S. U. Egelhaaf
A2	3D structure determination using electron cryo-microscopy <i>C. Sachse</i>
A3	Fluorescence-based Experimental Techniques: FCS, FRAP, and FRET J.K.G. Dhont
A4	Scattering <i>R. Zorn</i>
A5	Macromolecular Crystallography T. E. Schrader
A6	Neutron Imaging S. Förster
A7	Protein NMR Spectroscopy in Solution <i>P. Neudecker</i>
A8	Computational Image Analysis P. Kollmannsberger
A9	Integrative Structural Biology and Hybrid Modeling <i>G. F. Schröder</i>
A10	Simulation Methods in a Nutshell <i>M. Ripoll</i>
В	Molecules
B1	Protein Folding and Protein Stability J. Fitter, A. Stadler
B2	Protein Dynamics R. Biehl
B3	Crowded Protein Solutions: Dynamics, Clustering and Phase Behavior <i>G. Nägele</i>
B4	Membrane Channels & Pumps JP. Machtens
B5	DNA and Chromatin H. Schiessel

- B6 G-protein-coupled receptors A. Baumann
- B7 Amyloid aggregation B. Strodel
- B8 Optogenetics V. Gordeliy, I. Okhrimenko, K. Kovalev
- B9 Synthetic Biology Science and Engineering of Synthetic Biological Systems *F. C. Simmel*

C Membranes, Filaments & Networks

- C1 Membranes and vesicles *T. Idema*
- C3 Cytoskeletal Filaments Actinflaments, Microtubules and Intermediate Filaments S. Köster
- C4 Theory of Semiflexible Network Materials *C. Storm*
- C5 Motor Proteins and Cytoskeletal Filaments: from Motility Assays to Active Gels *T. Auth*
- C6 Hydrodynamics of the active cytoskeleton *K. Kruse*
- C7 Synapses C. Karus, C. R. Rose

D Cells

- D1 Physics and modelling of intracellular diffusion S. Kondrat
- D2 Macromolecular Crowding: Colloidal Suspensions as Models for Biological Systems *P. R. Lang, Y. Liu*
- D3 Neuronal signaling F. Müller
- D4 Sequential Bottom Up Assembly of Synthetic Cells I. Platzman, J. P. Spatz
- D5 Theory of biological force sensing B. Sabass
- D6 Mechanobiology of Animal Cells *R. Merkel*
- D7 Cell and Tissue Mechanics *R E. Leube*
- D9 Microelectrode devices for *in vitro* and *in vivo* recordings of cellular activity A. Offenhäusser, S. Weidlich, D. Kireev, K. Srikantharajah, V. Rincón Montes

E Multicellular & Collective Behavior

- E1 Rheology *M. P. Lettinga*
- E2 Modeling blood flow and primary hemostasis in microcirculation D. A. Fedosov
- E3 Principles of Active Matter *R. G. Winkler*
- E4 Motility of Cells and Microorganisms T. Auth, J. Elgeti, R. G. Winkler, G. Gompper
- E5 Bacterial biofilms Physical determinants of microbial architecture *B. Maier*
- E6 Tissue Growth J. Elgeti
- E7 Necessity and feasibility of large-scale neuronal network simulations *M. Schmidt, M. Diesmann, S. J. van Albada*

F Diseases and Systems Biology

- F1 How Physics Shaped our Senses U. B. Kaupp, L. Alvarez
- F2 Population Genetics and Evolution *J. Krug*
- F3 Biophysics of Killing *H. Rieger*
- F4 Monogenetic diseases C. Fahlke
- F5 Physics of the malaria parasite *U. S. Schwarz*
- F6 Alzheimer's disease J. Kutzsche, D. Willbold
- F7 Quantitative approaches to antibiotic resistance *T. Bollenbach*

Preface

The spring school "Physics of Life" provides an introduction to and an overview of current research topics in biophysics of living systems, with an emphasis on understanding biological structure, dynamics and function. Biomolecules, cells, tissues, and their multiscale interactions are the main building blocks of biological organisms. The physical understanding of their structural and dynamical properties and their organization and synergy is very challenging due to the enormous complexity and non-equilibrium behaviour of these systems. However, this knowledge is essential for linking the structure and dynamics of biosystems to their corresponding functions.

The goal of this spring school is not only to give an overview of selected topics from biophysics to students and postdocs in physics, chemistry and biology, but also to establish an interdisciplinary connection between these fields. This includes, in particular, the introduction of biologists and chemists to physical experimental methods and theoretical modelling, and the introduction of physicists to the large variety of fascinating biological phenomena.

Introductory lectures present the basics of biosystems and biophysics. These lectures are intended to establish a common level of basic interdisciplinary knowledge. Subsequent lectures treat more advanced topics within different disciplines and emphasize interdisciplinary aspects.

Topics of the lectures include:

- Experimental and Theoretical Methods

The study of biological systems is particularly challenging since very often the macromolecular building blocks are inherently complex and the relevant length and time scales in these systems span many orders of magnitude. Success requires a combination of preparative techniques (synthesis), and the elucidation of structural and dynamical properties by scattering, microscopy, rheology and single molecule techniques. Lectures on experimental methods are complemented by theoretical frameworks of classical statistical mechanics, continuum hydrodynamics, and scaling theory. Furthermore, since many biological phenomena are far too complex to be well-described by an analytical theory, simulation techniques, such as Molecular Dynamics, Monte Carlo, and mesoscale hydrodynamics simulations, are often necessary and will be introduced in the basic lectures.

- Basic Building Blocks: Bio-Macromolecules

Bio-macromolecules are the basic building blocks of any biological system. Examples include various proteins, DNA, and lipids, with their properties and mutual molecular interactions inducing the assembly of biomolecules into complex structures with a versatile biological functionality. Bio-macromolecules possess exquisitely designed structures required for fulfilling their specific tasks, and a malfunction in molecular structure or dynamics often leads to the development of disease.

- Membranes, Filaments, and Networks

Membranes, filaments, and networks are biomolecular assemblies with distinct functions in the cell. Membranes serve as the main barriers between different cellular compartments and cells, while filaments assemble into cytoskeletal networks defining cellular mechanics, motility, and function. The school will cover various biophysical aspects of lipid membranes, membrane-protein interactions, biological filaments (e.g. actin, microtubules, intermediate filaments), motor proteins, active cytoskeletal networks, and synapses.

- Biological Cells

Mechanics and the behaviour of cells determine the development and functionality of various tissues. It is well known that cells may act differently in different environments, a fascinating adaptation which is facilitated by various cell organelles and machinery. Here, the lectures will cover cellular mechanics and adhesion and how these properties and processes are modulated. In addition, cell division, motility, and signalling are reviewed. Finally, an in-vitro cell model and the bridging of cells and electronics (i.e. bioelectronics) are discussed.

- Multicellular Organization and Collective Behaviour

The next level of biological organization is multicellular assemblies which constitute tissues. Currently, one of the most fascinating research directions is tissue growth and repair, because it opens a variety of avenues for biomedical applications. Other topics include mechanical properties of tissues and rheology. In addition to multicellular organization, the lectures touch upon other areas of living matter, such as the collective behaviour of swimming micro-organisms and bacterial biofilms. Finally, current developments in the theoretical description of the collective behaviour of active systems are discussed.

- Systems Biology and Diseases

Systems biology aims to interconnect various biological components, in order to integrate information about specific biocomponents into a comprehensive picture for complex biological systems. Thus, it focuses on complex interactions within biological systems. Often, the malfunctioning or alteration of such interactions leads to different diseases and disorders. Several prominent examples are monogenetic diseases, Alzheimer's and malaria. Finally, the topics of antibiotics resistance, disease spreading, and genetics and evolution are addressed. Of course, in the end all the different levels and scales must work together, from molecules to information processing. This is illustrated for the important case of neurobiology.

This school could not take place without the help and dedication of many colleagues. We would like to express our sincere thanks to all of them for the effort and enthusiasm which they have put into the preparation and presentation of their lectures and manuscripts. Without the participation of all these colleagues, the program would not be as interesting, versatile, and attractive.

We are grateful to our colleagues from Forschungszentrum Jülich for their willingness to conduct these courses. The FZJ divisions involved in the lectures are: ICS-1/JCNS-1: Neutron Scattering, ICS-2/IAS-2: Theoretical Soft-Matter and Biophysics, ICS-3: Soft Matter, ICS-4: Cellular Biophysics, ICS-5: Molecular Biophysics, ICS-6: Structural Biochemistry, ICS-7: Biomechanics, ICS-8: Bioelectronics, INM-6: Computational and Systems Neuroscience, and JCNS-MLZ at Heinz Maier-Leibnitz Zentrum in Garching.

In particular, these colleagues are:

Dr. Thorsten Auth, Prof. Arnd Baumann, Dr. Ralf Biehl, Prof. Markus Diesmann, Prof. Jörg Fitter, Prof. Valentin Gordeliy, Dr. Peter Lang, Dr. Jan P. Machtens, Prof. Rudolf Merkel, Prof. Frank Müller, Prof. Gerhard Nägele, Dr. Philipp Neudecker, Dr. Marisol Ripoll, Dr. Benedikt Sabass, Dr. Tobias Schrader, Prof. Gunnar Schröder, Dr. Andreas Stadler, Prof. Birgit Strodel, Prof. Dieter Willbold, Prof. Roland G. Winkler, and Dr. Reiner Zorn.

We are very glad that several colleagues from universities and research laboratories have agreed to contribute to the program of the school:

Prof. Patricia Bassereau	Laboratoire Physico Chimie Curie, Institut Curie, PSL Research University, Paris, France
Prof. Tobias Bollenbach	Institute of Biological Physics, University of Cologne
Prof. Stefan U. Egelhaaf	Condensed Matter Physics Laboratory, Heinrich Heine University Düsseldorf
Prof. Jochen Guck	Biotechnology Center, Technische Universität Dresden
Prof. Timon Idema	Department of Bionanoscience, Kavli Institute of Nanoscience, Delft University of Technology, The Netherlands
Prof. U. Benjamin Kaupp	Molecular Sensory Systems, Center of Advanced European Studies and Research, Bonn
Prof. Philip Kollmannsberger	Center for Computational and Theoretical Biology, University of Würzburg
Prof. Svyatoslav Kondrat	Department of Complex Systems, Institute of Physical Chemistry, Warsaw, Poland
Prof. Sarah Köster	Institut für Röntgenphysik, Georg-August-Universität Göttingen
Prof. Joachim Krug	Institute for Theoretical Physics, University of Cologne
Prof. Karsten Kruse	Departments of Biochemistry and Theoretical Physics, NCCR Chemical Biology, University of Geneva, Switzerland
Prof. Rudolf E. Leube	Institute of Molecular and Cellular Anatomy, Uniklinik RWTH Aachen
Prof. Berenike Maier	Experimental Biophysics, Institute of Theoretical Physics, University of Cologne
Prof. Heiko Rieger	Center for Biophysics and Department of Physics, Saarland University, Saarbrücken
Prof. Christine R. Rose	Institute of Neurobiology, Heinrich Heine University Düsseldorf
Prof. Carsten Sachse	Structural and Computational Biology Unit, European Molecular Biology Laboratory, Heidelberg
Prof. Helmut Schiessel	Institute Lorentz for Theoretical Physics, Leiden University, The Netherlands
Prof. Ulrich S. Schwarz	Institute for Theoretical Physics and BioQuant-Center for Quantitative Biology, Heidelberg University
Prof. Friedrich C. Simmel	Physics of Synthetic Biological Systems, Physics Department, Technische Universität München
Prof. Joachim P. Spatz	Department of Cellular Biophysics, Max Planck Institute for Medical Research, Heidelberg
Prof. Cornelis Storm	Theory of Polymers and Soft Matter group, Department of Applied Physics, Eindhoven University of Technology, The Netherlands

We are very grateful to the board of directors of Forschungszentrum Jülich for the continuous organizational and financial support, which we have received for the realization of the IFF spring school and the production of this book of lecture notes. We acknowledge the financial support by the European Network of Excellence SoftComp. Finally, our special thanks go to Barbara Daegener for the general management, Meike Kleinen for her help in the preparation of the lecture notes, and Claire Ryalls, Ursula Funk-Kath, and Dorothea Henkel for the webpages of the school.

Gerhard Gompper, Jan K. G. Dhont, Jens Elgeti, Christoph Fahlke, Dmitry A. Fedosov, Stephan Förster, Pavlik Lettinga, Andreas Offenhäusser

January 2018

I Introduction: Physics of Life

Jens Elgeti, Dmitry A. Fedosov, Gerhard Gompper Theoretical Soft Matter and Biophysics Institute of Complex Systems and Institute for Advanced Simulation Forschungszentrum Jülich GmbH

Contents

1	Introduction	2
2	Macromolecules and their Interactions 2.1 DNA 2.2 Proteins 2.3 Biopolymer Filaments and Networks 2.4 Membranes	4 5 5 7 8
3	Biological Materials and their Properties	9
4	Rheology of Biofluids	11
5	Out-of-Equilibrium Physics and Living Matter5.1Active Matter, Self Propulsion, and Microswimmers5.2Tissue Growth	
6	Grand Challenges, and Major Goals	15

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Life is described by biology, isn't it? All the various plants and animals, all the different organs and cell types, and all the myriads of proteins and signaling molecules and chemical reactions. What is "physics of life" supposed to be? There is probably no unique or generally accepted answer to this question. However, there are certainly several aspects of life and living organisms, where physics plays a crucial role.

What probably comes to mind first are physical techniques used to study and understand living systems. X-ray imaging and X-ray scattering, nuclear magnetic resonance (NMR) and magnetic resonance imaging, electron microscopy, and positron emission tomography (PET) are all methods developed in physics, and later applied to living systems. This is an important aspect of the interaction of physics with biology and medicine, but not the essence of "physics of life". Instead, we see the following key issues:

• Materials and Material Properties

Organisms consist of macromolecules, like proteins and DNA, macromolecular aggregates, cells, tissues, organs, etc. These objects have material properties such as elasticity, deformability, surface and adhesive properties, as all other materials. For example, DNA has a bending and torsion modulus, and is electrically charged, properties which are important for understanding how very long DNA strands can be compacted into the tiny cell nucleus. Membranes surrounding cell organelles and cells have a bending rigidity and a shear modulus, but also a permeability for small molecules and water. The knowledge of material properties forms the basis for any quantitative understanding of living matter, in the same way as knowledge of the properties of various metals and plastics is required to construct a robot.

• Macromolecules and their Interactions

The key players for the functioning of a cell are macromolecules. For a description of processes in the cell, not only the proteins involved have to be known by name. In addition, their structure, the charge distribution on their surfaces, the internal dynamics, and their interactions have to be characterized to allow for quantitative predictions based on statistical theories.

• Rheology of Biofluids

Most organisms contain a large amount of fluids. These fluids are not static, but are moving and flowing. Examples are blood, the cerebrospinal fluid, and the cytoplasmic fluid. The behavior of these fluids is governed by the same hydrodynamic laws – the Navier-Stokes equation – as the flow of water in the oceans. Thus, the physics of fluids is crucial to understand the rheology of bioflows.

• Non-Equilibrium Systems and Active Matter

Living systems are persistently out of equilibrium – an organism at equilibrium is dead. This persistent out-of-equilibrium state is maintained by continuous energy consumption. The physical understanding of out-of-equilibrium systems is much more difficult than well-established thermodynamics – and holds many surprising phenomena to be discovered. This will be among the most exiting areas in physics in the coming decades.

• Structure versus Function

Condensed matter physics is mostly concerned with "structure" and "dynamics". What is

the structure of Ga interstitials in Si-matrices of semi-conductors? What is the structure of the fracture surface of a rod which breaks due to an overload? What is the conformation of a polymer in a compound material. In living systems, this is of course also important, but equally or even more important is "function". How does a signaling molecule cause a membrane channel to open, which then allows the influx of other molecules causing the cell to perform a task. For physicists, this requires a new way of thinking about dynamical systems, and opens the field to tackle completely new problems.

• Complex Systems

Biological and living systems are typically very complex, with a large number of different components, highly dynamical conditions, spatially very inhomogeneous distributions, etc. This is a grand challenge for physics, because the main physics approach is to boil things down to the essential ingredients, and then to understand the underlying mechanisms in detail. It will certainly not be possible to understand all living matter in this way. However, it will be a very interesting journey to find the questions, problems and systems, where a clever reductionist approach can give fundamentally new insights into the processes of life.

• Emergent Behavior

Ensembles of many interacting objects often show collective behavior, which does not depend on the detailed structural or dynamical properties of the individuals. An example is the formation of swirls in large schools of fish, which are seen for many different species. Another example is the turbulent behavior of many active bacteria on a substrate. These ensembles exhibit "emergent properties" that the smaller/simpler entities do not have.

A good example for emergent behavior is the brain. The individual neurons are rather simple cells, similar to a transistor in electronics. But if 100 billion or so are coupled together in the human brain, they start to become self-conscious and think about stupid jokes. This emergent behavior is not a property of the neuron itself. An even better example is life itself. An organism consists mainly of hydrogen, carbon, and oxygen. Mix 100 kilograms of it together, shake well, so that everything falls in place, and it starts to run around and uses a cell phone.

• Synthetic Biology

Richard Feynman (Nobel laureate in Physics 1965 for quantum electrodynamics) once said "what I cannot create, I do not understand". This statement can be interpreted in the context of biology and life sciences in such a way that we should aim at constructing artificial model systems, which recreate essential aspects of living systems. There is indeed a recent development in science, which is heading in this direction.

Synthetic biology is nowadays defined as the artificial design and engineering of biological systems and living organisms for purposes of improving applications for industry or biological research. In general, its purpose can be described as the design and construction of novel artificial biological pathways, organisms or devices, or the redesign of existing natural biological systems.

• Development and Tissue Growth

From the developing embryo to the caterpillar transforming into a butterfly, from an expanding bacterial colony to a growing tree, cells need to grow, divide, and rearrange to build the final, fully developed state. Besides biochemical and other regulation mechanisms, mechanical forces and constraints are key players in morphogenesis and tissue patterning. Understanding these mechanical forces and their feedback to growth and patterning is thus mandatory for understanding development and growth. At the same time, these feedbacks and growth result in novel physical phenomena and unprecedented behavior.

• Evolutionary Biology and Genomics

Chance and necessity in evolution is a fundamental theme of biology. How can this dynamics be understood, starting from its molecular basis, which lies in genes and their interactions? How do adaptation and functional innovation take place in the sea of stochastic changes of molecular evolution? How can genomic data and evolution experiments, for example in bacterial systems, be employed to develop and test statistical theories of evolution?

• Environment and Stimulation

Biological systems are not living in vacuum, but are constantly exposed to a large variety of external signals and stimuli. These include mechanical, optical, and electrical signals. Therefore, it is important to understand the response of living matter to such stimuli on all levels, from single molecules and cells to tissues. This has spawned the fields of biomechanics and bioelectronics. In the latter case, electronic devices are developed to stimulate nerve cells (neural prosthetics) or to detect reactions of cells to other environmental stimuli (biosensors). A particularly interesting recent development is optogenetics, where light signals are used to control genetically modified cells (typically neurons) in living tissues.

• Multiscale Modeling

A detailed understanding of living systems requires a quantitative description. For complex systems in living matter, phenomena on very different length- and time-scales are often intimately linked together. As an example, the simulation of a cell in a microfluidic flow simultaneously requires the resolution of cell-wall interactions (nano-scale), cell deformation (micro-scale) and fluid flow (micro- or millimeter scale). From the theoretical side, this requires multiscale modeling approaches, which cover a wide range of scales. A high-fidelity modeling will need computing resources of the next generation of supercomputers.

Some general discussion on various aspects of "Physics of Life" can be found in Refs. [1–11].

2 Macromolecules and their Interactions

The chemistry and physics of synthetic macromolecules has been studied intensively in the last decades, starting essentially with work of H. Staudinger (Nobel Laureate in Chemistry 1953) and his colleagues in the 1920s. This has lead to a wealth of knowledge about the behavior of linear and branched macromolecules (= polymers) in dilute and concentrated solutions, culminating in the Nobel Prize in Physics for P.-G. de Gennes in 1991. In particular, it has been shown that many polymer properties depend not so much on the properties of a single monomer, but are dominated by the large number N of (nearly) identical units in such a chain molecule.

This implies that many properties are universal. For example, the end-to-end distance R_e of a polymer in good solvent scales as $R_e \sim N^{\nu}$, with $\nu = 0.58$ in three spatial dimensions.

Macromolecules are also the main building blocks of all biological cells. DNA is a linear polymer, which in its sequence of nucleic acids stores the genetic information. A complex machinery in the cell transforms this information into amino-acid sequences, which constitute functional units – the proteins. There is a large variety of proteins in the cell, each of which performs a specific task. Many proteins have globular shapes, but some proteins assemble into long filaments, like actin filaments and microtubules. Another type of self-assembly occurs for lipid molecules, which due to their amphiphilic character form bilayer membranes. These membranes form the outer envelope of the cell – the plasma membrane – and are present in the interior of the cell to form compartments. A special class of proteins is membrane proteins, which play an important role in the regulation of the transport of nutrients, signaling molecules, and waste products through the membrane.

2.1 DNA

The best-known polymer in living systems is probably the carrier of the genetic information, *deoxyribonucleic acid (DNA)*, see Fig. 1. It consists of a sequence of four monomers, the nucleotides adenine, thymine, guanine and cytosine, which in pairs form the famous double helix. Three subsequent bases encode for one amino acid. The genetic code specifies 20 standard amino acids.

In eukaryotic cells, DNA is tightly packed within the nucleus in the form of chromosomes [12, 13]. The (linear) contour length of DNA within a chromosome is typically of the order of centimeters, and is contained within the volume of the nucleus which has a diameter of a few microns. If this amount of DNA were simply compressed in such a small space, the resulting pressure would rip apart any conceivable cage. Such a tight packing is enabled through the interactions of DNA with so-called chromatin-proteins, which are cylindrical segments around which the DNA is wound. It is still a matter of debate how the mechanism of tight packing functions. It cannot be too tight, because the stored information must still be accessible for protein synthesis. In prokaryotic cells (bacteria), DNA is present as ring-like molecules, which are super-coiled through interactions with enzymes. These helical super-coiled DNAs float freely in the cell's plasma. It is thus clear that DNA only functions with the help of complex interactions with several types of proteins. This is the case for almost all processes in living matter: many different types of molecules work together to maintain life.

2.2 Proteins

Proteins are the machinery and the building blocks of life. Most of biological functions, from molecular motors to the cytoskeleton, from DNA replication to anti-virus protection, are performed – with a very few exceptions – by proteins. Proteins (also known as polypeptides) are linear chains of amino-acids, which fold into a globular or fibrous form. The amino acids are co-valently connected by peptide bonds. Proteins are formed from a "library" of only 20 standard amino-acids. These amino-acids can be broadly differentiated into 8 hydrophobic amino-acids – normally buried inside the protein core –, 4 charged amino-acids – with their side chains often making salt bridges –, and 8 polar amino-acids – usually participating in hydrogen bonds as proton donors or acceptors.



Fig. 1: Structure of biopolymers. Left: DNA is a double-stranded helix. Right: Microtubules are assembled from dimers of α and β tubulin; 13 protofilaments from a hollow tube.

Proteins were first described already in 1838 by the Dutch chemist Gerhardus Johannes Mulder and named by the Swedish chemist Jöns Jakob Berzelius.

The extraordinary versatility of proteins is achieved by folding the linear amino-acid chain into a three-dimensional structure. The particular sequential order of amino-acids defines the "primary structure" of a protein. One of the most distinguishing features of polypeptides is their ability to fold into a globular state. This folding occurs on two hierarchical levels. The "secondary structure" consists of α -helices, pieces of the chain which form tight staircase-like cylindrical units, and β -sheets, which form ribbon-like units. One of the important forces that stabilizes the secondary structure are hydrogen bonds between the amino-acids. These elements then order into a three-dimensional arrangement, which defines the "tertiary structure". The function of a protein strongly relies on the correct tertiary structure. An example of the protein aquaporin 1 with a high content of α -helices is shown in Fig. 2.

The extent to which proteins fold into a defined structure varies widely. Some proteins fold into a highly rigid structure with small fluctuations and are therefore considered to be single structure. Other proteins – called "intrinsically disordered proteins" – undergo large rearrangements from one conformation to another. Thus, not all proteins require a folding process in order to function.

While cells generate proteins with high fidelity, a consistent prediction of protein structure on the basis of their sequence alone has not been possible so far. Conversion of the DNA sequence into an amino-acid sequence is well understood, and a programmable synthesis is possible; however, the question "how does the three-dimensional structure emerge?" remains unanswered.



Fig. 2: Aquaporin 1 is a member of the Aquaporin family. It was first discovered in red blood cells. In 2003, Peter Agre was awarded the Nobel prize for discovering this "water channel". This channel is used for rapid transportation of water across the cell membrane. It is very specific and only allows passage of water molecules, but is impermeable to hydrogen ions. Note the large content of α -helices. From Ref. [14].

2.3 Biopolymer Filaments and Networks

An important biopolymer is *actin*, one of the three major components of the cytoskeleton. It participates in many important cellular functions, including muscle contraction, cell motility, cell division and cytokinesis, vesicle and organelle movement, cell signaling, and the establishment and maintenance of cell junctions and cell shape. Actin filaments (F-actin) consist of globular monomeric protein subunits (G-actin), which form a double-helical structure.

Another structural component of the cytoskeleton is *microtubules*. For instance, they contribute to a structural integrity of the cell, serve as the routes for vesicle transport, and facilitate the separation of chromosomes during mitosis. The structure of microtubules consists of a hollow tube of dimers of the protein subunits α and β tubulin, which are arranged in protofilaments, see Fig. 1.

An important aspect of biopolymers is their mechanical properties, which can be characterized by the persistence length, *i.e.* the length below which a polymer behaves essentially like a stiff rod, whereas it behaves like a random coil on much larger scales. The persistence lengths of DNA, actin and microtubules are $50 \ nm$, $10 - 20 \ \mu m$ and about $1 \ mm$, respectively. These numbers can be connected to their functions and location in the cell. Very long DNA strands have to be curled up in the cell nucleus, which is only a few μm in diameter (the total length of a single DNA strand in the human genome is $1.8 \ m$); DNA therefore has to very flexible. On the other hand, microtubules have to be very stiff, because they serve as the "highways" for vesicle transport; their persistence length therefore exceeds the cell diameter.

In contrast to synthetic polymers, biopolymers often do not have permanently polymerized and fixed structures. Instead, they are dynamic states of an active polymerization process, which proceeds with different rates at the two ends of the *polar* filaments. Actin belongs to this class

of "tread-milling" filaments. Microtubules have another type of active dynamics, which consists of periods (in time) of steady growth, which are interrupted by sudden depolymerization events ("catastrophes").

2.4 Membranes

The main component of a cell membrane is lipids, amphiphilic molecules with two hydrocarbon tails (see Fig. 3), which form bilayers in water. This was first shown experimentally by Gorter and Grendel (1925). They first extracted lipid molecules from the plasma membrane of red blood cells. By pouring the lipids onto a water surface, and comparing the area of the resulting patch with the surface area of the original cell surfaces, they concluded that the cell membrane consists of a lipid bilayer. Later, it turned out that their arguments contained several errors, they had for example underestimated the size of the original cell surface. However, these errors canceled out each other to a large extent, so that their conclusion was nevertheless correct. Nowadays, it is well established that the major component of the cell membrane is a lipid bilayer. Within cell membranes, there are usually a variety of other molecules embedded (an artistic impression is given in Fig. 4). The macromolecules and macromolecular complexes, that are embedded in the membrane, regulate exchange of matter with its surroundings. For example, ion channels are complexes which regulate exchange of specific ions.



Fig. 3: *Phosphatidyl-choline (PC) is a typical example of a membrane lipid. The polar head (left) is connected to two hydrocarbon tails (right).*

Membranes consisting of lipids and (bio-)polymers are quite easily deformable. A typical example of such a biological membrane is the "skin" of a red blood cell (see Sec. 4 below), the size of which is in the few micrometer range. Due to the flexibility of the membrane, a red blood cell can pass through micro-vessels with a diameter four times smaller than its own size. In addition, they are easily deformed by solvent flow, which is important for blood flow in the microvasculature.

Biological membranes have a very complex structure, as they are composed of many different types of lipids, cholesterol, and proteins, which can assemble into rafts or other supramolecular structures. Thus, fundamental insights on membrane properties are often obtained through the study of more basic "physical" membranes, which have a much simpler composition in comparison to biological membranes.



Fig. 4: Artistic impression of a cell membrane in which many macromolecules and macromolecular complexes are embedded. From Ref. [15].

3 Biological Materials and their Properties

Biological materials are extremely complex composites, characterized by intricate structural and mechanical properties [16–18]. The main building blocks are macromolecules (e.g., DNA, proteins, lipids – see Sec. 2), which then are assembled into a fascinating variety of macromolecular aggregates and structures (e.g., filaments, membranes, cells). In turn, macromolecules and cells continue to organize into further complex assemblies, yielding biological tissues, organs, and finally complete organisms. Therefore, from the structural point of view, a number of hierarchical classes of composite biological materials (i.e., macromolecules, cells, tissues, etc.) are generally defined and investigated.

The elucidation of structural and mechanical properties of biological materials is a significant scientific endeavor. A starting point here is the knowledge of involved components, such as the constituent macromolecules, molecular components for supramolecular assemblies and so on. However, even a complete knowledge about involved constituents does not generally permit direct predictions of their hierarchical organization. One of the prominent examples of supramolecular structures is cell membrane [19], which is schematically shown in Fig. 4. Biomembranes are typically composed of many different lipids and proteins, which self-organize through a number of physical interactions, including hydrophilic/hydrophobic and ionic interactions, Van der Waals and electrostatic forces. Understanding of these interactions is essential for the prediction of structural characteristics of various supramolecular assemblies. Clearly, structural properties of biological materials are associated with their mechanical properties and function [16–18]. For instance, a lipid bilayer assembly of a cell membrane results in a bending elasticity of the overall membrane, while a non-uniform distribution of intra-membrane components may lead to the presence of spontaneous curvature (i.e., a preference to curve to the inside or the outside of the cell, as the two leaflets are usually not identical). This gives a cell the means to control membrane properties by adjusting the membrane composition. In addition, trans-membrane proteins within a biological membrane allow the control of its permeability characteristics for the exchange of water, ions, and small molecules between the cell plasma and the extracellular space.



Fig. 5: Illustration of the major structures inside an eukaryotic (animal) cell. From the Indian Museum, Kolkata, India.

Cell mechanics is a large sub-field of biophysics, which concerns the material properties and behavior of living cells [20,21]. Here, the main aim is to understand and predict cell's mechanical properties and behavior starting from cellular structures (see Fig. 5), and relate them to cellular functions. Mechanical properties, such as elasticity, adhesiveness, and viscosity, can be probed by a number of experimental techniques, including micropipette aspiration [22], optical tweezers [23], flow cytometry [24], etc. It is important to keep in mind that these overall mechanical properties arise from a number of contributions from different cellular components, including cell's membrane, cortex, bulk cytoskeleton, organelles, etc. In addition, biological cells are generally able to adapt or respond to various mechanical deformations, making their properties dependent on applied stress, strain, and strain rate. Such mechano-sensitive adaptations proceed through cytoskeletal re-organization and force generation, and therefore biological cells are often referred to as active materials. Furthermore, biological cells might be motile, can initiate cell division or programmed cell death called apoptosis.

Many biological materials have mechanical properties that are far superior to those in existing synthetic materials [18]. For example, spider silk has impressive mechanical properties, because its strength is comparable to steel, while its density is lower than that of cotton. This motivates large scientific efforts for the development of bioinspired (or biomimetic) materials in order to take advantage of the versatile properties of biological materials. Another major goal here is to establish synthetic "biocompatible" materials which can be introduced into the human body, leading eventually to the substitution of tissues and organs.

In conclusion, most biological materials possess outstanding and adaptive mechanical properties. These properties emerge from the complex hierarchical structures and are tightly coupled to biological functions of such materials. Elucidation of the material properties of biological systems requires understanding of structural organization and mechanical characteristics at



Fig. 6: Simulation snapshot of blood in shear flow. RBCs are shown in red and in orange, where orange color depicts the rouleaux structures formed due to aggregation interactions between RBCs. The image also displays several cut RBCs with the inside drawn in cyan to illustrate RBC shape and deformability. Reused with permission from Ref. [25].

multiple length scales starting from macromolecules, to supramolecular assemblies, cells, and finally organs, and whole organisms. Here, physics provides its indispensable contribution by the application of a variety of physical methods and the establishment of physical models and mechanisms in order to understand living systems.

4 Rheology of Biofluids

Fluids constitute a major part of all living organisms. Examples are blood, lymphatic liquid, cytoplasmic and cerebrospinal fluids. Biofluids are generally not simple viscous fluids like water, but are suspensions of cells and macromolecules, and are therefore often called "complex fluids". Complex fluids [26] may exhibit an intricate mechanical response to applied stresses or strains, such as viscoelasticity, shear-thinning or shear-thickening viscosity, yield stress, etc. Rheology concerns the flow of matter (primarily liquids) and aims to characterize the response of complex fluids to applied stresses. A particularly interesting question in rheology is how to connect macroscopic properties of complex fluids (e.g., shear-thinning) with the behavior and interactions of suspended particles or molecules in flow.

Rheology of blood has been subject to extensive investigations [27,28], since the flow properties of blood are clearly very important for organism functioning. Blood transports oxygen and nutrients to cells of the body, removes waste products, and circulates many important molecules and cells. Furthermore, changes in its flow properties are often correlated with various blood disorders and diseases, such as anemia, hypertension, malaria, etc. Blood is a suspension of blood cells, various proteins, and dissolved ions. The major cellular component (about 45% by volume as illustrated in Fig. 6) corresponds to red blood cells (RBCs), and therefore mechanical

and flow properties of RBCs determine rheological characteristics of blood.

RBCs in whole blood form aggregates called "rouleaux", which resemble stacks of coins (see Fig. 6) [27, 29]. The RBC aggregation is facilitated by the plasma proteins and results in a significant increase of blood viscosity at low shear stresses [27, 29]. Moreover, whole blood exhibits a yield stress (a threshold stress for flow to begin) due to the RBC aggregation [27]. This means that under very low stresses blood behaves similar to a solid. Overall rheological properties of blood are characterized by strong shear-thinning and weak elasticity. The shear-thinning characteristics of blood arise from the RBC ability to aggregate, deformability, and dynamics in flow [25, 30].

Understanding of fundamental rheological properties of biofluids is important for the prediction of their flow characteristics within an organism. This knowledge is also crucial for many biomedical and bioengineering applications, such as the development of blood substitutes, clinical tests, and drug delivery carriers.

5 Out-of-Equilibrium Physics and Living Matter

Classical non-equilibrium systems are driven out of equilibrium by external forces. Water is pumped through a pipe, heat flows from hot to cold, and currents follow a potential difference. The defining property of living matter is that the force driving it out-of-equilibrium originates from the microscopic constituents themselves. Bacteria swim in a fluid, motor proteins move the filaments in the actin cytoskeleton, and cells grow and divide. It is always that some microscopic components (e.g., motor proteins, flagella) are responsible for the activity. Typically, living matter is divided in two distinct classes. The first class is *microswimmers* or *active matter*, whose constituents generate forces or stresses, leading to active contributions in the force balance. The second class is called *growing matter* or *growing tissues*, and can be characterized by mass-balance violated through an active source or sink.

Examples range across all length scales: from swimming *Escherichia coli* that rotate helical flagella to swim [31], over active polymers in the cell that constantly polymerize and depolymerize and are connected by molecular motors to create stresses [32], to large schools of fish or groups of sperm that move collectively [33], as well as growing and developing biological tissues [34].

5.1 Active Matter, Self Propulsion, and Microswimmers

Motion, in one way or another, is a hallmark of life. Within the cell, motor proteins generate forces that move filaments and vesicles around. The whole cytoskeletal network is very dynamic and active. On a larger scale, cells move through soil, fluid, or tissues.

Motion is generated by motor proteins. Examples are myosin in the actin cytoskeleton that lets cells crawl, axonemal dyneins that allow the flagellum of sperm cells to beat, or myosin in sarcomeres to contract our muscles. The basis of activity and motion is always molecular motors converting chemical energy into forces. Thus, the system is driven out of equilibrium by locally generated forces. Sometimes, the activity can be described by an elevated effective temperature [35], which captures the increased kinetic energy of the particles due to their continuous propulsion. However, often this description is inappropriate, and the activity manifests itself in phenomena not considered by equilibrium physics. Nevertheless, equilibrium concepts and analogies can sometimes be useful to depict what is going on. For example, collections of self-



Fig. 7: Schematic of active Brownian spheres moving near a surface. From Ref. [42].

propelled disks with repulsive interactions spontaneously "phase separate" [36], motile tissues undergo "glass-like arrest" as density or adhesion increases [34,37], or non-equilibrium surface accumulation occurs. The physics of active matter and microswimmers has gained increasing attention in recent years; comprehensive reviews can be found in Refs. [33,38–41].

A particular illustrative example microswimmers in confinement is a sperm swimming between two planar surfaces. For many different cells and microorganisms, microswimmers are found to accumulate at surfaces. From a physics perspective, the simplest form of a microswimmer – the proverbial "spherical cow" – is a simple self-propelled Brownian sphere. It has a preferred direction of motion in the body-centered reference frame, along which it moves with a nearly constant velocity. However, this preferred direction performs rotational Brownian motion, like a passive particle. Passive Brownian particles would have a uniform distribution between the walls; however, if they run around in a bounded space, they are bound to hit a wall sooner or later (see Fig. 7).

The decoupling of the rotational degrees of freedom from translational motion allows for a theoretical treatment via a Fokker-Planck equation. The solution of this equation quantifies the following physical mechanism of wall accumulation. Particles are driven to one of the chamber walls depending on their initial orientation. Less particles remain in the center. Particles at the walls get stuck as long as their direction of motion still points toward the wall. Only when this direction has changed enough – a slow process driven by rotational diffusion – they are able to move again away from the wall.

This effect – that moving particles are bound to hit and accumulate at confining surfaces – is very generic [43], ranging from self-propelled rods and spheres to idealized run-and-tumble particles, and from swimming sperm cells and bacteria (like *E. coli*) to micro-algae (like Chlamy-domonas) [44]. Thus, the minimalistic approach of physics provides a useful basis for many biological systems.

5.2 Tissue Growth

"The ability of cells and tissues to respond to mechanical force is central to many aspects of biology" [46]. The prime example of growing matter is biological tissue. Cells grow and divide, they die and disappear. Other forms of growing materials are bacterial colonies or possibly polymers during polymerization. Biological tissues form functional parts of organisms, com-



Fig. 8: Illustration of homeostatic pressure. **A** The homeostatic pressure is the steady state pressure (cell growth and division balanced by cell death) of a tissue. **B** Two tissues can compete purely by mechanical forces if separated by a movable wall. From Ref. [45].

posed of cells. They develop during embryogenesis, and most are under constant renewal over the course of their lifetime. In the past decades, it has become increasingly clear that physics, and especially mechanics, play an important role in cellular and tissue growth. Nowadays, it is well accepted that "mechanical feedback regulates proliferation" [47]. With the development of new experimental techniques and *in vitro* models, a deeper understanding of how cell growth couples to forces emerges.

A simple, heuristic argument for the physics of growth is a thermodynamic analogy: Growth is a change of volume. In terms of thermodynamics, pressure is the conjugated variable to volume. Therefore, it seems natural to assume a feedback between pressure and cellular growth.

Mechanical feedback on growth has been implemented in many different ways [48, 49]. One intuitive approach is to expand the growth rate in powers of the pressure around the zero growth-rate pressure – the *homeostatic pressure* [45], illustrated by a simple *gedankenexperiment* (see Fig. 8). A tissue is grown in a finite compartment, bound in one direction by a movable piston connected to a spring. As the tissue grows, it compresses the spring. Eventually, the growth forces of the tissue are not strong enough to further compress the spring. This force (divided by the area of the piston) is the homeostatic pressure. This is an intrinsic "material" property of the tissue type. In the second gedankenexperiment, the spring is replaced by a second tissue with a higher homeostatic pressure. Now the pressure equilibrates to an intermediate pressure between the two homeostatic pressures. At this pressure, the first tissue shrinks, while the second one grows. Thus, the first tissue steadily shrinks, until it finally vanishes, due to the higher pressure exerted by the second tissue.

A second illustrative example relates to the material properties of a growing or self renewing tissue. As cells in a tissue divide and die, each such event locally relaxes some stress. For polymer melts or other complex fluids, stress relaxation leads to viscous behavior on long time scales. Consequentially, tissues should be regarded as fluids on time scales much longer than their self-renewal time [50].

6 Grand Challenges, and Major Goals

The importance of "Physics of Life" has been recognized by the *Federation of American Scientists* already in 2003, when the goals of "Understanding of Complex Systems" and "Applying Physics to Biology and Medicine" were identified as two out of seven Grand Challenges in Physics for the 21st century [51]. A similar conclusion was reached by a committee of the *National Research Council (USA)* in a 2007 study: "What are the prospects for Condensed-Matter and Materials Physics in the early part of the 21st century?", where "What is the physics of life?" and "What happens far from equilibrium and why?" were identified as two out of six Grand Challenges [52].

These grand challenges are challenges for the fields of biology and physics alike. Physics has thrived on a reductionist approach: from universality of critical phenomena in statistical physics, to entropy and free-energy minimization in thermodynamics. Unraveling the underlying concepts that describe a multitude of phenomena has been key to progress in the physical sciences. On the other hand, in biological systems details often matter, and even though life follows the same physical laws, the complexity often renders answers from basic principles alone unfeasible. Thus, new principles for complex systems – "thermodynamics of live matter" – are needed to truly understand physics of life.

In the pursuit of a full understanding "Physics of Life", many major goals and objectives are on the way, or close-by. This concerns, in particular:

- A cell is already a very complex system. This renders it difficult to understand its essential mechanisms and processes. A possible approach to resolve this problem is the construction of a "Minimal Cell", both from the experimental and theoretical sides. This can be done in a step-by-step process: start with a closed fluid membrane, integrate a passive actin network, add motor proteins to make the cytoskeleton active, integrate membrane channels and pumps for material exchange, add metabolic processes for lipid synthesis, etc. Alternatively, an existing, functioning cell can be simplified by successively removing more and more of its components. Both approaches are very interesting to pursue.
- Information processing is an essential aspect of biological function on the cellular, tissue, and organismal levels. How are different signaling events segregated spatially in a cell without any membranous separation? The cytosol is an amazingly crowded space with many trafficking pathways which are independently regulated. How is this possible? For example, in neurons many independent synaptic inputs converge together to generate a post-synaptic response. The relative importance of these signals and their collaboration or lack thereof and the role of spatial organization, remains to be elucidated [11].
- Biological structures operate at multiple levels, from nano-scale molecules to meter-scale systems. Understanding the individual scales is a requisite for deeper insights into the whole system. Most importantly, these insights have to be combined and integrated into a multi-scale framework.
- Network dynamics are omnipresent in biological systems. This ranges from complex metabolic networks, signaling cascades and gene regulatory networks on the microscale to food chains and complex environmental networks on the macro-scale. Can these networks be characterized to an extent that quantitative predictions or optimizations are possible? Is spatial organization of metabolism relevant or essential? How does the network respond to external signals or perturbations?

- Activity is a hallmark of life. Unraveling the "thermodynamics of active matter" is clearly a challenge for the physical sciences. There has been significant progress recently to generalize thermodynamic concepts to non-equilibrium systems in the theoretical framework of "stochastic thermodynamics" and "fluctuation theorems", and their applications to molecular machines [53]. Can these concepts be generalized for other kinds of non-equilibrium systems?
- The continuity equation (mass conservation) has been at the basis of many physical descriptions of matter. How do we deal with changing numbers as cells grow, divide, and die, as proteins are synthesized, and filaments are polymerized dynamically? Are fundamentally new theories with additional growth terms required, as suggested by the description of tissue growth? Is it alternatively better to handle these systems as multi-component systems, where one component is converted into another?
- Biological systems are not always in a "healthy" state, but are often hampered by diseases. Diseases have two aspects in research: on the one hand, it is important to understand their origins and contribute to the design of new treatment options; on the other hand, they offer new avenues to better understand the mechanisms and processes of the organism itself. In many cases, diseases have important physical aspects; e.g., a blood clot physically blocks the flow of blood in a blood vessel, a growing tumor presses mechanically on neighboring organs. The "Physics of Disease" deals with many important diseases, such as cardiovascular diseases, cancer, neurodegenerative diseases, bacterial and viral infections, etc.
- Nature has generated all kinds of micro-machines. Can similar or novel approaches be used to construct artificial microrobots, which autonomously fulfill some tasks in medicine or environment? Can many of such microrobots work together in swarms to perform even more complex tasks?
- The coupling of the nervous system to electronics opens up the possibility of neuroprosthetics, i.e. prosthetic devices which can be directly controlled, or whose input can directly interact with the brain ("artificial eye that can see").
- Neuromorphic engineering, also known as neuromorphic computing, is the idea to employ very-large-scale integration (VLSI) systems containing electronic analog circuits to mimic neuro-biological architectures present in the nervous system [54,55]. Such designs may help to reduce the large energy consumption of present-day chip designs, and furthermore may lead to unprecedented algorithmic options and performance. Can the new insights into brain organization and function contribute to the design of new neuromorphic computer chips?

References

- [1] J. Knight, Nature **419**, 244 (2002).
- [2] K. E. Kasza, A. C. Rowat, J. Liu, T. E. Angelini, C. P. Brangwynne, G. H. Koenderink, and D. A. Weitz, Curr. Opin. Cell Biol. **19**, 101 (2007).
- [3] P. Friedl and D. Gilmour, Nat. Rev. Mol. Cell Biol. 10, 445 (2009).
- [4] J.-L. Maître and C.-P. Heisenberg, Curr. Opin. Cell. Biol. 23, 508 (2011).
- [5] N. Goldenfeld and C. Woese, Annu. Rev. Condens. Matter Phys. 2, 375 (2011).
- [6] B. Ladoux and A. Nicolas, Rep. Prog. Phys. 75, 116601 (2012).
- [7] M. Gruebele and D. Thirumalai, J. Chem. Phys. 139, 121701 (2013).
- [8] H. G. Garcia, R. C. Brewster, and R. Phillips, Integr. Biol. 8, 431 (2016).
- [9] D. Needleman and Z. Dogic, Nat. Rev. Mater. 2, 17048 (2017).
- [10] R. Austin, Nature 550, 431 (2017).
- [11] T. C. Südhof, Neuron 96, 536 (2017).
- [12] H. Schiessel, J. Phys. Condens. Matter 15, R699 (2003).
- [13] H. Schiessel, *Biophysics for Beginners: A Journey through the Cell Nucleus* (Taylor & Francis, Boca Raton, 2013).
- [14] Http://users.soe.ucsc.edu/ pinal/P1.html (University of California, Santa Cruz).
- [15] D. Voet and J. G. Voet, Biochemistry (Wiley, New York, 1995), 2nd ed.
- [16] Y. C. Fung, *Biomechanics: Mechanical properties of living tissues* (Springer-Verlag, New York, 1993), 2nd ed.
- [17] M. A. Meyers, P.-Y. Chen, A. Y.-M. Lin, and Y. Seki, Prog. Mater. Sci. 53, 1 (2008).
- [18] P.-Y. Chen, J. McKittrick, and M. A. Meyers, Prog. Mater. Sci. 57, 1492 (2011).
- [19] E. A. Evans and R. Skalak, *Mechanics and thermodynamics of biomembranes* (CRC Press, Inc., Boca Raton, Florida, 1980).
- [20] D. H. Boal, Mechanics of the cell (Cambridge University Press, Cambridge, 2002).
- [21] G. Bao and S. Suresh, Nat. Mater. 2, 715 (2003).
- [22] R. M. Hochmuth, J. Biomech. 33, 15 (2000).
- [23] J. Sleep, D. Wilson, R. Simmons, and W. Gratzer, Biophys. J. 77, 3085 (1999).
- [24] A. L. Givan, ed., Flow Cytometry: First Principles (John Wiley & Sons, New York, 2001), 2nd ed.
- [25] D. A. Fedosov, H. Noguchi, and G. Gompper, Biomech. Model. Mechanobiol. 13, 239 (2014).
- [26] R. G. Larson, *The structure and rheology of complex fluids* (Oxford University Press, Oxford, NY, 1999).
- [27] E. W. Merrill, E. R. Gilliland, G. Cokelet, H. Shin, A. Britten, and R. E. Wells Jr, Biophys. J. 3, 199 (1963).
- [28] S. Chien, S. Usami, H. M. Taylor, J. L. Lundberg, and M. I. Gregersen, J. Appl. Physiol. 21, 81 (1966).
- [29] S. Chien, S. Usami, R. J. Dellenback, and M. I. Gregersen, Am. J. Physiol. 219, 143 (1970).
- [30] L. Lanotte, J. Mauer, S. Mendez, D. A. Fedosov, J.-M. Fromental, V. Claveria, F. Nicoud, G. Gompper, and M. Abkarian, Proc. Natl. Acad. Sci. USA 113, 13289 (2016).
- [31] H. C. Berg, E. coli in Motion (Springer, New York, 2004).
- [32] J.-F. Joanny and J. Prost, HFSP J. 3(2), 94 (2009).
- [33] M. C. Marchetti, J. F. Joanny, S. Ramaswamy, T. B. Liverpool, J. Prost, M. Rao, and R. A. Simha, Rev. Mod. Phys. 85, 1143 (2013).

- [34] T. E. Angelini, E. Hannezo, X. Trepat, M. Marquez, J. J. Fredberg, and D. A. Weitz, Proc. Natl. Acad. Sci. USA 108, 4714 (2011).
- [35] M. E. Cates and J. Tailleur, EPL 101, 20010 (2013).
- [36] M. E. Cates and J. Tailleur, Annu. Rev. Condens. Matter Phys. 6, 219 (2015).
- [37] D. Bi, J. H. Lopez, J. M. Schwarz, and M. L. Manning, Nat. Phys. 11, 1074 (2015).
- [38] E. Lauga and T. R. Powers, Rep. Prog. Phys. 72, 096601 (2009).
- [39] T. Ishikawa, J. Royal Soc. Interface 6, 815 (2009).
- [40] D. L. Koch and G. Subramanian, Annu. Rev. Fluid Mech. 43, 637 (2011).
- [41] J. Elgeti, R. G. Winkler, and G. Gompper, Rep. Prog. Phys. 78, 056601 (2015).
- [42] J. Elgeti, R. G. Winkler, and G. Gompper, SoftComp Newsletter 12 (2015), URL http://www.eu-softcomp.net.
- [43] J. Elgeti and G. Gompper, Eur. Phys. J. Special Topics 225, 2333 (2016).
- [44] L. Alvarez, B. M. Friedrich, G. Gompper, and U. B. Kaupp, Trends Cell Biol. 24, 198 (2014).
- [45] M. Basan, T. Risler, J.-F. Joanny, X. Sastre-Garau, and J. Prost, HFSP J 3(4), 265 (2009).
- [46] M. A. Schwartz, Science 323(5914), 588 (2009).
- [47] M. A. Wozniak and C. S. Chen, Nat. Rev. Mol. Cell Biol. 10(1), 34 (2009).
- [48] D. Drasdo and S. Hoehme, New J. Phys. 14 (2012).
- [49] A. Taloni, A. A. Alemi, E. Ciusani, J. P. Sethna, S. Zapperi, and C. A. M. La Porta, PLoS One 9(4), e94229 (2014).
- [50] J. Ranft, M. Basan, J. Elgeti, J.-F. Joanny, J. Prost, and F. Jülicher, Proc. Natl. Acad. Sci. USA 107, 20863 (2010).
- [51] FAS Public Interest Report, The Journal of the Federation of American Scientists, Volume 56, Number 2 (Summer 2003). http://www.fas.org/faspir/2003/v56n2/challenge.htm.
- [52] FYI: The AIP Bulletin of Science Policy News, Number 63: June 21, 2007. http://www.aip.org/fyi/2007/063.html.
- [53] U. Sefert, Rep. Prog. Phys. 75, 126001 (2012).
- [54] D. Kuzum, R. G. D. Jeyasingh, B. Lee, and H.-S. P. Wong, Nano Lett. 12, 2179 (2012).
- [55] A. Calimera, E. Macii, and M. Poncino, Funct. Neurol. 28, 191 (2013).

A 1 Super-resolution in Optical Microscopy

S. U. Egelhaaf

Condensed Matter Physics Laboratory Heinrich Heine University, Düsseldorf

Contents

1	Introduction		
2	What Detail can be Observed in a Microscope? 2.1 Resolution 2.2 Localization	2 2 4	
3	How can Details be Observed in a Microscope? 3.1 Localization microscopy 3.2 Stimulated emission depletion microscopy 3.3 Structured illumination microscopy	4 6 8 9	
4	Summary and Outlook	12	
A	Confocal microscopy	13	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Our eyes are an important tool to explore the world around us. The visual inspection of samples also plays an important role in science. Without any tool this is limited to features larger than about 0.1 mm. Conventional optical microscopy allows us to observe structures or objects that are almost three orders of magnitude smaller, about half the wavelength of visible light. This limit has been pushed to much smaller structural details by recent technical advances. Super-resolution microscopy not only allows the visualization of structures with a very high spatial resolution, down to a few nanometers, but can also provide quantitative information. This renders optical microscopy and its super-resolution variants increasingly powerful to investigate biological and soft matter systems.

The resolution of an optical microscope, i.e. the possibility to distinguish two nearby features, mainly depends on the objective. This is limited by diffraction. The recent super-resolution techniques circumvent the diffraction limit by using photochemistry, photophysics and special illumination schemes. In the following, first the resolution of a microscope and the accuracy of the localization of individual objects will be discussed. Subsequently, different super-resolution techniques are presented with a focus on the underlying principles rather than photophysical, photochemical or technical aspects. For details we refer to the literature [1-10].

2 What Detail can be Observed in a Microscope?

2.1 Resolution

Resolution is the ability to distinguish two nearby features. It is quantified by the minimum separation at which two features can just be resolved. If the image of a point source was a point, two features could be infinitely close and would still be resolved. However, the resolution of a microscope is limited due to the wave-like nature of light. Because of diffraction, the image of an object is broadened compared to the original object [2-4]. Diffraction occurs when light encounters an edge or inhomogeneity, for example a detail of the sample or the aperture of a lens. The diffraction pattern consists of a central maximum and several higher order maxima whose positions depend on the wavelength of the light and the size of the feature or aperture. Thus, diffraction leads to some blurring that renders it impossible to obtain a sharp image.

Illuminating a circular aperture, e.g. the aperture of an objective, or imaging a point source results in a central spot surrounded by a series of concentric rings of decreasing intensity, the so-called Airy pattern (Fig. 1a). The intensity profile of an aperture with radius a is given by

$$I(\theta) = I_0 \left(\frac{2J_1(x)}{x}\right)^2 \tag{1}$$

where I_0 is the intensity at the center, $J_1(x)$ is the Bessel function of the first kind of order one, $x = (2\pi/\lambda)a\sin\theta$ with λ the wavelength of the light and θ the observation angle. Thus, in a perfect but finite microscope the image of a point source is an Airy pattern. The pattern is similar but broader in a real microscope due to, e.g., aberrations. This extended three-dimensional image of a point source is represented by the point spread function (PSF) [5]. The PSF takes into account diffraction, aberrations and other effects and hence characterizes the response of a specific microscope with all its optical elements to a point source. Typically the PSF can be well approximated by a Gaussian distribution.


Fig. 1: (top) One-dimensional cross-section and (bottom) two-dimensional representation of the point spread function (PSF) of (a) an individual point source, (b) two well-separated point sources, (c) two point sources at a separation fulfilling the Rayleigh criterion, (d) the Sparrow criterion, (e) the Abbe criterion and (f) two point sources that are too close to be distinguished. The individual PSFs are shown as blue lines, the sums of two PSFs as red lines.

In the case of two point sources, each produces a blurred image, the individual PSF (Fig. 1a,b). If they are close together, the two PSFs overlap and hence it might be difficult or even impossible to resolve them, i.e. to confidently conclude the existence of two distinct sources and to determine their positions. According to the Rayleigh criterion, the images of two point sources are considered just resolvable if the central maximum of the PSF of one source coincides with the first minimum of the PSF of the second source and thus the total intensity shows an intensity minimum between the two maxima (Fig. 1c). This requirement results in a minimum separation of the two sources, $d_{\rm res} = 0.61\lambda/{\rm NA}$, where NA = $n\sin\theta$ the numerical aperture of the objective, *n* the refractive index of the medium between the object and objective and 2θ the angle of the light cone accepted by the objective [2]. The equation assumes that the objective also serves as condenser, as is usual in fluorescence microscopy, or otherwise NA = $(NA_{obj} + NA_{cond})/2$. Another criterion, the Sparrow criterion, requires that the minimum between the two maxima in the sum of the two PSFs just disappears (Fig. 1d). This results in $d_{\rm res} = 0.47\lambda/{\rm NA}$ and hence a similar relationship [2].

Another argument describes the sample as a complex grating with various distances between the individual 'slits'. Two parallel slits represent the simplest component of this complex grating. (But any structure, even a non-periodic structure, can be decomposed into such basic components using a Fourier transform.) When these two slits are illuminated, the intensity distribution shows a central (zero order) maximum and, on both sides, further maxima of decreasing intensity (and increasing order). If only one maximum contributes to the image, there is no contrast in the image and hence the slits are not 'visible'. Therefore, to determine the distance of the two slits, at least two orders, usually the central and first maxima, must be detected and hence enter the objective. This criterion leads to the minimum resolvable separation between the two slits, $d_{\rm res} = 0.5\lambda/\rm NA$ (Fig. 1e) [4]. This is Abbe's diffraction limit. This relation for the minimum resolvable separation between two objects, $d_{\rm res}$, is similar to the ones discussed above

$$d_{\rm res} = 0.47 \frac{\lambda}{\rm NA} \qquad \text{Sparrow criterion}$$

$$d_{\rm res} = 0.50 \frac{\lambda}{\rm NA} \qquad \text{Abbe's limit} \qquad (2)$$

$$d_{\rm res} = 0.61 \frac{\lambda}{\rm NA} \qquad \text{Rayleigh criterion.}$$

Further measures for the resolution limit consider the width of the PSF, e.g., its full width at half maximum (FWHM) or at $1/e^2$ of the maximum, and result in similar relations.

All criteria indicate that the resolution can be increased by a higher numerical aperture NA and a lower wavelength λ . They only differ in the prefactor due to the different definitions of the detection limit of two neighboring objects. Typical values NA < 1.5 and $\lambda \approx 500$ nm suggest that for diffraction-limited optics the spatial resolution is limited to about 200 nm in lateral direction (and about 500 nm in axial direction) [4, 5].

2.2 Localization

To localize an object, the centre of its PSF needs to be determined [11]. The localization accuracy in lateral direction, d_{loc} , is given by

$$d_{\rm loc}^2 \approx \frac{d^2}{N} + \frac{a^2}{12N} + \frac{8\pi d^4 b^2}{a^2 N^2}$$
(3)

where N is the number of photons collected from an individual source, typically a fluorescent molecule, d and b are the standard deviations of the PSF and the background signal, respectively, and a is the pixel size of the camera [2, 3]. The terms describe the effects of the photon noise and the PSF, the finite pixel size and the background noise, respectively. Typically, the accuracy of the localization is dominated by the width of the PSF, $d \approx d_{\rm res}$, and the number of detected photons N [4]. The accuracy increases with the number of detected photons, which can reach a million and hence $d_{\rm loc} \approx d_{\rm res}/\sqrt{N} \ll d_{\rm res}$; localization accuracies $d_{\rm loc} \leq 10$ nm have been achieved. Thus, individual objects can be localized with an accuracy much higher than the resolution limit, i.e. the minimum separation that allows one to distinguish two objects (Sec. 2.1).

The localization accuracy can approach the separation between neighbouring fluorophores. Thus, the image quality can be limited by the labelling density n. Too low a labelling density and hence too large a distance between fluorophores results in artifacts, for example homogeneous structures might appear heterogeneous. The effect of the labelling density n is quantified by the Nyquist criterion [3]; features smaller than twice the separation between fluorophores cannot reliably be resolved. Therefore, the minimum resolvable feature size is

$$d_{\text{Nyquist}} = \frac{2}{n^{1/D}} \tag{4}$$

with D the dimension of the feature to be imaged. For example, to obtain a 10-nm resolution in a two-dimensional image, a labelling density $n = 40 \times 10^3 / \mu m^2$ is required. In the following it is assumed that the labelling density is high enough to not limit image quality.

3 How can Details be Observed in a Microscope?

The above considerations suggest that the position of *one* single object, i.e. the centre of the PSF, can be determined with an accuracy of about $d_{loc} \approx 10$ nm or even better, whereas *two* objects can only be resolved if their distance is at least $d_{res} \approx 200$ nm, the diffraction limit or essentially the width of the PSF. To achieve a high accuracy in the location of objects, therefore, the width of their PSF, i.e. the size of their individual images, needs to be smaller than the separation

between two objects. Thus, the aim is to reduce the effective width of the PSF or to increase the effective object separation.

This is possible by exploiting the special features of fluorescence. Fluorescent molecules absorb light over a certain wavelength range, the excitation wavelengths, and emit at longer wavelengths, the emission wavelengths, where the emission might be spontaneous or stimulated through radiation of a specific wavelength. The diffraction limit applies to the excitation and emission; the excitation beam cannot be focussed tighter than to a diffraction-limited spot, while the emission from a point source expands to a diffraction-limited PSF. Therefore, both, appropriate illumination patterns and detection schemes, can be exploited to achieve super-resolution. The main methods will be explained below. Some of them allow for wide-field imaging whereas others require to scan the sample. Since confocal microscopy also requires scanning and hence shares some features with these methods, a brief description of confocal microscopy is provided in the Appendix A (although it only provides diffraction limited images and hence is not a super-resolution technique).



Fig. 2: Principle of spectral precision distance microscopy (SPDM). (a) A diffraction-limited image does not allow the separation of individual objects. (b) Thus, the objects are labelled by different kinds of fluorophores, here fluorophores with different emission spectra. (c) Filters are used to separate the contributions of the different fluorophores. (d) Now they can be localized with high accuracy. (e) Subsequently the locations are combined to yield a super-resolution image.

3.1 Localization microscopy

If the distance between two objects is smaller than the resolution criterion (Eq. 2) their images overlap and the objects cannot be resolved (Fig. 1f). However, if the two objects can be distinguished, e.g. by their different absorption or emission spectra or fluorescent lifetimes, their images can be separated, e.g. by a filter. This is the basic idea of spectral precision distance microscopy (SPDM) [6]. For example, if the emission spectra are different, i.e. the fluorophores can be distinguished by their colours (Fig. 2a,b), filters can be used to separate their signals (Fig. 2c). The densities of fluorophores with the same characteristics need to be low enough that any two fluorophores of the same kind are separated further than the diffraction limit. Depending on the object density, this may require many different fluorophores. If the densities are low enough, the fluorophores can be localized with high accuracy (Fig. 2d) and subsequently the images can be combined to yield the complete image (Fig. 2e).



Fig. 3: Principle of single molecule localization microscopy (SMLM), e.g. photoactivated localization microscopy (PALM) or stochastic optical reconstruction microscopy (STORM). (a) Objects are labelled by identical fluorophores and hence a diffraction-limited image does not allow the separation of the individual objects. (b) However, if the fluorophores are activated at different times, the effective fluorophore concentration is decreased and hence their effective separation increased. (c) Thus they can localized with high accuracy. Repeating the photoactivation or photoswitching process many times, all fluorophores can be imaged and localized. (d) Subsequently the locations are combined to yield a super-resolution image.

Instead of distinguishing fluorophores by their spectral characteristics, they can also be separated in time, i.e. different fluorophores but with the same characteristics are active at different times. Hence, at any instant, only a subset of the fluorophores contributes to the image (Fig. 3a,b) [3, 4, 7]. This results in a lower density of fluorophores in the images and hence a larger separation between them. Localization in the individual images thus can be very accurate (Fig. 3c) and subsequently composed to yield a super-resolution image (Fig. 3d). Thus, a high resolution is achieved through the localization of single molecules which are randomly activated. This is exploited in single molecule localization microscopy (SMLM).

The central point is the temporal separation of fluorescence. This can be achieved by using photoactivation, photoswitching, or the spontaneous blinking of fluorophores (Fig. 4) [7-10]. Photoactivatable fluorophores initially are in a non-active state but can be activated through illumination. Photoswitchable fluorophores can reversibly switch between a non-active and an active state with the conversion also controlled through illumination. These fluorophores can also stochastically transition between non-active and active emission states by exposure to light of specific wavelengths. Their switch probability can be small enough that a low level of light results in only a single active fluorophore within a diffraction-limited area at any instant. Furthermore, the emission during a single fluorescence burst can be high enough to allow for an accurate localization and subsequently photobleaching can be fast enough to allow for the activation of other fluorophores at a high rate. In any case, initially the fluorophores are in a non-active state and, after activation by light, a small fraction of them converts into the active state. Photoactivation and subsequent photobleaching is applied in photoactivated localization microscopy (PALM), while photoswitching and photoblinking are exploited in stochastic optical reconstruction microscopy (STORM). Typical spatial resolutions are about 20 to 50 nm.



Fig. 4: *Key concept of single molecule localization microscopy (SMLM), e.g. photoactivated localization microscopy (PALM) or stochastic optical reconstruction microscopy (STORM). (a) Initially the fluorophores are in a non-fluorescent state (black circles). (b) A low level of activation light (red arrow) is applied so that a few individual fluorophores are stochastically activated or switched to the fluorescent state. The activated fluorophores produce diffraction-limited images (blue) which can be analyzed to accurately localize the activated fluorophores. This process is repeated to create a super-resolution image.*

Because only a very small fraction of the fluorophores are activated at any instant and hence many individual images, typically hundreds to thousands, are required to complete the entire image, the time resolution is very poor, a few minutes for a complete image.

3.2 Stimulated emission depletion microscopy

In stimulated emission depletion microscopy (STED), the possibility to turn on and off fluorescence is used to reduce the volume that can contribute to the image. Fluorophores in the periphery of the diffraction-limited focal spot are depleted and hence they cannot interfere with the localization of the central fluorophore [4, 8]. Fluorescence in the periphery is suppressed by stimulated emission through a donut-shaped illumination with a wavelength that corresponds to the energy gap of the fluorophore (Fig. 5). Using stimulated emission, the fluorophores are forced from the excited to the ground state within their fluorescence lifetime, i.e. before spontaneous emission can occur [10]. Since the wavelengths of the stimulating illumination and the stimulated emission are the same, they can be separated from spontaneous emission using a filter.



Fig. 5: *Key concept of stimulated emission depletion microscopy (STED). (a) Point spread function (PSF, blue) as observed upon illumination of a fluorophore, (b) concentric donut-shaped depletion beam (red) that suppresses spontaneous emission and (c) combination of the excitation and depletion beams to produce a small central active region.*

The donut-shaped illumination only depletes the peripheral signal, whereas the center is not illuminated and hence fluorescence is preserved. Upon increasing the intensity of the depletion beam, the depleted region expands but the center remains largely unaffected. Thus, the size of the unaffected central region and hence the resolution can be tuned. The full width at half maximum (FWHM) of the effective focal spot is given by [2, 3, 8]

$$d_{\rm STED} = \frac{\lambda}{2\,\rm NA}\sqrt{1+I/I_{\rm s}}\tag{5}$$

where I is the maximum intensity of the stimulating illumination and I_s the saturation intensity of the fluorophore. In principle, the resolution of STED is unlimited (as long as the sample is not severely damaged). Typically a resolution of 20 to 100 nm is achieved. Since only one spot is recorded at any instant, the whole sample must be scanned (Fig. 6), which lowers the acquisition speed. The imaging approach hence is conceptually closely related to confocal imaging (Appendix A). As in confocal microscopy, no further image analysis is needed.

Depletion can also be achieved by photoswitching; a donut-shaped beam switches fluorophores into the non-active state instead of inducing stimulated emission. This process, reversible saturable optical linear fluorescence transitions (RESOLFT), can reduce the damaging effect a depletion beam with a high intensity may have [3, 8].



Fig. 6: Principle of (a) stimulated emission depletion microscopy (STED) and (b) confocal laser scanning microscopy (CLSM). The diffraction-limited excitation beam (blue) and, only in STED, the concentric donut-shaped depletion beam (red) are scanned across the sample (objects are represented as black circles) resulting in super-resolution and diffraction-limited images, respectively. Note the different scanning densities (arrows) reflecting the different effective point spread functions (PSFs).

3.3 Structured illumination microscopy

In structured illumination microscopy (SIM) also the illumination is modified to increase the resolution [10]. The sample is illuminated with a periodic sinusoidal pattern (Fig. 7). The im-



Fig. 7: Key concept of structured illumination microscopy (SIM). The sample structure, represented by a grating with (a) a periodic and (b) a non-periodic structure (black), is illuminated with (c,d) a periodic pattern with periodicity ΔL_p (red) which yields (e,f) an interference pattern, a so-called Moiré pattern. The interference pattern has stripes with a larger separation, $\Delta L' > \Delta L_p$. Since this interference pattern is imaged, smaller structures in the sample can still be resolved yielding super-resolution images.



Fig. 8: Principle of structured illumination microscopy (SIM). (a) The sample (objects represented as black circles) is illuminated by a periodic sinusoidal pattern (red) that is shifted and rotated at least three times. (b) This results in several recordings with diffraction-limited images of the objects (blue). (c) The information is extracted and yields a super-resolution image. (d) For comparison, the image obtained by conventional microscopy is also shown.

age is an interference pattern between the imposed periodic pattern and the sample, a so-called Moiré pattern. (As above, the sample can be thought of as being composed of periodic patterns with different frequencies.) Due to this interference, high frequency or short wavelength, ΔL_s , information is shifted to lower frequencies or larger wavelength, $\Delta L'$, which might now pass through the optical system. By shifting and rotating the imposed pattern and taking several images, an image containing the high frequency information shifted back to their original frequencies can be reconstructed (Fig. 8). Its resolution depends on the highest frequency present, which is determined by the frequency of the pattern. The frequency of the pattern is also limited by diffraction. Thus, the maximum improvement of the resolution is a factor of two (Fig. 7), resulting in a resolution of about 100 nm. Although the improvement is relatively modest, SIM has the advantage that it does not require specific labelling but only a grating or a spatial light modulator in one of the image planes of a standard epifluorescence microscope with a laser. Furthermore, wide-field observation allows for a relatively fast image acquisition, although several (typically nine or fifteen) images with different orientations of the periodic pattern have to be taken [10].

The contrast of the imposed pattern is only high in the focal plane (Fig. 9). Thus, the fluorophores that lie in the focal plane are illuminated by a pattern, whereas the fluorophores outside the focal plane essentially are homogeneously illuminated. Shifting the pattern will hence only change the illumination of the fluorophores in the focal plane whereas the fluorophores outside the focal plane are not affected by a shift (Fig. 9a). The axial locations of the fluorophores thus can be distinguished by considering images with different positions of the pattern. Hence optical sectioning can be achieved by taking this into account in the analysis. Moreover, if the imposed pattern is also modulated in the axial direction, for example by three-beam interference, the axial resolution can also be improved by a factor of two, to about 300 nm.



Fig. 9: Principle of axial resolution improvement by structured illumination microscopy (SIM) with optical sectioning. In contrast to previous figures, this cut contains the axial direction (vertical direction). (a) The sample is illuminated by a periodic sinusoidal pattern (red) that is shifted; here it is shifted horizontally by a third of a period (as in Fig. 8, top row). The contrast of the pattern is only high in the focal plane. In the different images, hence objects close to the focal plane have different intensities (indicated by the brightness of the diffraction-limited images and the black circles for the non-fluorescent objects), whereas objects outside the focal plane have the same intensity in all images. In-focus and out-of-focus objects thus can be distinguished. (b) The locations of the objects in the focal plane are determined and the out-of-focus objects are removed. This yields a super-resolution image with optical sectioning.

4 Summary and Outlook

Very recent advances in optical microscopy render super-resolution images possible. Their spatial resolution is significantly improved compared to conventional optical microscopy, in principle it is unlimited as Abbe's diffraction limit can be circumvented. Various experimental realizations indeed achieved an improvement by about an order of magnitude. Conventional microscopy is limited to a resolution of about 200 nm, whereas super-resolution techniques allow the localization of objects with an accuracy of about 20 nm and even 1 nm has been reached [1, 5]. Super-resolution techniques exploit the possibility to localize individual objects very accurately, as long as they can be distinguished from their neighbours; the centre of the point spread function (PSF) can be determined very accurately if its width is smaller than the separation between objects. This requires to reduce the width of the PSF or to increase the effective separation between objects. In single-molecule based stochastic methods, like PALM or STORM, the effective separation of objects is increased, whereas in deterministic ensemble-level methods, such as STED or SIM, the width of the PSF is reduced by applying a suitable illumination.

Super-resolution microscopy is developing into a very powerful tool to obtain high resolution images of biological and soft matter systems. It allows the investigation of structures and processes from the cellular or bulk level down to the molecular scale and hence of a broad range of length scales. Beyond the visual information, also quantitative structural and dynamic data can be obtained through advanced image processing [4, 11], where the complex image analysis of large quantities of data significantly profits from increasing computing power.

So far the molecular, nanometer scale was not accessible by optical microscopy and, e.g., electron microscopy had to be applied. For optical microscopy, usually, the sample preparation is less invasive and hence the samples can be studied closer to their native states. Nevertheless, most of the super-resolution methods require to introduce fluorophores, whose effects on the behaviour of the samples has to be considered carefully. The samples furthermore might need some fixation to avoid motion-induced blur if the movements of the observed structures during the measurement time are close to the desired (small) resolution. In addition, the samples often are exposed to continuous and long illumination with relatively high laser intensities. Thus, care has to be taken to avoid modifications of the sample structure or functionality as well as radiation damage during the measurements. It is crucial that what is observed with the microscope is indeed representative of the sample.

These exceptional advances could be achieved already shortly after the introduction of superresolution microscopy. Thus, further significant developments are expected in the future. They are expected to concern the imaging methods with improvements in the spatial as well as the temporal resolution but also, e.g., the contrast, the availability of several colours, the noninvasiveness, the imaging depth or the accessibility to a wider research community and non-experts. Moreover, advances in related areas might improve the super-resolution techniques. For example, super-resolution microscopy might benefit from the development and modification of new and improved fluorescent labels or the ability to combine super-resolution microscopy with other measurement or manipulation techniques, like correlative electron microscopy or optical tweezers. These advancements will make information available which so far is not accessible. They hence will enable new discoveries not only in biological and soft matter systems but also in other areas of science, e.g., in semiconductors and material sciences.

Appendix

A Confocal microscopy

A significant improvement over conventional optical microscopy is provided by confocal mi-It rejects out-of-focus light and thus achieves a better discrimination of features croscopy. along the optical axis, known as optical sectioning. Hence it mainly improves the axial resolution [12, 13]. Consider three object points which are located inside (Fig. 10, red point) and outside (green and grey points) the focal plane. The objective and tube lens accordingly form images inside and outside the intermediate image plane. In the intermediate image plane, only the image of the point in the focal plane is in focus (red). If an aperture is introduced in the intermediate image plane, the in-focus image is unaffected, while the out-of-focus images are essentially blocked (dotted green and grey lines) and thus the depth of field is significantly reduced. The aperture is in the *conjugate* plane to the objective *focal* point; it is a *'confocal'* aperture. The size of the confocal aperture ideally matches the size of the central part of the Airy pattern to collect a large fraction of the in-focus light while blocking most of the out-of-focus light. The axial resolution is significantly improved compared to conventional microscopy. Since the excitation and detection are point-like, also the lateral resolution is increased by a factor $1/\sqrt{2}$. However, although confocal microscopy improves the resolution, it does not provide resolution beyond the diffraction limit.

While the confocal aperture improves the resolution, only a single point can be observed at any instant. A complete image must thus be combined from individual points by scanning the sample (Fig. 6b). Lateral scanning in the sample plane is typically performed using two movable mirrors. In addition, the sample is also scanned in the axial direction to obtain a stack of individual sections with each of the sections having a significantly improved axial resolution. Based on these sections, a three-dimensional image of the sample is constructed. This renders confocal microscopy slower than conventional microscopy.



Fig. 10: Key concept of confocal laser scanning microscopy (CLSM). The confocal aperture removes most of the light from out-of-focus objects (grey and green), whereas the light from in-focus objects (red) can pass.

References

- [1] S.-W. Chu, Top. Appl. Phys. 129, 495 515 (2015)
- [2] E. Sezgin, J. Phys.: Condens. Matter 29, 273001 (2017)
- [3] D. M. Shcherbakova, P. Sengupta, J. Lippincott-Schwartz, V. V. Verkusha, Annu. Rev. Biophys. 43, 303 – 329 (2014)
- [4] T. J. Gould, S. T. Hess, J. Bewersdorf, Annu. Rev. Biomed. Eng. 14, 231 254 (2012)
- [5] L. Schermelleh, R. Heintzmann, H. Leonhardt, J. Cell Biol. 190, 165 175 (2010)
- [6] C. Cremer et al., Biotechnol. J. 6, 1037 1051 (2011)
- [7] G. Patterson, M. Davidson, S. Manley, J. Lippincott-Schwartz, Annu. Rev. Phys. Chem. 61, 345 – 367 (2010)
- [8] S. W. Hell, Science 316, 1153 1158 (2007)
- [9] S. Habuchi, Front. Bioeng. Biotech. 2 20 (2014)
- [10] A. Jost, R. Heintzmann, Annu. Rev. Mater. Res. 43, 261 282 (2013)
- [11] M. C. Jenkins, S. U. Egelhaaf, Adv. Coll. Interface Sci. 136, 65 92 (2008)
- [12] C. J. R. Sheppard, D.M. Shotton, *Confocal Laser Scanning Microscopy* (Bios Scientific Publishers, Springer, New York, 1997)
- [13] J. Pawley, Handbook of Biological Confocal Microscopy (Springer, Berlin, 2006)

A 2 3D structure determination using electron cryo-microscopy

C. Sachse Structural and Computational Biology Unit EMBL (European Molecular Biology Laboratory) Heidelberg and Ernst-Ruska-Centre 3 Forschungszentrum Jülich GmbH

Contents

1	Intr	oduction	2
2	Sam	nple preparation	2
	2.1	Negative stain	2
	2.2	Specimen embedding in vitreous ice	3
3	Ima	ges from the electron microscope	3
	3.1	Hardware components of electron microscopes	3
	3.2	Image formation	4
4	3D i	image reconstruction and atomic models	4
	4.1	Classification and initial model	4
	4.2	Projection matching and 3D reconstruction	5
	4.3	Interpretation of cryo-EM maps	5
Ref	erence	2S	5

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

On October 4th in 2017, the Nobel prize for chemistry was awarded to Jacques Dubochet, Joachim Frank and Richard Henderson for developing cryo-electron microscopy for the highresolution structure determination of biomolecules in solution. The recognition came at a time when several 100s of atomic resolution cryo-EM structures have been deposited to the EM databank every year [1] (Fig. 1 left). At this day, cryo-EM is capable to elucidate biological structures of very large assemblies well as single molecules. Traditionally, very large assemblies were imaged by cryo-EM made up from regular helical structures or 2D crystals of MDa in size whereas now single particles structures as small as ~60 kDa have been resolved at atomic resolution. Despite the power of the technique, the number of structures resolved every year is still low in comparison with other structural techniques such as X-ray crystallography and NMR (Fig. 1, right). In the current lecture, I will provide the principal foundations of cryo-EM including sample preparation, electron microscopes and the 3D image reconstruction techniques.



Fig. 1: Left. Number of cryo-EM maps deposited per year achieving different resolutions (from <u>www.pdbe.org/emstats</u>). In 2016, approx. 200 structures were determined below 4 Å in resolution, which allows de novo atomic model building. Right. Total number of structures deposited to Protein Data Bank (PDB) is 135,000 with 90 %, 9 % and 1 % from the disciplines of X-ray crystallography, NMR and EM respectively.

2 Sample preparation

2.1 Negative stain

After the development of the electron microsope by Ernst Ruska in the 1930s [2], nanostructures from biological samples could be readily visualized using these microscopes. In 1959, Sidney Brenner showed that large helical plant viruses like tobacco mosaic virus (TMV) were imageable at high contrast when dried in negative stain solution like uranyl acetate [3]. Although the approach dehydrates the molecule of interest and thus leads to collapse of the native fine structure of the molecule, it reveals overall structural information of the embedded molecule by forming a negative imprint in the stain (Fig. 2, left). Negative stain is now commonly used for quick and routine initial sample characterization.

2.2 Specimen embedding in vitreous ice

Although negative stain embedding is a simple preparation procedure, specimen preservation is poor and the molecular structural details are lost due to dehydration. In the 1980s, Jacques Dubochet and colleagues devised a method to embed macromolecules in vitrified water by plunge-freezing them in liquid ethane [4]. This way, the biological molecules remain hydrated in their near-native state while being immobilized within a glass-like layer of ice, which enables direct imaging of molecules in the vacuum of electron microscopes (Fig. 2, right).



Fig. 2: Electron micrographs of tobacco mosaic virus (TMV). Left. Embedded in negative stain. Right. Image of cryo-frozen specimen.

3 Images from the electron microscope

Individual images of biomolecular structures from the electron microscope are noisy when taken from negatively stained or vitrified specimens. High doses as they would be required to produce higher contrast images of biological molecules result in the destruction of the molecule by the electron beam. To avoid radation damage, individual images need to be acquired at very low doses and consequently at poor contrast. Enhanced detail can only be retrieved by averaging many identical views of the molecule [5].

3.1 Hardware components of electron microscopes

The construction principle of electron microscopes resembles that of a light microscope with the difference that electrons are used for imaging instead of light [6]. Due to the property of electrons that show strong interaction with matter, the inside of the electron microscope needs to be kept in vacuum. As electrons have much smaller wavelengths than light, they are highly suitable for imaging of molecular up to atomic details. They are accelerated at the gun and focussed into a parallel beam by the condenser lens. The electron beam hits the sample, is scattered and converted into an image by the object lens. Finally, the image contains a projection of the specimen with contributions from the inside and outside of the molecule, is recorded by an electron detector and available in digital format.

3.2 Image formation

Due to the poor visibility of biomolecules in ice, most of the image constrast has to be generated by defocussing the specimen. In this way, the actual image is predominantly formed by the interference of waves from the scattered and unscattered beam by what is referred to as phase contrast. The mathematical description of the contrast transfer function (CTF) [7] provides a way to correct for the iamge distortions that are induced by defocussing.

4 3D image reconstruction and atomic models

Once digitized many cryomicrographs are required to extract molecular information from the data. Depending on the sample arrangements different image processing software is available for two-dimensional crystals, helical assemblies, single particles and cellular structures. Most of the structures are currently determined by the single-particle approach, which I will outline in more detail.

4.1 Classification and initial model

The initial steps of characterization require isolation of particles from the micrograph. The micrographs or particle images need to be corrected for the CTF of the microsope. In order to characterize the sample, the resulting stack of particles will be classified. In the ideal case, it is expected that particles occur in random orientations on the grid. In reality, they often have a different molecular composition and/or conformation. Classification by reference-free algorithms such as multivariate statistical analysis (MSA) or k-means algorithms is commonly used to initially characterize the sample. Although caution should be taken, in case when references are available, cross-correlation based approaches can be used to identify the particles of interest in the data set. In addition, maximum likelihood approaches are now routinely used in electron microscopy image processing [8]. Obtaining an initial three-dimensional model from two-dimensional projections is not always straightforward and methods of random conical tilt and angular reconstitution have been developed [9], [10].



Fig. 3: Left. Summary of work-flow of image processing in cryo-EM: helix segmentation, CTF correction, classification and iterative projection matching cycle (Figure from [11]). Right. Cryo-EM structure of TMV at 3.3 Å resolution in side and top view [12]. Detail of α-helical density T107-E131.

4.2 Projection matching and 3D reconstruction

Once coarse orientations have been assigned to the particles using the initial model, they need to be further refined by iterative cycles of projection, alignment and reconstruction until convergence [13] (Fig. 3, left). In order to obtain high-resolution structures, the accurate assignment of five parameters of x and y translation and three Euler angles is essential. The Fourier shell correlation is used as a metric for resolution [14]. In order to make structural details visible in the reconstructions, the maps need to be corrected for decay of amplitudes by the B-factor [15].

4.3 Interpretation of cryo-EM maps

Regardless of the obtained resolution, the aim of the cryo-EM experiment is the interpretation of the density by an atomic model. At 20 Å resolution, the overall envelope of the molecule can be discerned and atomic structures obtained from NMR and X-ray crystallography can be fitted as rigid bodies. This resolution is considered the practical limit for 3D reconstructions obtained from negative stain EM specimens. As molecules remain hydrated in cryo-EM, it can go to resolutions as high as 1.9 Å [16]. Below 10 Å resolution, secondary structure such as α -helices can be identified and atomic models can be adjusted to best fit the EM density. At resolutions below 4.0 Å, de novo atomic model building becomes possible by threading the amino acid sequence into the density (Fig. 3, right). Model building is generally followed by coordinate refinement including prior knowledge of atomic bonds and angles from standard amino acids and secondary structures in proteins as restraints to improve the accuracy of the presentation of these structures. Atomic models are frequently used by non-structural biologists to test biological function in vivo as the structures often directly reveal molecular mechanisms.

References

- [1] A. Patwardhan, "Trends in the Electron Microscopy Data Bank (EMDB).," *Acta Crystallogr D Struct Biol*, vol. 73, no. 6, pp. 503–508, Jun. 2017.
- [2] E. Ruska, ... LECTURES, PHYSICS 1986, THE DEVELOPMENT OF THE ELECTRON MICROSCOPE AND OF ELECTRON MICROSCOPY, DECEMBER 8, 1986 BY Nobel Lectures, 1995.
- [3] S. Brenner and R. W. Horne, "A negative staining method for high resolution electron microscopy of viruses.," *Biochim Biophys Acta*, vol. 34, pp. 103–110, Jul. 1959.
- [4] M. Adrian, J. Dubochet, J. Lepault, and A. W. McDowall, "Cryo-electron microscopy of viruses.," *Nature*, vol. 308, no. 5954, pp. 32–36, Mar. 1984.
- [5] R. Markham, S. Frey, and G. J. Hills, "Methods for the enhancement of image detail and accentuation of structure in electron microscopy," *Virology*, vol. 20, pp. 88–102, 1963.
- [6] R. M. Glaeser, "Electron Crystallography of Biological Macromolecules Chapter 5," p. 496, 2007.
- [7] H. P. Erickson and A. Klug, "Fourier transform of an electron micrograph: effects of defocussing and aberrations, and implications for the use of underfocus contrast enhancement," *Bericht Bunsen Gesell*, vol. 74, no. 11, pp. 1129–&, 1970.
- [8] F. J. Sigworth, P. C. Doerschuk, J. M. Carazo, and S. H. W. Scheres, An Introduction to Maximum-Likelihood Methods in Cryo-EM, 1st ed., vol. 482. Elsevier Inc., 2010, pp. 263–294.
- [9] M. Radermacher, "Three-dimensional reconstruction of single particles from random

and nonrandom tilt series," *J Electron Microsc Tech*, vol. 9, no. 4, pp. 359–394, Aug. 1988.

- [10] M. van Heel, "Angular reconstitution: a posteriori assignment of projection directions for 3D reconstruction.," *Ultramicroscopy*, vol. 21, no. 2, pp. 111–123, 1987.
- [11] S. A. Fromm and C. Sachse, "Cryo-EM Structure Determination Using Segmented Helical Image Reconstruction.," *Meth Enzymol*, vol. 579, pp. 307–328, 2016.
- [12] S. A. Fromm, T. A. M. Bharat, A. J. Jakobi, W. J. H. Hagen, and C. Sachse, "Seeing tobacco mosaic virus through direct electron detectors.," *J Struct Biol*, vol. 189, no. 2, pp. 87–97, Feb. 2015.
- [13] P. Penczek, M. Radermacher, and J. Frank, "Three-dimensional reconstruction of single particles embedded in ice," *Ultramicroscopy*, vol. 40, no. 1, pp. 33–53, 1992.
- [14] G. Harauz and M. van Heel, "Exact filters for general geometry three dimensional reconstruction," *OPTIK.*, vol. 73, no. 4, pp. 146–156, 1986.
- [15] P. B. Rosenthal and R. Henderson, "Optimal determination of particle orientation, absolute hand, and contrast loss in single-particle electron cryomicroscopy.," *J Mol Biol*, vol. 333, no. 4, pp. 721–745, Oct. 2003.
- [16] A. Merk, A. Bartesaghi, S. Banerjee, V. Falconieri, P. Rao, M. I. Davis, R. Pragani, M. B. Boxer, L. A. Earl, J. L. S. Milne, and S. Subramaniam, "Breaking Cryo-EM Resolution Barriers to Facilitate Drug Discovery.," *Cell*, vol. 165, no. 7, pp. 1698– 1707, Jun. 2016.

A 3 Fluorescence-based Experimental Techniques: FCS, FRAP, and FRET

J.K.G. Dhont

Soft Condensed Matter Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Introduction		
	1.1	The dynamics of macromolecules	2
	1.2	What is fluorescence?	5
2	Fluorescence Correlation Spectroscopy (FCS)		7
	2.1	The principles of FCS	8
	2.2	Diffusion of spherical colloids and proteins through isotropic and nematic net-	
		works of rods	12
	2.3	Complex formation and protein-membrane adsorption	15
3	Fluo	rescence Recovery After Photo-bleaching (FRAP)	16
	3.1	The principles of FRAP	17
	3.2	Measuring second moments of the fluorescence intensity:	
		a geometry-independent FRAP analysis	18
	3.3	FRAP experiments using an oscillatory reading intensity:	
		the glass transition in binary colloids	22
4 Förster Resonance Energy Transfer (FRET)		ter Resonance Energy Transfer (FRET)	23
	4.1	The principles of FRET	23
	4.2	FRET-efficiencies obtained from the simultaneous measurement of donor and	
		acceptor fluorescence intensities	25
	4.3	Unfolding of phosphoglycerate kinase (PGK) induced by guanidine hydrochlo-	
		ride (GndHCl)	27

Lecture Notes of the $49^{\rm th}$ IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

This chapter presents a discussion on the principles of various experimental fluorescence-based techniques, which are used to investigate the dynamics of macromolecules. The experimental techniques that will be addressed are Fluorescence Correlation Spectroscopy (commonly abbreviated as FCS), Fluorescence Recovery After Photo-bleaching (FRAP), and Förster Resonance Energy Transfer (FRET). These techniques probe, for example, the self- and collective diffusive motion of macromolecules, reaction-binding and adsorption rates, and the internal dynamics of molecules with a complex architecture, like polymers and proteins. In this chapter, the principles of each of these three fluorescence techniques will be disussed, and a few examples will be given of measurements on particulate systems.

Some knowledge on diffusion and fluorescence is necessary in order to quantitatively interpret experimental data from FSC, FRAP and FRET. In the following two subsections we will therefore present an elementary treatment of diffusion of macromolecules and fluorescence.

1.1 The dynamics of macromolecules

Macromolecules are by definition molecules that are much larger in size as compared to the molecules of the solvent (usually water) in which they are embedded. Thermal collisions of solvent molecules with a macromolecule leads to a random motion of its center-of-mass, its orientations, and possibly its conformational state. The corresponding translational and orientational motion is commonly referred as *Brownian motion* (after the Scottish botanist who observed the thermal motion of large molecules for the first time in 1827, although at that time the origin of this phenomenon was not at all understood). Brownian motion is nothing else than the thermal motion of the relatively large macromolecule. This introduction is restricted to rigid, spherical particles for which orientational Brownian motion is irrelevant. Such a macromolecule is commonly referred to as a *colloid*.

The degree of thermal motion can be quantified by the so-called *mean-squared displacement*. Consider a colloids that resides with its center at the origin at, say, time t = 0. Let $\mathbf{r}(t)$ (with t > 0) denote the vector, pointing to the center of the colloid, with a magnitude that is equal to distance of the center from the origin. The instantaneous value of $\mathbf{r}(t)$ depends on the initial configuration of solvent molecules that cause the Brownian displacements. For each repetition of an experiment that measures the instantaneous location $\mathbf{r}(t)$ of the colloid, the initial configuration of solvent molecules is different, so that a different value for $\mathbf{r}(t)$ will be recorded. Averaging any function $f(\mathbf{r}(t))$ of $\mathbf{r}(t)$ over many such experiments is an average over all possible realizations of thermal collisions of solvent molecules with the colloid. Such a thermal average will be denoted by the brackets $< f(\mathbf{r}) > (t)$ (where t is the time at which the average is constructed). Since thermal collisions are random, so that displacements in each direction are equally likely, we have $< \mathbf{r} > (t) = 0$. The most simple thermal average that is commonly used to characterize the Brownian motion of the colloid, is the mean-squared displacement W(t), defined as,

$$W(t) \equiv \langle r^2 \rangle(t) ,$$

although the average $\langle r \rangle (t)$ would also be a candidate to characterize Brownian motion. The thermal average is equivalent with a so-called ensemble average, where each realization of **r** is weighted by the probability that that realization occurs. Let $P(\mathbf{r}, t)$ be the probability



Fig. 1: The mean-squared displacement as a function of time. For short times, displacements of the colloid are limited to within a "cage" of neighbours (as indicated by the cartoon), while at long times the colloid moves from one cage to the other. Since inter-colloid interactions diminish displacements, we have $D_0 > D_s^s > D_s^l$.

density function (pdf) for r at time t. The thermal average of a function $f(\mathbf{r})$ can thus also be written as,

$$\langle f(\mathbf{r}) > (t) = \int d\mathbf{r} f(\mathbf{r}) P(\mathbf{r}, t) ,$$

where the volume integral ranges over all values for r. The mean-squared displacement (for which $f(\mathbf{r}) \equiv r^2$) can thus be calculated as a function of time, once the pdf $P(\mathbf{r}, t)$ is known. The equation of motion for this pdf, for very dilute solutions of rigid spherical macromolecules, reads [1, 2],

$$\frac{\partial P(\mathbf{r},t)}{\partial t} = D_0 \nabla^2 P(\mathbf{r},t) , \qquad (1)$$

where,

$$D_0 = \frac{k_B T}{\gamma}, \qquad (2)$$

is Einstein's diffusion coefficient, with $\gamma = 6\pi\eta_0 a$ the Stokes friction coefficient (k_B is Boltzmann's constant, T the temperature, η_0 the shear viscosity of the solvent, and a the radius of the colloids). It thus follows that,

$$\frac{dW(t)}{dt} = D_0 \int d\mathbf{r} \ r^2 \, \nabla^2 P(\mathbf{r}, t) = D_0 \int d\mathbf{r} \ P(\mathbf{r}, t) \, \nabla^2 r^2 = 6 \, D_0$$

where in the second equation two partial integrations have been performed, while in the third equation it is used that $\nabla^2 r^2 = 6$ and that the integral of P is unity. Since W(t = 0) = 0, it follows that,

$$W(t) = 6 D_0 t . (3)$$

The mean-squared displacement thus varies linearly with time, which is due to the many random interactions with surrounding solvent molecules. Such a linear time dependence is, however, no longer valid for concentrated dispersion, where the colloids interact with each other, or when the colloid moves through a crowded environment. For such cases, and for sufficiently small times, the colloid will "rattle in a cage" formed by the neighbouring colloids (or, for example, move within a mesh of the network). For longer times, the mean-squared displacement is determined by the motion from one cage to the other, as depicted in Fig.1. For times shorter than the escape time of the colloid from its cage, there are random displacements of the colloid, just like in an unbounded fluid, so that for these times the mean-squared displacement is still given by eq.(3) but with a different diffusion coefficient (see Fig.1). The diffusion coefficient is affected through hydrodynamic interactions of the colloid with the cage, and is referred to as the *short-time self diffusion coefficient*, denoted by D_s^s (the subscript stands for "self", the superscript for "short"). The nomenclature "self diffusion" relates to the motion of a single colloid, in contrast to concerted motion of an assembly of colloids (which will be discussed below). The equation of motion (1) now reads,

$$\frac{\partial P(\mathbf{r},t)}{\partial t} = D_s^s \nabla^2 P(\mathbf{r},t) , \text{ for } t < \tau_c , \qquad (4)$$

where τ_c denotes the cage escape time. As before this equation of motion leads to $W(t) = 6 D_s^s t$. For long times, longer than the escape time, there is a random displacement from cage to cage, so that eq.(3) again applies, but with yet another diffusion coefficient which is referred to as the *long-time self diffusion coefficient*, denoted by D_s^l (the superscript stands for "long"). For these long times we have,

$$\frac{\partial P(\mathbf{r},t)}{\partial t} = D_s^l \nabla^2 P(\mathbf{r},t) , \text{ for } t > \tau_c , \qquad (5)$$

and hence $W(t) = 6 D_s^l t$ (see Fig.1). The long-time self diffusion coefficient is given by a similar Einstein equation (2) as for the diffusion at very low concentrations, except that the friction coefficient is now affected by interactions of the colloid under consideration with its neighbours (or with the network through which it moves).

Brownian motion is also at the origin of migration of colloids due to their spatially varying concentration. Consider an arbitrary initial concentration profile, but with very small spatial gradients in the concentration $\rho(\mathbf{r}, t)$. The flux of colloids for such small gradients will be proportional to the $\nabla \rho(\mathbf{r}, t)$. For steeper spatial variations, higher order gradients also come into play. The most general form for the flux **j** of colloids that is linear in spatial gradients of the concentration reads,

$$\mathbf{j}(\mathbf{r},t) = -\int d\mathbf{r}' \int_{-\infty}^{t} dt' D_c(\mathbf{r}-\mathbf{r}',t-t') \nabla' \rho(\mathbf{r}',t') , \qquad (6)$$

where D_c is referred as to as the *collective diffusion kernel* (a minus sign is introduced to render D_c positive). The position dependence of the diffusion kernel D_c is due to interactions between the colloids: the flux at a point **r** is affected by concentration gradients at neighbouring points **r'** as a result of interactions between the colloids. The diffusion kernel is thus only non-zero for distances $|\mathbf{r} - \mathbf{r'}|$ smaller than the range of interactions between the colloids, which is of the order of the size of the colloids. The time dependence of the diffusion kernel accounts for possible memory effects, where the flux at a certain instant of time may depend on earlier states

of the system: microstructural order may not adapt instantaneously to the changing density for diffusion processes that are faster than the typical time on which microstructural relaxation occurs. For smooth concentration variations, however, the diffusion process is slow, so that one may expect that memory effects are of no importance. The collective diffusion kernel can thus be approximated as arbitrarily sharply peaked function around $\mathbf{r} = \mathbf{r}'$ and t = t', so that eq.(6) reduces to,

$$\mathbf{j}(\mathbf{r},t) = -D_c \nabla \rho(\mathbf{r},t') ,$$

with,

$$D_c = \int d\mathbf{r} \int_{-\infty}^0 dt \ D_c(\mathbf{r}, t) ,$$

where D_c is the *collective diffusion coefficient*, also referred to as *Fick's diffusion coefficient* (after the German physiologist Fick, who introduced the diffusion equation in 1855). It thus follows that,

$$\frac{\partial \rho(\mathbf{r},t)}{\partial t} = -\nabla \cdot \mathbf{j} = \nabla \cdot [D_c \nabla \rho(\mathbf{r},t)]$$

where the first equation is the *continuity equation*, which is an exact equation that expresses conservation of the number of colloids. Notice that the diffusion coefficient is inside the divergence, since it is generally concentration dependent as a result of interactions between the colloids. For small deviations from the mean concentration $\bar{\rho}$, however, the diffusion coefficient can be taken outside the divergence (since $\nabla \cdot [D_c \nabla \rho] = D_c \nabla^2 \rho + (dD_c/d\rho) |\nabla \rho|^2$, and the latter term is of second order in deviations of the density from its average value $\bar{\rho}$). For such small deviations of the actual concentration from the average concentration, we thus finally find that,

$$\frac{\partial \rho(\mathbf{r},t)}{\partial t} = D_c \nabla^2 \rho(\mathbf{r},t) .$$
(7)

which is usually referred to as Fick' diffusion equation.

The *diffusion equations* (4,5,7) will be essential for the interpretation of some of the FCS and FRAP experiments that will be discussed later.

1.2 What is fluorescence?

Molecules can emit light after excitation of electrons to higher energy levels by means of absorption of electromagnetic radiation or by mechanical means. When the excitation results from absorption of electromagnetic radiation, the subsequent emission of light is referred to as *luminescence*. There are two types of luminescence processes, depending on the electronic transitions: *fluorescence* and *phosphorescence*. Both luminescence processes will be discussed in the present section on a descriptive level. In case of fluorescence, the molecule is referred to as a *fluorophore* or simply a *dye*.

Luminescence occurs in several distinct steps, and along various pathways:

Excitation:

Excitation of an electron only occurs when the energy $E_e = h \nu_e$ of an incident photon is equal to the energy difference between the quantum states of the electron after and before excitation



Fig. 2: A schematic representation of electronic- and superimposed vibrational energy levels (a so-called Jablonksi diagram), with an indication of the several transitions that may occur, leading to either fluorescence or phosphorescence.

(here, h is Planck's constant and ν_e is the frequency of the incident light leading to excitation). Each quantum state comprises an electronic state (such as s-, σ -, and π -orbitals) and a vibrational state related to the relative motion of nuclei. The energy differences between vibrational states are much smaller than those for electronic states. This can be schematically depicted by a horizontal thick line, representing the energy of the electronic state, and a number of super-imposed thinner lines, representing the additional energy of different vibrational states (see Fig.2). Such a representation is commonly referred to as a *Jablonski diagram*. The electronic state with the lowest energy is the ground state, denoted by S(0), while the subsequent states with higher energies are denoted by $S(1), S(2), \cdots$. The energy diagram in Fig.2 is an idealized representation, since the quantum-states and the corresponding energies occupied by the luminescent molecule vary with time due to interactions with its surrounding solvent molecules. This leads to a continuum of frequencies of the incident light that can be absorbed and emitted. The time an excitation requires is of the order of a $10^{-15} s$, which is the time for the incident light to traverse its own wavelength. A possible excitation is indicated in Fig.2 by the left light grey arrow.

Vibrational relaxation:

The most likely process immediately after absorption is a transition to a lower-level vibrational state, as indicated by the top-left open arrow in Fig.2. Vibrational relaxation takes about $10^{-13} - 10^{-12} s$, which is much shorter than the typical life time of an exited electronic state, that is at least of the order of 10^{-9} .

Fluorescence:

After fast vibrational relaxation, several further relaxation mechanisms back to the electronic ground state S(0) may occur. One of the mechanisms is the transition from the exited electronic state directly to the electronic ground state (as depicted by the black-white dotted arrow in Fig.2). The process where light is emitted as a result of this transition is referred to as fluorescence. Such an electronic relaxation takes about $10^{-9} - 10^{-7} s$. Right after this electronic transition, vibrational relaxation occurs within the ground state (see the open arrow just below the black-white dotted arrow in Fig.2). Due to relatively fast vibrational relaxation right after

excitation, the energy difference E_f between the quantum state before and after electronic relaxation is smaller than the energy E_e of absorption. Since $E_f = h \nu_f$ (with ν_f the frequency of the fluorescent light), while the frequency is trivially related to the wavelength λ of light as $\nu = c/\lambda$, where the speed of light c is similar for all frequencies, it follows from $E_e > E_f$ that the wavelength λ_e of the absorbed light is smaller than the wavelength λ_f of the emitted light. This phenomenon is referred to as the *Stokes shift* (after the British physicist Stokes who first described fluorescence in 1852). The Stokes shift is of immense experimental importance, since it allows to detect fluorescent light intensities without the interference of incident light, by the use of, for example, band-filters or dichroic mirrors.

Phosphorescence:

A different pathway for electronic relaxation to the ground state involves an intermediate quantum state, referred to as a *triplet state* (the state T(1) in Fig.2). In the S(n)-states, the excited electron has an opposite spin orientation to the electron that remains in the S(0)-state, just like in the ground state before excitation (see the boxed scheme in the right-lower corner of Fig.2). The energy of such states with opposite spins corresponds to a single wave function, and are therefore referred to as *singlet states* (hence the abbreviation S(n)). The transition from the S(1)-state to the triple T(1)-state (indicated by the dashed arrow in Fig.2) occurs when the exited electron changes its spin orientation (see the right-lower box in Fig.2). There are now three different wave functions that have the same energy; hence the name "triplet state". The transition from a singlet state to a triplet state is referred to as *inter-system crossing*, and takes typically $10^{-10} - 10^{-8} s$. The decay from the triplet state to the singlet ground state can occur, like for fluorescence, under the emission of light or through the generation of heat. The emission of light due to this transition if referred to as *phosphorescence*. The T(1)-to-S(0) transition is improbable, and therefore takes typically a relatively quite long time of $10^{-3} - 10^{+2} s$.

Delayed fluorescence:

Since the triplet state is long-lived, there can be a reversed inter-system crossing $T(1) \rightarrow S(1)$ (*not* indicated in Fig.2), by adsorption of energy from interactions with surrounding molecules. After such a transition, the excited singlet state can return to the ground state, just like for fluorescence. This process is referred to as *delayed fluorescence*. The time it takes for this type of fluorescence is of the order of the life-time of the triplet state, which can be in the time window of the macromolecular diffusive processes of interest, and may therefore interfere with measurements that probe such processes. In the section on FCS we will show how these contributions can be dealt with.

When the illuminating intensity is sufficiently large, multiple excitations of an electron to quite high energy states can occur. Such multiple-excited dyes are susceptible to react with surrounding molecules, thereby loosing their fluorescent properties. Such a process is referred to as **photo-bleaching**. As will be discussed later, FRAP takes advantage of photo-bleaching.

2 Fluorescence Correlation Spectroscopy (FCS)

FCS can be used to measure diffusion coefficients of macromolecules as well as binding constants of associating/dissociating macromolecules. In the following subsections, the principles of FCS are discussed, and experimental examples are presented concerning diffusion of colloids and proteins through isotropic and nematic networks of rods, as well as the measurement of protein-membrane binding constants.

2.1 The principles of FCS

In a FCS experiment, the time dependent fluorescence intensity is measured, resulting from fluorescent macromolecules that diffuse in-and-out a small illuminated volume. The faster the intensity fluctuates, the faster the Brownian motion of the macromolecules, implying a larger diffusion coefficient. A small illuminated volume is created by focusing a laser beam, where the small illuminated volume is referred to as the *confocal volume* (see Fig.3a). The spatial distribution of the intensity of a focused laser beam, with its propagation direction in the *z*-direction, is given by (see the dotted lines in Fig.3a),

$$I(\mathbf{r}) = I_0 \exp\left\{-\frac{x^2 + y^2}{\sigma_1^2} - \frac{z^2}{\sigma_2^2}\right\} , \qquad (8)$$

where I_0 is the overall intensity of the laser beam, and with σ_1 and σ_2 the extent of the confocal volume in the x, y- and z-direction, respectively. Typically, $\sigma_1 \approx 300 nm$ and $\sigma_2 \approx 1500 nm$. The fluorescence intensity originating from dyes residing at a position **r** within the confocal volume is both proportional to the local concentration $\rho_d(\mathbf{r}, t)$ of dyes and the local incident intensity $I(\mathbf{r})$, so that the fluorescence intensity $I_f(t)$ is proportional to,

$$I_f(t) \sim \int d\mathbf{r} \ I(\mathbf{r}) \ \rho_d(\mathbf{r}, t) ,$$
(9)

The time dependence of the fluorescence intensity is entirely due to Brownian motion of the colloids in and out of the confocal volume. What is neglected in eq.(9) is the dependence of excitation probabilities and fluorescent intensities on the orientation of dye molecules. This is probably a good approximation when there are many dye molecules with random orientations attached to the colloid.

An implicit assumption in eq.(9) is that emitted fluorescent intensities originating from different dye molecules add up, that is, interference effects from light emitted by different dyes is neglected. The rational behind this assumption is that the life time of an exited state is much larger than the time it takes a photon to traverse over a distance of the order of a wavelength, and that life times have a very broad distribution due to strong interactions of dye molecules with their environment. There is thus no phase relation between photons emitted by different dye molecules, so that interference effects play no role, or in other words, fluorescent light is *incoherent*.

The fluctuations of the fluorescence intensity are quantified by the so-called *fluorescence intensity correlation function* $C_f(t)$, which is defined as,

$$C_f(t) \equiv \langle I_f(t) I_f(t=0) \rangle$$
, (10)

where the brackets $\langle \cdots \rangle$ refer to ensemble averaging (as discussed in section 1.1). This correlation function measures the time span during which intensities are still correlated. A random Brownian displacement of a colloid has no memory of earlier displacements after some time (say, τ seconds), during which several displacements occurred. The fluorescence intensity originating from that colloid at time t' is therefore no longer correlated to fluorescent intensities from the same colloid at times earlier than $t' - \tau$. For times $t > \tau$, the correlation function becomes therefore equal to,



Fig. 3: (a) A sketch of a focused laser beam, giving rise to a small volume of high intensity (the confocal volume). The laser beam propagates from top to bottom. The dashed elliptically shaped lines indicate equi-intensity surfaces with equal differences in intensity, and thus depicts the confocal volume. The black dots indicate the fluorescent macromolecules which move diffusively through the confocal volume. (b) Macromolecules diffusing through a nematic network of long and thin rod-like colloids, where the rods are on average oriented in a direction perpendicular to the propagation direction of the laser beam.

The correlation function is a monotonically decreasing function of time, with a final value equal to that given in eq.(11). The decay rate will be seen to be related to the diffusion coefficient of the colloids. Which diffusion coefficient is measured (the short-time or long-time self diffusion coefficient, or the collective diffusion coefficient; see section 1.1) will be discussed in the subsequent subsection.

From eqs.(9,10) it is found that,

$$C_f(t) \sim \int d\mathbf{r} \int d\mathbf{r}' I(\mathbf{r}) I(\mathbf{r}') < \rho_d(\mathbf{r}, t) \rho_d(\mathbf{r}', t = 0) > .$$
(12)

In a further evaluation of the integrals, it is convenient to work with Fourier transforms. The Fourier transform $f(\mathbf{k})$ of a function $f(\mathbf{r})$ is defined as (where the dependence on the conjugate Fourier variable k indicates that this is the Fourier transform),

$$f(\mathbf{k}) \equiv \int d\mathbf{r} f(\mathbf{r}) \exp\{i\,\mathbf{k}\cdot\mathbf{r}\}, \qquad (13)$$

where the Fourier variable \mathbf{k} is commonly referred to as the *wave vector*. With the help of Parseval's theorem (which states that the integral of a product of two functions of \mathbf{r} and \mathbf{r}' with respect these variables is equal to the wave vector integral of the product of their Fourier transforms), together with the representation,

$$\delta(\mathbf{k} + \mathbf{k}'') = \int d\mathbf{r} \exp\{i(\mathbf{k} + \mathbf{k}'') \cdot \mathbf{r}\},\$$

of the delta function $\delta(\cdot)$, eq.(12) can be rewritten as,

$$C_f(t) \sim \int d\mathbf{k} \int d\mathbf{k}' I(\mathbf{k}) I(\mathbf{k}') < \rho_d(\mathbf{k}, t) \rho_d(\mathbf{k}', t = 0) > , \qquad (14)$$

where the Fourier transform of the intensity in eq.(8) is equal to (with $\mathbf{k} = (k_x, k_y, k_z)$),

$$I(\mathbf{k}) = I_0 \exp\left\{-\frac{1}{2}\left[(k_x^2 + k_y^2)\sigma_1^2 + k_z^2\sigma_2^2\right]\right\}.$$
(15)

Since the dye molecules are attached to the macromolecules, the Fourier transform of the dye concentration ρ_d , according to the definition (13) of the Fourier transform, is equal to,

$$\rho_d(\mathbf{k},t) = B_f(\mathbf{k}) \sum_{j=1}^N \exp\{i\mathbf{k} \cdot \mathbf{r}_j(t)\}, \text{ where } B_f(\mathbf{k}) = \int d\mathbf{r}' \,\chi_d(\mathbf{r}') \,\exp\{i\mathbf{k} \cdot \mathbf{r}'\},$$

where it is assumed that all macromolecules are identically labeled with dyes, according to the characteristic function χ_d for the distribution of dye molecules on the surface and/or within the core of the macromolecules (the characteristic function is defined as $\chi_d(\mathbf{r}') = 1$ when there is a dye molecule at \mathbf{r}' , and = 0 otherwise). Furthermore, $\mathbf{r}_j(t)$ is the position coordinate (the position of the center-of-mass) of the j^{th} macromolecule. Note that when the distribution of dyes is non-spherically symmetric, the fluorescence amplitude B_f will be time dependent through the rotational Brownian motion of the macromolecule. Here we will assume a spherically symmetric dye distribution within the macromolecule. A possible asymmetry of the dye distribution, however, does not affect the correlation function, provided that the macromolecule is much smaller that the linear dimensions of the confocal volume. The assumption of a symmetric dye distribution is therefore not as restrictive as it may seem. For translationally invariant systems, where the macroscopic concentration is constant throughout the confocal volume, one can show that,

$$<\rho_d(\mathbf{k},t)\rho_d(\mathbf{k}',t=0)>\sim \delta(\mathbf{k}+\mathbf{k}')\sum_{i,j=1}^N<\exp\{i\mathbf{k}\cdot(\mathbf{r}_i(t)-\mathbf{r}_j(t=0))\}>$$

It thus follows from eq.(14) that,

$$C_f(t) \sim \int d\mathbf{k} \ I^2(\mathbf{k}) \ |B_f(\mathbf{k})|^2 \ \sum_{i,j=1}^N < \exp\{i\mathbf{k} \cdot (\mathbf{r}_i(t) - \mathbf{r}_j(t=0))\} > .$$
 (16)

This is the fundamental expression for any FCS experiment, which for each different type of sample must be explicitly evaluated in order to interpret FCS data quantitatively.

When the macromolecules are very small compared to the confocal volume, the fluorescence amplitude B_f is essentially a constant equal to $B(\mathbf{k} = \mathbf{0})$. We will restrict ourselves here to the motion of identical macromolecules that do not mutually interact, but may be moving through a crowded environment. In such cases we have,

$$\sum_{i,j=1}^{N} < \exp\{i\mathbf{k} \cdot (\mathbf{r}_{i}(t) - \mathbf{r}_{j}(t=0))\} > = 0 , \ i \neq j , \text{ for mutually non-interacting colloids }.$$

For such small, mutually non-interacting, and identical macromolecules (possibly in a non-fluorescent crowded environment), we thus find from eq.(16) that,

$$C_f(t) \sim \int d\mathbf{k} \ I^2(\mathbf{k}) < \exp\{i\mathbf{k} \cdot (\mathbf{r}(t) - \mathbf{r}(t=0))\} > ,$$
(17)

where $\mathbf{r}(t)$ is the position of a macromolecule. The correlation function appearing in the integral can be obtained from the diffusion equations (4,5) as follows. First note that,

$$<\exp\{i\mathbf{k}\cdot(\mathbf{r}(t)-\mathbf{r}(t=0))\}>=\int d\mathbf{r}\int d\mathbf{r}_0 \ P(\mathbf{r},\mathbf{r}_0,t)\,\exp\{i\mathbf{k}\cdot(\mathbf{r}-\mathbf{r}_0)\}\,,\tag{18}$$

where $P(\mathbf{r}, \mathbf{r}_0, t)$ is the probability density function (pdf) for the occurrence of \mathbf{r} at time t and \mathbf{r}_0 at time t = 0. Next, the conditional pdf $P(\mathbf{r} | \mathbf{r}_0, t)$ is introduced, which is the pdf for \mathbf{r} at time t, given that the position coordinate is equal to \mathbf{r}_0 at time t = 0. A little thought shows that,

$$P(\mathbf{r},\mathbf{r}_0,t) = P(\mathbf{r}|\mathbf{r}_0,t) \times P_0(\mathbf{r}_0) ,$$

where $P_0(\mathbf{r}_0)$ is the pdf for \mathbf{r}_0 . For a sample that does not macroscopically change with time, and is translation invariant (which was already assumed before), the latter pdf is a constant, independent of \mathbf{r}_0 . The conditional pdf is the solution of the equations of motion (4) or (5) (depending whether the short- or long-time diffusive properties are probed for the particular sample under consideration), but with the initial condition that $P(\mathbf{r} | \mathbf{r}_0, t = 0) = \delta(\mathbf{r} - \mathbf{r}_0)$, with $\delta(\cdot)$ the delta function, as before. Since the conditional pdf is a function of $\mathbf{R} = \mathbf{r} - \mathbf{r}_0$, that is, $P(\mathbf{r} | \mathbf{r}_0, t = 0) \equiv P(\mathbf{R}, t)$, it thus follows from eq.(18) that,

$$\langle \exp\{i\mathbf{k}\cdot(\mathbf{r}(t)-\mathbf{r}(t=0))\}\rangle = \int d\mathbf{R} P(\mathbf{R},t) \exp\{i\mathbf{k}\cdot\mathbf{R}\}.$$

The correlation function is therefore the Fourier transform $P(\mathbf{k}, t)$ of $P(\mathbf{R}, t)$, which is immediately found from the Fourier transformed equations of motion eq.(4) or eq.(5). Since under Fourier transformation $\nabla^2 \rightarrow -k^2$, these equations of motion transform to,

$$\frac{P(\mathbf{k},t)}{\partial t} = -D_s^{s,l} k^2 P(\mathbf{k},t) .$$

д

The superscript is either "s" or "l" in case of short-time or long-time diffusion, respectively. Hence,

$$<\exp\{i\mathbf{k}\cdot(\mathbf{r}(t)-\mathbf{r}(t=0))\}> = P(\mathbf{k},t) = \exp\{-D_s^{s,l}k^2t\}.$$
 (19)

The explicit time dependence of the fluorescence correlation function finally follows from eqs.(15,17,19) and the evaluation of the wave-vector integral,

$$C_f(t) \sim (1 - F + F \exp\{-\tau/\tau_{trip}\}) \times \left[1 + \frac{D_s^{s,l} t}{\sigma_1^2}\right]^{-1} \left[1 + \frac{D_s^{s,l} t}{\sigma_2^2}\right]^{-1/2} .$$
(20)

Here, the prefactor accounts in an empirical fashion for the fraction F of dyes that emit light after being for a typical time τ_{trip} in a triplet state giving rise to delayed fluorescence (see subsection 1.2). This is the standard equation used in fitting FCS correlation functions to obtain the diffusion coefficient.

The above standard expression for the fluorescence correlation function assumes that the diffusion coefficient is isotropic, that is, the diffusion coefficient is identical in all three directions. This is different when the macromolecule diffuses through a crowded environment with an anisotropic structure. An example of this will be discussed in the next subsection, where diffusion of macromolecules through a transient nematic network of long and thin rods is discussed. In such a nematic network, the rods have a preferential orientation along the so-called nematic director. This situation is sketched in Fig.3b. The diffusion coefficient perpendicular to the nematic director is different from that along the nematic director. Instead of eqs.(4,5) for the pdf of the position coordinate of a macromolecule, we now have,

$$\frac{\partial P(\mathbf{r},t)}{\partial t} = D_{s,\parallel}^{s,l} \frac{\partial^2 P(\mathbf{r},t)}{\partial x^2} + D_{s,\perp}^{s,l} \begin{bmatrix} \partial^2 P(\mathbf{r},t) \\ \partial y^2 + \frac{\partial^2 P(\mathbf{r},t)}{\partial z^2} \end{bmatrix},$$

where the nematic director is taken along the x-direction, and where $D_{s,\parallel}^{s,l}$ and $D_{s,\perp}^{s,l}$ ate the diffusion coefficients parallel and perpendicular to the nematic director, respectively. Repeating the above analysis leads to (apart from the prefactor that accounts for the occurrence of triplet states),

$$C_f(t) \sim \left[1 + \frac{D_{s,\parallel}^{s,l} t}{\sigma_1^2}\right]^{-1/2} \left[1 + \frac{D_{s,\perp}^{s,l} t}{\sigma_1^2}\right]^{-1/2} \left[1 + \frac{D_{s,\perp}^{s,l} t}{\sigma_2^2}\right]^{-1/2} .$$
(21)

Apart from the diffusion coefficients of interest, the unknown quantities in eqs.(20,21) are the linear dimensions $\sigma_{1,2}$ of the confocal volume. These are determined from an FCS experiment with a dilute suspension of colloids for which the diffusion coefficient is known from, for example, dynamic light scattering.

The prefactor in eqs.(20,21) depends on the average number of macromolecules residing within the confocal volume. Fluctuations of the fluorescence intensity are small when there are many macromolecules within the confocal volume, giving rise to very low experimental values of the correlation function. When the concentration of macromolecules is so low that there are large time spans where there are no macromolecules in the confocal volume, the experimental errors in the correlation function are also quite large. The ideal number of fluorescent macromolecules in the confocal volume to be able to measure accurate correlation functions is about 2-5.

Note that the expression in eq.(16) for the FCS correlation function is a k-weighted average of the correlation function that is measured in a dynamic light scattering experiment. This expression shows that for large concentrations, where the fluorescently labeled macromolecules interact with each other, FCS measures the collective diffusion coefficient rather than the self diffusion coefficient. It is sometimes erroneously stated that FCS always probes self diffusion due to the incoherence of the fluorescence light.

2.2 Diffusion of spherical colloids and proteins through isotropic and nematic networks of rods

As an example of FCS experiments, we consider here the self diffusion of spherical colloids and proteins through networks of very long and thin rods. The rods are fd-viruses, which consist of a 880 nm long dsDNA strand which is covered by 2700 proteins. The thickness of a fd-virus particle is 6.8 nm. The coat proteins render the fd-virus particles quite stiff (the persistence length is about 2500 nm), so that these particles are often used model systems for very long, thin, and stiff colloidal rods. At low fd-rod concentrations, the rods are randomly oriented. The corresponding phase is referred to as the *isotropic phase*. On increasing the concentration of fd-rods, interactions between the rods give rise to a phase transition where the rods tend to align in a preferred direction. This phase is the so-called *nematic phase*, and the preferred direction is the *nematic director* (see a schematic picture of the nematic phase in Fig.3b, where the rods



Fig. 4: (a) The self diffusion coefficient of a spherical silica colloid with a radius of 210 nm in an isotropic network consisting of fd-virus rods, as a function of the rod-concentration. The diffusion coefficients are normalized to the diffusion coefficient D_0 of a very dilute solution of the colloid in water. The middle curve is obtained from FCS, while the short-and long-time self diffusion coefficients (the upper and lower curves, respectively) have been obtained from dynamic light scattering (DLS). (b) The mean squared displacement W(t) of the colloid obtained from light scattering experiments as a function of time, for various fd-rod concentrations, as indicated in the figure in mg/ml. Data are taken from Ref.[3].

are depicted as thick lines). The diffusion of spherical colloids and proteins in such a nematic network is anisotropic, that is, the diffusion coefficients are different for diffusion along the nematic director and perpendicular to it.

Figure 4a shows self diffusion coefficients of a spherical silica colloidal particle with a radius of 210 nm as a function of the concentration of the fd-rods, within the isotropic phase [3]. The upper and lower curves in this figure correspond to the short- and long-time diffusion coefficients, respectively, as measured with dynamic light scattering (DLS). As can be seen, FCS measures a cross-over from short-time to long-time diffusion up to fd-concentrations of about 3 mq/ml. FCS probes long-time diffusion only for fd-concentrations larger than 3 mq/ml. The reason for this is as follows. When the colloid traverses the confocal volume in typical times during which the mean-squared displacement (MSD) is characterized by the short-time diffusion coefficient (for times below the cage escape time in Fig.1), FCS probes the corresponding short-time diffusion coefficient. If the residence time of the colloid in the confocal volume is much larger than the cage escape time, FCS probes predominantly long-time diffusion. The MSD as measured by dynamic light scattering is given in Fig.4b, for various rod-concentrations in mq/ml indicated in the figure [3]. These experimental MSDs comply with the sketched MSD time dependence in Fig.1. For a MSD equal to about $\sigma_1^2 \approx 0.1 \, \mu m^2$ (as before, σ_1 is the smallest linear dimension of the confocal volume), it can indeed be seen from these plots that long-time diffusion dominates for $W(t) \le 0.1 \,\mu m^2$ for concentrations larger than about $2.5 - 3 \,mq/ml$. For fd-concentrations below 1 mq/ml, the MSD for $W(t) \le 0.1 \mu m^2$ is within the short-time regime, so that FCS probes short-time diffusion. For intermediate concentrations, FCS probes a mix of short- and long-time diffusion. A comparison of FCS and FRAP measurements (as well as simulations) in Ref.[4] also shows that one should be careful with the interpretation of experimental results



Fig. 5: (a) A FCS correlation function that probes the motion of the protein apoferritin in a nematic network of fd-virus rods (of which the concentration is 55 mg/ml). The lower curve marked "background", is the correlation function from the network, without the protein being present. (b) The long-time diffusion coefficients of the protein as a function of the rod-concentration, within the isotropic phase and the nematic phase. The vertical bar indicates the two-phase isotropic/nematic coexistence region. In the nematic phase, the upper and lower curves correspond to diffusion parallel and perpendicular to the nematic director, respectively. Data are taken from Ref.[5].

in terms of short- and long-time diffusion (or a mix of the two), depending on the size of the illuminated region relative to the size of the macromolecule as well as the matrix within which the macromolecules are embedded.

Figure 5a shows a typical example of a FCS correlation function of a protein (apoferritin), embedded in a network of fd-rods in the nematic state (see Fig.3b) [5]. The solid curve drawn through the data points is a fit to eq.(21). The lower curve marked "background" is the contribution to the measured correlation function resulting from residual fluorescence of the fd-rods and possible contaminations. Since many molecules fluoresce to some extent, it is always important to correct for such background contributions. Figure 5b shows the diffusion coefficients extracted from such FCS measurements [5]. Within the isotropic phase (for fd-concentrations less than about 20 mq/ml) there is just a single diffusion coefficient, where the FCS correlation function can be fitted to eq.(20) (the triplet correction is not necessary here, since the decay time due to delayed fluorescence is much smaller than $10 \,\mu s$). For the relatively small protein, it is verified that the long-time self diffusion coefficient is probed with FCS. As expected, the diffusion coefficient decreases with increasing concentration in the isotropic phase, since Brownian displacements of the protein are hindered by the increasing concentration of obstacles from the rod network. In the nematic phase of the rod network, there are two diffusion coefficients, as discussed before: a diffusion coefficient that relates to the MSD for motion parallel to the nematic director (the preferred direction of alignment of the rods), an a diffusion coefficient related to motion perpendicular to the director. These two diffusion coefficients are shown in Fig.5b, where it is seen that motion parallel to the director is less hindered by the network as compared to perpendicular motion. The physics underlying diffusion of spherical Brownian particles in rod networks has been discussed in Refs. [3, 5, 6].
2.3 Complex formation and protein-membrane adsorption

FCS can also be used to probe binding constants of two complex-forming macromolecules. In such an experiment, one of the species is fluorescently labeled. Since the self-diffusion coefficient of the free non-associated macromolecule and the association complex are different, the FCS correlation function is a superposition of two modes,

$$C_f(t) \sim A_0 F_0(t) + Q A_c F_c(t) ,$$
 (22)

where F_0 is the FCS correlation function of the free, unbounded macromolecule, and F_c of the associated complex. The amplitudes A_0 and A_c measure the relative number of non-associated and associated macromolecules. The fraction f_c of associated macromolecules is given by,

$$f_c = \frac{A_c}{A_0 + A_c}$$

Furthermore, the "quantum-yield ratio" Q in eq.(22) is the ratio of the fluorescence quantum yield (the probability that a dye returns to its ground state through fluorescence) of the complex and the free macromolecule. Both correlation functions F_0 and F_c are given by eq.(20), with the respective diffusion coefficients for the free macromolecule and the association complex. These diffusion coefficients are long-time self diffusion coefficients, averaged over the orientation of the macromolecule and the complex, provided that both are sufficiently small as compared to the linear dimensions of the confocal volume, and that fluorescent species do not mutually interact (except for the interaction leading to complex formation). In principle, binding constants can be obtained in crowded environments, also in concentrated solutions of both species by labeling only a small fraction of one of the species.

A fit of the FCS correlation function to two modes is only possible when the two diffusion coefficients of the free macromolecule and the complex are sufficiently different. As shown in Ref.[7], the size difference between the free macromolecule and the complex to perform quantitative measurements must be larger than about 2.

A convenient method to measure the complexation of soluble proteins with membrane proteins is to use nanodisks [9]. Nanodisks are disk-like pieces of phospholipid bilayers which are stabilized by scaffold proteins that wrap around the circumference of the disk (a schematic is given in Fig.6a). Nanodisks form spontaneously by first mixing the scaffold proteins and detergent-solubilized lipids and subsequently slow removal of the detergent (for example by dialysis). Membrane proteins can be embedded during the formation of the disks, provided that suitable detergent-solubilization of these proteins can be achieved, and the corresponding detergents are also removed during dialysis. The diameter of nanodisks is quite narrow distributed, and can be varied from about 10 to 20 nm, depending on the scaffold protein and the mixing ratio of the proteins and lipids. The thickness of nanodisks is close to 5 nm. Since the disks are quite monodisperse in size, and they are quite a bit larger than common proteins, they are ideal mimetic systems to probe association/dissociation reactions of soluble proteins and membrane proteins with FCS.

As an example of such experiments, Fig.6b shows FCS correlation functions, where the fluorescently labeled protein IAPP adsorps to nanodisks of mixed lipid membranes (without any incorporated membrane proteins) [9]. IAPP (islet amyloid polypeptide) is a peptide hormone secreted by the pancreas, which readily forms amyloid fibers, and thereby contributes to diabetes. Due to the large size difference between the protein and the nanodisk, there is almost a decade difference in the decay time of the correlation functions for the unbound and bound



Fig. 6: (a) A sketch of the structure of a nanodisk (taken from [8]. The green chains represent adsorpted molecules through the yellow-indiated head group. (b) Normalized FCS correlation functions for free diffusing IAPP (lower curve), labeled nanodisks (upper curve), and a mixture (middle curve). (b) The fraction of bound IAPP as a function of the concentration of nanodisks. The solid line is a fit to a Langmuir adsorption isotherm. Data are taken from Ref.[9].

protein. The fraction of bound protein as a function of the total nanodisk concentration, as obtained from fits of FCS correlation functions as discussed above, is shown in Fig.6b [9]. Here, the quantum-yield ratio Q is found from intensity measurements with the same incident intensity from samples containing only protein, and where essentially all proteins are bound at high concentrations of nanodisks. The solid line is a fit to a Langmuir adsorption isotherm,

$$f_c = \frac{c_{nd}}{K_D + c_{nd}} \,,$$

where c_{nd} is the total nanodisk concentration and K_D is the so-called dissociation constant. The value for K_D from the fit in Fig.6c is found to be equal to 50 nM [9]. A Langmuir adsorption isotherm assumes that adsorbed proteins do not interact with each other, so that they for example do not form oligomers, that the proteins adsorb within a monolayer, and that adsorption does not change the membrane structure that affects adsorption of other proteins. It is thus very well possible that Langmuir-adsorption fails for higher concentrations of adsorbed proteins. Oligomer formation of adsorpted proteins could in this way be probed by FCS.

3 Fluorescence Recovery After Photo-bleaching (FRAP)

As discussed at the end of subsection 1.2, photo-bleaching is the destruction of the fluorescence of a dye molecule with a very high-intensity light beam. FRAP uses the photo-bleaching process in order to probe self diffusion coefficients through the diffusion driven recovery of the fluorescence intensity right after the application of a high-intensity light pulse that gives rise to the photo-bleaching process.

We limit the discussion here to the determination of translation diffusion coefficients, although FRAP can also be used to probe rotational diffusion, taking advantage of the dependence of the absorption and fluorescence quantum efficiency of dye molecules depending on their orientation relative to the polarization direction of the incident and detected light [10, 11].



Fig. 7: (a) The concentration variation of unbleached macromolecules before (a constant concentration) and right after the bleach pulse (sinusoidaly varying concentration). A low concentration occurs at positions where the bleach-pulse intensity is high. (b) The concentration of bleached macromolecules after the bleach pulse. (c) The laser beam with low intensity is oscillatory shifted over the originally bleached pattern. A high intensity is measured when the maximima of the laser intensity coincide with the maximum in the concentration of unbleached macromolecules, and the other way around. (d) A schematic of the experimental set up, where a mirror is mounted on a piezo element which oscillatory changes the position of the mirror in order to oscillate the laser interference pattern as depicted in (c). The Pockels cell is used to switch between the large bleaching intensity and the relatively low reading intensity. This schematic is taken from Ref.[12].

After discussing the principles of FRAP, two experimental realizations of FRAP experiments are presented, together with experimental examples.

3.1 The principles of FRAP

After bleaching part of the sample, the dynamics of recovery of fluorescence due to intermixing of bleached and intact macromolecules is probed. Provided that bleaching of dyes does not change the interactions between the macromolecules, the time dependence of the recovering intensity after photo-bleaching is governed by self diffusion (see section 1.1). The quantitative interpretation of FRAP data relies on eq.(9) for the fluorescence intensity as well as the diffusion equations (4,5) which determine the time dependence of the concentration of unbleached macromolecules.

There are typically two ways used to create a bleached volume within the sample. For local measurements (like in part of a living cell), a high-intensity bleach pulse of a focused laser beam is applied to the sample. The bleaching now occurs in a region with an extent that can be as small as that of the confocal volume discussed before in subsection 2.1 (see also Fig.3a). The same laser beam is used at a much lower intensity (the "reading intensity"), avoiding further bleaching, to monitor the recovery of the fluorescence as unbleached macromolecules diffuse

into the bleached region. For samples which are homogeneous over larger length scales (of about 1 mm or more), the signal-to-noise ratio of a FRAP measurement can be significantly improved as follows. Two laser beams are now crossed under a small angle, creating a standing, sinusoidaly varying interference pattern of alternating high and low intensity. Applying a short pulse of a high intensity now creates a sinusoidaly varying concentration variation of bleached and unbleached macromolecules (see Figs.7a,b). After the bleach pulse, the same interference pattern with a much smaller intensity is used to probe the recovery of the fluorescent intensity. The amount of unbleached macromolecules, after the bleach pulse is applied, gives rise to a background fluorescence intensity, on top of which the time dependent recovering intensity is measured. To eliminate the background intensity, and thus maximizing the signal-to-noise ratio, the interference pattern can be spatially shifted in an oscillatory fashion over the bleached pattern, between full overlap with the original bleached pattern and a shift of half a wavelength (as depicted in Fig.7c). In the former position, a small intensity is measured (since the large bleaching intensity coincides with the region where relatively many of the macromolecules are bleached), while in the latter the largest intensity is measured (large bleaching intensity coincides with the region of high concentration of unbleached macromolecules). This is depicted in Fig.7c. The timely oscillating fluorescence intensity with the given frequency of the oscillating intensity pattern is filtered-out using a lock-in amplifier. The lock-in amplifier signal resulting from the spatial variation of the concentration of bleached and unbleached macromolecules is now a "null measurement", that is, the measured FRAP signal is now zero when there are no concentration gradients. A schematic of the experimental set up is given in Fig.7d. A mirror is mounted on a piezo element, which sinusoidaly shifts the mirror up-and-down, thus oscillatory shifting the location of the interference pattern of the two crossed laser beams in the sample as shown in Fig.7c.

In the next two subsections, examples of both types of experiments are discussed.

3.2 Measuring second moments of the fluorescence intensity: a geometry-independent FRAP analysis

In case a small region is bleached, a quantitative analysis of the recovery dynamics requires the solution of the diffusion equation (4) in case of short-time diffusion and eq.(5) for longtime diffusion. As mentioned above, the size of the bleached region is usually much larger than the size of the macromolecules, so that long-time diffusion is commonly probed (see, for example, Ref.[4]). The solution of the diffusion equation (5) is non-trivial, and depends on the precise geometry of the initially bleached region. For a square bleached region, for example, the solution of the diffusion equation involves error functions [13]. There is, however, a method to obtain diffusion coefficients that does not require prior knowledge of the geometry of the bleached region [14, 15] (the analysis given here differs in details from those in the original publication in Ref.[14], with a slightly different result as presented in Refs.[14, 15]). Since the probability density function $P(\mathbf{r}, t)$ to find an unbleached macromolecule at position \mathbf{r} at time tis directly proportional to the concentration $\rho^{(-)}(\mathbf{r}, t)$ of unbleached macromolecules, the latter satisfies the same diffusion equation (5). We write the solution of the diffusion equation in terms of the so-called Green's function $G(\mathbf{r}, t)$ as,

$$\rho^{(-)}(\mathbf{r},t) = \int d\mathbf{r}' \ G(\mathbf{r} - \mathbf{r}',t) \ \rho^{(-)}(\mathbf{r}',t=0) \ .$$
(23)

The Green's function is defined to be the solution of the diffusion equation with the initial condition,

$$G(\mathbf{r} - \mathbf{r}', t = 0) = \delta(\mathbf{r} - \mathbf{r}'), \qquad (24)$$

where $\delta(\cdot)$ is, as before, the delta function. Consider an experiment where diffusion in 2D is probed, by either using a sample cell with a height (in the z-direction, say) that is sufficiently small so that bleaching is uniform along the z-direction, or by bleaching a sufficiently deep pattern such that diffusion in the z-direction does not contribute to the recovery of the fluorescence intensity. In that case, the relevant Green's function is the solution of the diffusion equation in 2D (where the z-dependence is lost), which reads [1, 2],

$$G(x, y, t) = \frac{1}{4 \pi D_s^l t} \exp\left\{-\frac{x^2 + y^2}{4 D_s^l t}\right\}.$$
 (25)

The quantity that is of interest is the second moment of the local fluorescence intensity $I_f(x, y)$,

$$I_2(t) \equiv \int dx \int dy \, (x^2 + y^2) \, I_f(x, y) = \tilde{C} \int dx \int dy \, (x^2 + y^2) \, I(x, y) \, \rho^{(-)}(x, y, t) \,, \quad (26)$$

where the constant \tilde{C} is proportional the number of dyes per macromolecule and the fluorescence quantum efficiency of the dyes. The illumination intensity I(x, y) is chosen such that it is a constant throughout the bleached region: $I(x, y) \equiv I$. We now define $\Delta \rho^{(-)}(x, y, t) = \rho^{(-)}(x, y, t) - \rho$, with ρ the uniform concentration of macromolecules before the bleach pulse is applied. Since $\Delta \rho^{(-)}(x, y, t)$ is only non-zero within the bleached region, this allows us to rewrite eq.(26) as,

$$I_2(t) - I_2^{(0)} = C \int dx \int dy \, (x^2 + y^2) \,\Delta\rho^{(-)}(x, y, t) , \qquad (27)$$

with $I_2^{(0)}$ the second moment of the intensity before the bleach pulse (corresponding to the concentration ρ), and where $C = \tilde{C}I$. Since $\Delta \rho^{(-)}(x, y, t)$ satisfies the diffusion equation, we can substitute eqs.(23,24,25) with $\rho^{(-)}$ replaced by $\Delta \rho^{(-)}$ into eq.(26). Changing to the integration variables X = x - x' and Y = y - y' leads to,

$$I_{2}(t) - I_{2}^{(0)} = \frac{C}{4\pi D_{s}^{l} t} \int dX \int dY \int dx' \int dy' \left([X + x']^{2} + [Y + y']^{2} \right) \\ \times \exp\left\{ -\frac{X^{2} + Y^{2}}{4 D_{s}^{l} t} \right\} \Delta \rho^{(-)}(x', y', t = 0) .$$
(28)

Assuming a symmetric bleaching intensity, so that $\Delta \rho^{(-)}(x', y', t) = \Delta \rho^{(-)}(-x', y', t)$ and $\Delta \rho^{(-)}(x', y', t) = \Delta \rho^{(-)}(x', -y', t)$, the resulting Gaussian integrals can be evaluated,

$$C \\ 4\pi D_s^l t \int dX \int dY \int dx' \int dy' \left(X^2 + Y^2\right) \exp\left\{-\frac{X^2 + Y^2}{4 D_s^l t}\right\} \Delta \rho^{(-)}(x', y', t = 0) \\ = 4 D_s^l t C \int dx' \int dy' \Delta \rho^{(-)}(x', y', t = 0) = 4 D_s^l t \left[I_f(t = 0) - I_f^{(0)}\right],$$

where $I_f(t = 0)$ is the total fluorescence intensity right after the bleach pulse, and $I_f^{(0)}$ the intensity just before the bleach pulse is applied, and,

$$C \int dx \int dy \int dx' \int dy' (x'^2 + y'^2) \exp\left\{-\frac{X^2 + Y^2}{4 D_s^l t}\right\} \Delta \rho^{(-)}(x', y', t = 0)$$

= $C \int dx' \int dy' (x'^2 + y'^2) \Delta \rho^{(-)}(x', y', t = 0) = I_2(t = 0) - I_2^{(0)}.$

We thus finally find that the so-called normalized second intensity moment $\mu_2(t)$ is equal to,

$$\mu_2(t) \equiv \frac{I_2(t) - I_2(t=0)}{I_f(t=0) - I_f^{(0)}} = 4 D_s^l t .$$
⁽²⁹⁾

As before, "t = 0" refers to the time right after the bleach pulse, and $I_f^{(0)}$ is the fluorescence intensity just before the bleach pulse. This result is independent of the geometry of the bleached region, and can thus be employed to determine the diffusion coefficient for those experiments where there is an uncertainty concerning the bleached geometry. There are a few differences between the results stated in Refs.[14, 15]. In these references, it is stated that $\mu_2(t) = [I_2^{(0)} - I_2(t)]/[I_f^{(0)} - I_f(t)] = 4D_s^l t$ (in the absence of immobile macromolecules), which differs from the expression in eq.(29). In addition, no mention is made in these references that the reading intensity should be uniform throughout the bleached region

The result in eq.(29) predicts that I_2 diverges at long times, which indicates that this relation is only valid for sufficiently small times. The reason for this is as follows. We stated just above eq.(27) that " $\Delta \rho^{(-)}(x, y, t) = \rho^{(-)}(x, y, t) - \rho$ is only non-zero within the bleached region." This allowed us to assume a constant incident intensity in order to be able to explicitly evaluate the integrals. During homogenization, however, bleached macromolecules diffuse far outside the bleached region (see, for example, Fig.2 in Ref.[15]), while the final uniform concentration of unbleached macromolecules is somewhat smaller than before the bleach pulse was applied (depending on the volume of the bleached region relative to the sample volume as a whole). Taking I(x, y) constant in eq.(27) thus implicitly assumes this to be also the case far outside the bleached region for longer times, which invalidates the above analysis for long times. The weight factor distance-squared in the definition of I_2 then leads to the unphysical divergence of I_2 at long times. The conclusion is therefore that eq.(29) can only be applied for times where $4D_s^l t$ is less than approximately the squared linear dimension of the bleached region. Such times are sufficiently long to perform accurate measurements of the diffusion coefficient.

Diffusion of proteins like BSA in solution and network-forming bio-polymers (resembling the polymers from which the bone of the larynx is constructed) have been probed in Ref.[15] using the above described FRAP method. The second moment for a FRAP measurement with the protein BSA, labeled with the dye fluoresceine isothiocyanate (FITC), is given in Fig.8a as a function of time. As can be seen, the second moment indeed varies linearly with time, and the expected deviation from the linear time dependence is observed at long times, as discussed above. Note, however, that the intercept at zero time is non-zero, which is not due to the presence of immobile macromolecules (a full recovery of the fluorescence intensity is reported). The non-zero intercept may be connected to the different expression for the normalized second moment used in Ref.[15], as mentioned above. A second example is given in Fig.8b from Ref.[14]. Here, the diffusion of lipids within the membrane of a living cell has been probed by labeling a fraction of the lipids that constitute the membrane with the dye NBD-HPC. The filled



Fig. 8: (a) and (b) The second intensity moment $\mu_2(t)$ for the protein BSA in a buffer solution (taken from Ref.[15]) and for lipids in the membrane of a living cell (taken from Ref.[14]), respectively. (c) The (normalized) fluorescence intensity as a function of time after a bleach pulse that completely bleaches a green fluorescent protein within part of a living cell (taken from Ref.[17]).

circles and the fitted solid line refer to measurements of the upper part of the membrane, and the filled squares and dotted line to the lower part. To within the rather large errors, the expected linear dependence is indeed observed. The respective diffusion coefficients are found to equal to $0.15 \pm 0.06 \,\mu m^2$, and $0.25 \pm 0.06 \,\mu m^2$. In the same Ref.[14], the diffusion coefficients have been obtained by the standard recovery analysis, which requires the solution of the diffusion equation involving Bessel functions [16]. The corresponding diffusion coefficients are $0.28 \pm 0.03 \,\mu m^2$ and $0.24 \pm 0.02 \,\mu m^2$. In addition, the diffusion coefficients are determined from the time dependence of the total fluorescence intensity, again both for the upper and lower part of the membrane, which are found to be equal to $0.19 \,\mu m^2$ and $0.23 \,\mu m^2$ (no errors are given). The various results for the diffusion coefficients seem to be in the same range. However, an expression that is used to obtain the second moment from experiments (see their eq.(13)) is the same as used in Ref.[15], but differs from our expression in eq.(29).

In the above analysis it is assumed that all macromolecules are mobile. For those cases where a significant fraction of the macromolecules are immobile, the Green's function in eq.(25) should be multiplied by the fraction k of mobile macromolecules, and a time-independent contribution must added to account for the immobile macromolecules. Repeating the above analysis then leads to,

$$\mu_2(t) \equiv \frac{I_2(t) - I_2(t=0)}{I_f(t=0) - I_f^{(0)}} = 4 k D_s^l t + (1-k) \mu_0 , \qquad (30)$$

where μ_0 is a time independent constant. It is therefore sometimes important (also for other FRAP experiments) to determine the fraction of immobile macromolecules. The fraction k of immobile macromolecules can be obtained by a measurement of the final recovery of the fluorescence intensity after a bleach pulse that bleaches all the macromolecules, so that right after the pulse the fluorescence intensity is zero. The recovery is incomplete due to bleached immobile macromolecules which remain at their fixed position. Figure 8c shows the result of such an experiment, where a green fluorescence protein is bleached in part of a living cell [17]. After long times, where the fluorescence intensity attains a time independent value, recovery is seen to be incomplete. The fraction k of immobile proteins is found from a plot of the





Fig. 9: (a) FRAP-signals for a double-glass $G_L(G_s)$, where both the large (L) and small (s) spheres are kinetically arrested, as evidenced by the not-fully decaying signals at large times. (b) A single-glass, where only the large spheres are in the glass state. (c) The phase diagram. The vertically axis represents the volume fraction φ_s of the small spheres, the horizontal axis of the larger spheres φ_L . The data are taken from Refs.[19, 20].

fluorescence intensity normalized to the intensity before the bleach pulse is applied, as shown in Fig.8c.

Protein transport in cells is not always purely diffusive. There may be temporary association to the cell network, transport of proteins attached to molecular motors may occur, and proteins may be expressed during a FRAP experiment. These different contributions to the FRAP signal can some times be separated (see, for example, Ref.[18]).

3.3 FRAP experiments using an oscillatory reading intensity: the glass transition in binary colloids

The bleached pattern in the FRAP experiment discussed in subsection 3.1 (see in particular Fig.7) is a sinusoidaly varying function of position $\sim \sin\{\mathbf{k} \cdot \mathbf{r}\}$, where \mathbf{k} is the so-called wave vector. This is a vector in the direction in which the intensity varies with position, with a magnitude equal to $2\pi/\Lambda$ with Λ the wavelength of the sinusoidaly varying interference pattern, also referred to as the *fringe spacing*. The fringe spacing can be calculated once the angle between the two crossed laser beams has been determined. The FRAP-signal analysis is now straight forward, since $\nabla^2 \sin\{\mathbf{k} \cdot \mathbf{r}\} = -k^2 \sin\{\mathbf{k} \cdot \mathbf{r}\}$, so that the concentration of the unbleached macromolecules remains sinusoidal for all times, and the solution of the diffusion equation (5) is a single exponential function, that is $\rho^{(-)}(\mathbf{r}, t) \sim \exp\{-D_s^l k^2 t\}$. A simple exponential fit to the FRAP signal therefore immediately leads to a value for the long-time self diffusion coefficient.

As an example of such a FRAP experiment we discuss here the diffusion in a glass-forming binary mixture of spherical colloids [19, 20]. Glasses are dynamically arrested non-equilibrium states, where a particle is kinetically trapped in cages that are formed by their neighbours. Strong repulsive interactions between the particles renders the system essentially frozen in a

non-equilibrium state, which only very slowly evolves toward equilibrium (referred to as "ageing"). The size ratio of the two species of colloidal spheres considered here is sufficiently large (1 to 9), such that the small colloidal spheres can move through the spaces in between tightly packed arrangements of the larger spheres. In the glass state, the spheres are confined to move within their cages, and are therefore immobile: they can only move over distances of the order of the linear dimension of their cage. The FRAP signal will therefore not fully decay to zero at long times, contrary to spheres in the liquid state, while the initial decay is related to the motion of the spheres within their cages. The dynamics of each of the two species can be probed independently by fluorescently labeling just one of the species, embedded in a matrix of unlabeled spheres of the other species. Figure 9a shows FRAP-signals as a function of time, where both species of spheres are in the glass state, the $G_L(G_s)$ -state, as evidenced by the non-fully decaying FRAP signals for both species. In Fig.9b the large spheres are in the glass state while the small spheres are still in the fluid phase. In this $G_L(F_s)$ -state, the small spheres can move diffusively to any position, and are not confined by a cage of neighbouring particles. These glass states are referred to as a double-glass (both species glassified) and a single-glass (only one species glassified). The transition between both glassy states is indicated by the wiggly line in the upper right corner of the phase diagram in Fig.9c, where the vertical and horizontal axes indicate the volume fraction of the small and large spheres, respectively. In addition to the glass states at relatively high concentration of both species, there is fluid F-phase where both species are in the fluid state, there is a coexistence F + C-phase where the large spheres crystallized, and a metastable M-state where the stable state is a F + C-state, but where the system remains fluid like during the duration of the experiment.

4 Förster Resonance Energy Transfer (FRET)

FRET probes the distance between two dyes, in the range of typically 1 - 20 nm. The two dyes can either be attached to different parts of a single macromolecule or to two associating/dissociating macromolecules. FRET can thus be used, for example, to study the distribution of distances between two domains within a single protein with a flexible connection between them, as well as the kinetics of protein-ligand binding.

In the subsections below we first discuss the principles of FRET, then describe in some detail one possible realization of a FRET experiment, and finally present an example where the unfolding of a protein is studied with FRET.

4.1 The principles of FRET

One of the dyes, the "donor", is excited by external illumination. The second dye, the "acceptor", is then excited by *radiationless* transfer of the donor excitation energy to the acceptor. The non-radiative transfer of energy from the donor to the acceptor proceeds through dipoledipole interaction, where the dipoles have a similar resonance frequency. The acceptor then falls back into its ground state by fluorescence. The fluorescence of the acceptor after radiationless transfer of the excitation energy from the donor is referred to a *FRET-event*. The corresponding fluorescence intensity depends on the probability of radiationless transfer of energy from the donor to the acceptor, which is commonly referred to as the *FRET-efficiency*. The FRET-efficiency in turn sensitively depends on the instantaneous distance between the two dyes. The connection between the fluorescence intensity of the acceptor after radiationless energy transfer and the distance between the donor and acceptor underlies the principle of FRET experiments to determine intra- and inter-molecular distances.

The quantum mechanical analysis of radiationless energy transfer (which is beyond the scope of this chapter) reveals that the FRET-efficiency E(r) depends on the instantaneous distance r between the donor and acceptor as (see the original papers by Förster in Refs.[21, 22], as well as Ref.[23]),

$$E(r) = \frac{R_0^6}{R_0^6 + r^6}, \qquad (31)$$

where $R_0 \approx 1 - 20 nm$ is the so-called *Förster radius*, which depends on the fluorescence quantum yield of the donor (in the absence of the acceptor), the refractive index of the medium, the relative orientation of the dipoles of the donor and acceptor, and the degree of spectral overlap of the donor-acceptor pair, that is, the spectral overlap of the fluorescence spectrum of the donor and the absorption spectrum of the acceptor, as depicted by the grey area in Fig.10a. Note that the FRET-efficiency is a quite sensitive function the donor-acceptor distance, varying like $\sim r^{-6}$, so that accurate measurements can be performed on the distribution of donoracceptor distances.

In the design of a FRET experiment, the first step is to chose an appropriate donor-acceptor pair. The absorption and fluorescent spectra of the two dyes should fulfill certain conditions, depending on the type of FRET experiment that is planned. Here, one should keep in mind that the absorption and fluorescence spectrum can significantly differ for free dyes in solution and dyes attached to the macromolecule. The chemical attachement of the donor and acceptor to appropriate positions on the macromolecules is highly non-trivial. There are nowadays several procedures to specifically label e.g. proteins (both co- and post-translational), the discussion of which is beyond the scope of this chapter (see Ref.[24] and references therein).

The type of FRET experiment of course depends on the particular phenomenon and system that should be addressed. First of all, there is a distinction between FRET experiments where an entire assembly of many macromolecules is probed simultaneously, and so-called single-molecule FRET experiments (smFRET) where fluorescence intensities are measured that originate from a single macromolecule. In this chapter we focus on smFRET measurements, which can for example be realized using a confocal microscope for the measurements of intra-macromolecular distances, such that at most one macromolecule resides within the confocal volume, and for association/dissociation reactions using evanescent-wave illumination, with one of the reacting molecules attached to a substrate.

Consider a smFRET experiment on a protein with two domains, where one domain is labeled with the donor and the other with the acceptor. The distance between these domains will fluctuate due to thermal motion. Repeated measurements of the FRET-efficiency E allows for the construction of a histogram, where the number of times that a certain value for E is measured is plotted as a function of E. Such a *FRET-histogram* is a representation of the probability that a certain value for E is found, and according to eq.(31) also represents the probability distribution for the distances between the two dyes. In case there are two low-energy states between which the protein switches, the FRET-histogram is typically a superposition of two probability distributions as sketched in Fig.10b. It should be mentioned, that in many cases (especially in cell-biology applications) the signal-to-noise ratio is quite small, so that the construction of FRET-histograms as depicted in Fig.10b is not feasible. In such cases one can only distinguish between high or low FRET signals, which binary information of course limits the conclusions that can be drawn from FRET experiments.



Fig. 10: (a) The absorption spectra (dashed curves) and fluorescence spectra (solid curves) of the donor and the acceptor. Radiationless transfer of the excitation energy of the donor to the acceptor requires that the fluorescence spectrum of the donor overlaps with the absorption spectrum of the acceptor (indicated by the grey area). Intensities are recorded simultaneously at two wavelengths λ_A and λ_D . The incident light only excites the donor (so that $I_A^{fluor}(\lambda_{A,D}) = 0$), while λ_A is red-shifted with respect to the fluorescence spectrum of the donor (so that $I_D^{fluor}(\lambda_A) = 0$). (b) A sketch of a typical FRET-histogram for a protein with the donor and acceptor attached to two different domains. Here, the number of events with which a certain value for E is measured is plotted against E.

There are several ways to experimentally probe FRET-efficiencies, examples of which are:

- Simultaneous measurement of the acceptor and donor fluorescence intensity at different wellchosen wavelengths (which will be discussed in more detail in subsection 4.2),

- Probing the lifetime of, or fluorescence intensity from, the donor in the presence and absence of the acceptor. The corresponding difference in the fluorescence life time and intensity measures the probability for a FRET-event.

- Measuring donor life times or intensities before and after bleaching the acceptor. Before bleaching both direct fluorescence and fluorescence followed by a FRET-event occurs, while after bleaching no FRET-events take place anymore. The donor intensity, for example, is therefore enhanced after bleaching to an extent that is determined by the FRET-efficiency.

- Similarly to the previous method, bleaching the donor eliminates FRET-events and thus diminishes the acceptor fluorescence intensity depending on the FRET-efficiency.

Clearly, a detailed discussion of all these different types of FRET experiments is beyond the scope of this chapter (for more detailed overviews, see Refs.[25, 26, 27]). Instead, the principles of one particular type of smFRET experiment is discussed in the subsequent subsection 4.2, and an experimental example is presented concerning the unfolding of a protein in subsection 4.3.

4.2 FRET-efficiencies obtained from the simultaneous measurement of donor and acceptor fluorescence intensities

Here we discuss one method to measure smFRET-histograms in some detail (the derivation given below is adapted from Refs.[26, 27]). A laser is used with a wavelength within the ab-

sorption spectrum of the donor (see Fig.10a). A time series of intensities is measured for two wavelengths λ_A and λ_D by two detectors, within the range of the fluorescence spectrum of the acceptor and that of the donor, respectively (as depicted in Fig.10a). Only those intensities are accounted for in the data analysis where the fluorescence intensity at λ_A exceeds a preset threshold value, in combination with a relatively low fluorescence intensity at λ_D , in order to assure that a considerable fraction of the intensity at λ_A is connected to FRET-events. The experimentally measured intensity $I^{exp}(\lambda_A)$ at wavelength λ_A has four different contributions: the intensity $I^{FRET}_A(\lambda_A)$ originating from FRET-events, $I^{fluor}_A(\lambda_A)$ and $I^{fluor}_D(\lambda_A)$ due to fluorescence after direct excitation of the acceptor and the donor by the laser light, respectively, and a background intensity $I^{back}(\lambda_A)$ that is due to light originating from single donors or acceptors (either macromolecules with only a donor or acceptor bounded, or freely diffusion unbounded dyes in solution), or any other spurious light of wavelength λ_A that is detected (like light of unwanted wavelength that leaks through beam splitters and dichroic mirrors). Hence,

$$\tilde{I}^{exp}(\lambda_A) = R(\lambda_A) \left[I_A^{FRET}(\lambda_A) + I_A^{fluor}(\lambda_A) + I_D^{fluor}(\lambda_A) \right],$$
(32)

where $\tilde{I}^{exp}(\lambda_A) = I^{exp}(\lambda_A) - I^{back}(\lambda_A)$ is the intensity corrected for the background intensity. Furthermore, $R(\lambda_A)$ is the response factor of the experimental set up, that is, the detector unit read-out produced by a single photon. Similarly, the experimentally measured intensity at the wavelength λ_D is equal to,

$$\tilde{I}^{exp}(\lambda_D) = R(\lambda_D) \left[I_D^{fluor}(\lambda_D) + I_A^{fluor}(\lambda_D) \right],$$
(33)

where $R(\lambda_D) \neq R(\lambda_A)$ is the response factor of the set up at wavelength λ_D . Note that the intensities on the right-hand side in eqs.(32,33) are in units of photons per second, while the experimental intensities on the left-hand side are in detector units (like the voltage of the detectors, which defines the units in which the response functions $R(\lambda_D)$ and $R(\lambda_A)$ are expressed). The numerical values of the response functions for the two detectors at the two wavelengths are to be independently determined.

The wavelength λ_A is chosen above the fluorescence spectrum of the donor, and the excitation wavelength is chosen such that the acceptor is not excited (as sketched in Fig.10a). This implies that,

$$I_D^{fluor}(\lambda_A) = 0$$
 , $I_A^{fluor}(\lambda_D) = 0$, and $I_A^{fluor}(\lambda_A) = 0$. (34)

The connection between experimentally determined intensities and FRET-efficiencies is established as follows. First consider the intensity at λ_A . Let σ_D denote the absorption cross section of the donor, which has the dimension of $length^2$, and is proportional to the probability that a photon that interacts with the donor leads to excitation. The incident photon flux I^{ext} from the external illuminating source is defined as the number of photons that pass a unit area per unit time. From these definitions it follows that $\sigma_D I^{ext}$ is the number of photons per unit time that interact with the donor and subsequently lead to excitation. The probability that a subsequent FRET-event occurs is by definition equal to E. The subsequent probability that after energy transfer the acceptor emits a photon with wavelength λ_A due to fluorescence is denoted by Φ_A , the so-called fluorescence quantum yield. The intensity at λ_A due to FRET-events is thus equal to,

$$I^{FRET}(\lambda_A) = \sigma_D I^{ext} E \Phi_A ,$$

and hence, using eqs.(32,34),

$$\tilde{I}^{exp}(\lambda_A) = R(\lambda_A) \,\sigma_D \, I^{ext} \, E \, \Phi_A \,. \tag{35}$$

Next consider the intensity at λ_D . When contributions from triplet states can be neglected, so that 1 - E is the probability that an excited donor relaxes through fluorescence, analogous considerations for the intensity arising from fluorescence of the donor at wavelength λ_D gives,

$$I_D^{fluor}(\lambda_D) = \sigma_D I^{ext} (1-E) \Phi_D ,$$

where Φ_D is the fluorescence quantum yield of the donor at wavelength λ_D . Hence, from eqs.(33,34),

$$\tilde{I}^{exp}(\lambda_D) = R(\lambda_D) \,\sigma_D \,I^{ext} \,(1-E) \,\Phi_D \,. \tag{36}$$

Taking the ratio of eqs.(35,36) and solving for the FRET-efficiency finally leads to,

$$E = \frac{\tilde{I}^{exp}(\lambda_A)}{\tilde{I}^{exp}(\lambda_A) + \Lambda \tilde{I}^{exp}(\lambda_D)},$$
(37)

where,

$$\Lambda = \frac{R(\lambda_A) \Phi_A}{R(\lambda_D) \Phi_D}$$
(38)

It is important to note that the intensities in eq.(37) are in units of the detector read-out. Expressing the intensities in eq.(37) in units of photons per second by dividing by the corresponding response functions $R(\lambda_A)$ or $R(\lambda_D)$, the same expression in eq.(37) holds, but with $\Lambda = \Phi_A/\Phi_D$. In summary, the assumptions made in the above derivation of eqs.(37,38) are:

- Direct excitation of the acceptor by the incident laser can be neglected.

- The fluorescence intensity at λ_A arising from the donor can be neglected.

- Triplet states are essentially absent.

- The excitation of the acceptor due to fluorescence photons from the donor can be neglected.

- The intensities at the two wavelengths λ_A and λ_D must be appropriately corrected for background contributions.

An issue in FRET-data analysis is the relative orientation of the donor and acceptor, which is one of the factors that sets the value of the Förster radius, and is commonly referred to as the *orientation factor*. When the relative orientation can take all values, and changes on a time scale that is short compared to the time required for FRET energy transfer, the orientation factor can be orientationally averaged in order to obtain the correct Förster radius. This usually requires long and flexible spacers between the macromolecule and the two dyes. The disadvantage of such long spacers is that they give rise to significant dye displacements which are not connected to changes of configurations of the macromolecule. A compromise should always be made between the uncertainty of the value of the orientation factor, and hence the Förster radius, and the uncertainty in the determination of donor-acceptor distance due to spacer-related displacements.

4.3 Unfolding of phosphoglycerate kinase (PGK) induced by guanidine hydrochloride (GndHCl)

As mentioned before, FRET can be used to study the distribution of the distance between domains within a single protein as well association/dissociation reactions. Here we will discuss an





Fig. 11: (*a*) The three dimensional structure of native PGK, and the positions where the donorand acceptor dyes are connected (Alexa 488 and Alexa 647, respectively). The Förster radius of this donor-acceptor pair is 4.9 nm. This structure is taken from Ref.[28], where the image is obtained by using PyMOL (W. L. DeLano, The PyMOL Molecular Graphics System, 2002, DeLano Scientific, San Carlos, CA, USA). (b) FRET-histograms, showing the number of acceptor fluorescence events versus the FRET efficiency E, for different GndHCl concentrations as indicated in the different histograms. The solid lines are Gaussian fits, using either a single or two Gaussians. The data are taken from Ref.[28].

example of the domain-distance distribution in a large two-domain protein, phosphoglycerate kinase (PGK) [28], that is expressed in many different cells. PGK catalyses the phosphorylation of a phosphoglycerate, where in addition ATP is reduced to ADP, which is one step in the carbohydrate-cascade degradation. The domain distance distribution is strongly affected by the presence of guanidine hydrochloride (GndHCl), which unfolds PGK in a reversible fashion. The study of the domain-distance distribution of PGK as a function of the GndHCl concentration may shed light on the role played by the presence of domains on the unfolding transitions of proteins.

Figure 11a shows the native (fully folded) structure of PGK, together with the positions where the donor and acceptor dyes are connected, and Fig.11b shows FRET-histograms for various GndHCL concentrations (indicated in the corresponding panels), as obtained from similar experiments as described in the previous subsection. As can be seen from the left upper panel in Fig.11b, in the native state (no GndHCl added), there is a single (Gaussian-like) peak, indicating that there is a single ground state. At a GndHCl concentration of 0.5 M, there is a second peak appearing at low values for E, corresponding to a state with a larger distance between the domains. At 0.7 M GndHCl, the native folded state and the unfolded states occur with approximately equal probability. The coexistence between these two states shows that unfolding is not an abrupt transition from one state to the other. On increasing the GndHCl concentration, the folded state becomes increasingly improbable, and disappears at about 1 M. Note, however, that the peak positions of the two states remain unchanged up to GndHCL concentrations of about 1 M (to within experimental error), implying that the average distance between the domains in the two states remains the same, independent of the GndHCl concentration. It thus seems that the protein structure in each of the two states remains unchanged. What does change on increasing the GndHCl concentration is the relative probability of occurrence of the two states. This behaviour is analogous to a macroscopic thermodynamic phase transition, where within the two-phase region the intrinsic properties of the two phases do not change, but only the relative volumes of the two phases change. Here, the "two phases" are the two states which the protein takes, and the "relative volumes" are the times that the protein spends in either state. Unpublished experimental results, where the domains are individually labeled, show that the domains are fully unfolded [29]. As can be seen from the two last panels in Fig.11b, for concentrations of GndHCL larger than 1 M, the peak position of the fully unfolded protein nevertheless shifts to lower values of E. This is attributed to the polymeric nature of the unfolded state, and is analogous to the globule-coil transition known for synthetic polymers.

Acknowledgement

I am grateful to Prof. J. Fitter for very useful discussions concerning the section on FRET.

References

- W.B. Russel, D.A. Saville, W.R. Schowalter, *Colloidal Dispersions*, Cambridge, University Press, 1989.
- [2] J.K.G. Dhont, An Introduction to Dynamics of Colloids, Elsevier, Amsterdam, 1996.
- [3] K. Kang, J. Gapinski, M.P. Lettinga, J. Buitenhuis, G. Meier, M. Ratajczyk, J.K.G. Dhont, J. Chem. Phys. 122, 044905 (2005).
- [4] C. Lellig, J. Wagner, R. Hempelmann, S. Keller, D. Lumma, W. Härtl, J. Chem. Phys. 121, 7022 (2004).
- [5] K. Kang, A. Wilk, J. Buitenhuis, A. Patkowski, J.K.G. Dhont, J. Chem. Phys. 124, 044907 (2006).
- [6] K. Kang, A. Wilk, A. Patkowski, J.K.G. Dhont, J. Chem. Phys. 126, 214501 (2007).
- [7] U. Meseth, T. Wohland, R. Rigler, H. Vogel, Biophys. J. 76, 1619 (1999).
- [8] L. Shi, K. Howan, A.-T. Shen, Y.J. Wang, J.E. Rothman, F. Pincet, Nature Protocols, 8, 935 (2013).
- [9] A. Nath, A.J. Trexler, P. Koo, A.D. Miranker, W.M. Atkins, E. Rhoades, in *Methods in enzymology*, Elsevier, Vol. 472, 89, Chapter 6 (2010).
- [10] M.P. Lettinga, G.H. Koenderink, B.W.M. Kuipers, E. Bessels, A.P. Philipse, J. Chem Phys. 120, 4517 (2004).
- [11] B.W.M. Kuipers, M.C.A. van de Ven, R.J. Baars, A.P. Philipse, J. Phys.: Condens. Matter 24, 245101 (2012).
- [12] A. Imhof, A. van Blaaderen, G. Maret, J. Mellema, J.K.G. Dhont, J. Chem. Phys. 100, 2170 (1994).
- [13] N.B. Simeonova, W.K. Kegel, Faraday Discuss. 123, 27 (2003).

- [14] H. Kubitscheck, P. Wedekind, R. Peters, Biophys. J. 67, 946 (1994).
- [15] P. Gribbon, T. E. Hardingham, Biophys. J. 75, 1032 (1998).
- [16] D.M. Soumpasis, Biophys. J. 41, 95 (1983).
- [17] J. Lippincott-Schwartz, N. Altan-Bonnet, G.H. Patterson, Nat. Cell. Biol. Sept. 2003, S7.
- [18] M. Fritzsche, G. Charas, Nature Protocols, Vol. 10. No.5, 660 (2015).
- [19] A. Imhof, J.K.G. Dhont, Phys. Rev. Lett. 75, 1662 (1995).
- [20] A. Imhof, J.K.G. Dhont, Coll. & Surf. A; Physicochem. Eng. Aspects 122, 53 (1997).
- [21] Th. Förster, Annalen der Physik 2, 55 (1948) (in German).
- [22] Th. Förster, Discuss. Faraday Soc. 27, 7 (1959).
- [23] R.W. van der Wiel, in FRET Förster resonance energy transfer: from theory to applications, Eds. I. Medintz, N. Hildebrandt, Wiley-VCH Verlag GmbH & Co. KGaA, Chapter 3 (2014).
- [24] M. Sadoine, M. Cerminara, N. Kempf, M. Gerrits, J. Fitter, A. Katranidis, Anal. Chem. 89, 11278 (2017).
- [25] N. Hildebrandt, in FRET Förster resonance energy transfer: from theory to applications, Eds. I. Medintz, N. Hildebrandt, Wiley-VCH Verlag GmbH & Co. KGaA, Chapter 5, (2014).
- [26] E. Sisamakis, A. Valeri, S. Kalinin, P.J. Rothwell, C.A.M. Seidel, in *Methods in enzymology*, Elsevier, Volume 475, Chapter 18 (2010).
- [27] T. Pons, in FRET Förster resonance energy transfer: from theory to applications, Eds. I. Medintz, N. Hildebrandt, Wiley-VCH Verlag GmbH & Co. KGaA, Chapter 8 (2014).
- [28] T. Rosenkranz, R. Schlesinger, M. Gabba, J. Fitter, ChemPhysChem 12, 704 (2011).
- [29] Personal communication with Prof. J. Fitter.

A 4 Scattering

R. Zorn

Jülich Centre for Neutron Science Forschungszentrum Jülich GmbH

Contents

1	Introduction	2
2	Scattering from static systems 2.1 Structure factor from density 2.2 Structure factor from pair correlation function	3 3 5
3	Scattering from dynamic systems	7
4	Scattering probes: photons, neutrons, electrons	12
5	Scattering methods for biological matter investigations5.1Small-angle scattering5.2Neutron spin echo spectroscopy	14 16 18

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Scattering methods play an important rôle in the investigation of biological systems. In contrast to real space (microscopy) methods, the result of a scattering experiment is related to the Fourier transform of the structure and dynamics of the system. This is why scattering methods are often said to work in 'reciprocal space'. This has certain advantages and disadvantages. While the main disadvantage – the absence of an immediately comprehensible three-dimensional 'picture' of the system – is immediately clear, the advantages need some explanation:

First of all, the direct-space observation of biological objects on a molecular level is often impossible due to technical reasons. Their size is often too small for light microscopy. Electron microscopy may be impossible in their native environment (solution, melt). Atomic force microscopy may be difficult because of their softness. In these situations scattering methods are the only way to obtain structural information¹. But even if microscopy methods are available scattering methods have the advantage of intrinsic averaging. Average properties (e.g. the average molecular weight of polymer molecules) can be obtained without doing a statistics of individual observations.

In the literature, a general overview of scattering methods focused on biological applications is rarely given. There is some literature specialised for soft matter, e.g. ref. 1 and (for polymers) ref. 2. Most textbooks focus on individual methods without specifying the application further: neutron scattering [3–7], light scattering [8–12], and x-ray scattering [13–15].

The central statement of this lecture is that scattering experiments indirectly measure correlation functions. The usual derivation of scattering laws is based on the fact that the scattering law (for neutrons, photons or any other radiation) is essentially the absolute square of the Fourier transform of a scattering density. In Fig. 1 this is shown as the left way from the density $\rho(\mathbf{r})$ to S(Q). The Wiener-Khintchine theorem

$$\left|\int_{-\infty}^{\infty} f(x) \exp(\mathrm{i}Xx) \mathrm{d}x\right|^2 = \int_{-\infty}^{\infty} \langle f(0)f(x)\rangle \exp(\mathrm{i}Xx) \mathrm{d}x \tag{1}$$

now states that the absolute square of a Fourier transform is the Fourier transform of the autocorrelation function. This opens another way (the right one in Fig. 1) to calculate the scattering law. Apart from elucidating the meaning of the scattering law in another way, this gives an alternative to calculate it even if the density itself is not known.

This lecture will in the first section treat the results from a static system where the scattering is completely elastic. In this situation the scattering will only contain information about the structure. Strictly speaking, this is a fictitious assumption because all materials show some dynamics (quantum-mechanically even at zero temperature). Nevertheless, the broad range of diffraction methods is covered with sufficient accuracy. The next part of the lecture will deal with inelastic scattering. In this experiment, scattering gives information about the structure via the momentum transfer and about the dynamics via the energy transfer. For inelastic scattering a Fourier transform in time has to be carried out in addition leading from the time correlation function to the scattering function. Finally, the differences between probe beams (light, x-rays, neutrons) will be explained and a selection of scattering instruments for biological investigations will be presented.

¹In this lecture "structure" will be understood as the exact positions of atoms in space. Of course, in a more loose sense also methods as NMR and even spectroscopic methods give information on the structure.



Fig. 1: The two ways to calculate the scattering law from the microscopic density, left: as the absolute-squared Fourier transform of the density, right: as the Fourier transform of the correlation function.

2 Scattering from static systems

In this section it will be assumed that the scattering system is *static*. It is either represented by fixed positions of point scatterers in space, \mathbf{r}_j , or a time-independent density, $\rho(\mathbf{r})$. The former case can be included in the latter by considering the microscopic density as a sum of delta functions:

$$\rho(\mathbf{r}) = \sum_{j=1}^{N} \delta(\mathbf{r} - \mathbf{r}_j).$$
⁽²⁾

From the fact that the scatterers are fixed follows that the scattering will be *elastic*, i.e. the energy of the scattered particles will not change due to the scattering process. This is clear from classical mechanics because a system which is static before and after the scattering process cannot exchange energy. The equivalent argument from the wave picture would be that upon scattering by fixed centres there is no Doppler shift of the frequency.

2.1 Structure factor from density

The result of an elastic scattering experiment is usually expressed in terms of the *differential* cross-section which is the probability density that a particle is scattered into a solid angle element $d\Omega$ normalised to the intensity of the incident beam:

$$\frac{\mathrm{d}\sigma}{\mathrm{d}\Omega} = \left\langle \left| \sum_{j=1}^{N} b_j \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}_j) \right|^2 \right\rangle \,. \tag{3}$$

 b_j is a measure of the 'scattering power' of the particle. From the dimensions of the other quantities it is obvious that it has the dimension [length]. Therefore, b_j is called the *scattering length*. Note that the scattering length is not necessarily positive. $b_j < 0$ just means that scattering leads to a reversal of the amplitude, in other words a phase shift π . The scattering length may even be complex. In that case, the imaginary part corresponds to absorption of the scattered particle by the scatterer.

It can be seen that expression (3) does not contain the scattering angle 2θ directly but a *scattering vector* **Q**. It is the vectorial difference of the wave vector **k'** after scattering and that before scattering, **k**. The wave vectors are defined by having the length $|\mathbf{k}| = k = 2\pi/\lambda$ and



Fig. 2: Definition of the scattering vector \mathbf{Q} in terms of the incident and final wave vectors \mathbf{k} and \mathbf{k}' . The black (isosceles) triangle corresponds to elastic scattering. The blue and red ones correspond to inelastic scattering with energy loss or gain of the scattered radiation, respectively.

the direction of the propagation of the wave. For elastic scattering k' = k, and the definition of **Q** is graphically demonstrated by the black (isosceles) triangle in Fig. 2 resulting in

$$Q = \frac{4\pi}{\lambda} \sin \theta \,. \tag{4}$$

From this equation one can see that scattering depends on a combination of the scattering angle and the wavelength of the scattered radiation. The same Q can be obtained by different combinations of 2θ and λ .

At this point it is necessary to explain the meaning of the average $\langle ... \rangle$ in (3) and justify it. Of course for a completely arrested system and completely coherent radiation, (3) would be valid without the average. Experimentally, this situation is only realised in laser light scattering from rigid objects. There, the experiments as well as the calculation do not yield a smooth function $d\sigma/d\Omega$ but an assembly of so-called speckles. For two reasons this situation is exceptional and the observed scattering is actually an average:

- If a dynamics exists, even if it is sufficiently slow not to cause a noticeable inelasticity, the particles will rearrange over the duration of the experiment. In this sense, (...) expresses a temporal average over the experimental time.
- 2. If the radiation used is not highly coherent, the sum over the amplitudes in (3) has to be restricted to the coherence volume which is usually much smaller than the sample volume. The results from the individual regions have to be added as intensities, i.e. after the absolute-square. This implies the same average but to be interpreted as a thermodynamic average over different realisations of the particle positions. In the case of ergodic systems both averages have the same result.

With the assumption that all scatterers are identical (for neutron scattering implying that they are the same isotope and have the same spin orientation) one can factor out the material-specific properties N and b:

$$\frac{\mathrm{d}\sigma}{\mathrm{d}\Omega} = |b|^2 NS(\mathbf{Q}) \,. \tag{5}$$

with the remaining term

$$S(\mathbf{Q}) = \frac{1}{N} \left\langle \left| \sum_{j=1}^{N} \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}_j) \right|^2 \right\rangle$$
(6)

which depends solely on the statistics of the positions of the scatterers. $S(\mathbf{Q})$ is called *structure factor*.

If the scattering is not effected by individual scatterers but by a field (e.g. the magnetic field for neutrons) or a distribution (e.g. the electron density for x-rays) one has to use a continuum description instead of (6):

$$S(\mathbf{Q}) = \frac{1}{N} \left\langle \left| \int_{V} \mathrm{d}^{3} r \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}) \rho(\mathbf{r}) \right|^{2} \right\rangle.$$
(7)

It is easy to verify that this expression corresponds to (6) with the definition (2) of the microscopic density inserted. Expression (7) is the absolute square of the Fourier transform of the density and thus represents the 'left way' in Fig. 1.

But even if the individual scatterers are point-like, the continuum description may be useful if their exact positions are not known but only their mesoscopic densities. This is often the case for biological systems. Then, expression (7) will be a good approximation as long as the length scale defined by Q is large compared to the distances between the scatterers, $Q \ll 2\pi/\text{distance}$ (e.g. for small-angle x-ray or -neutron scattering).

At that point a simple way to introduce mixed scatterers is to start with the *scattering length density*

$$\rho_b(\mathbf{r}) = \sum_{j=1}^N b_j \delta(\mathbf{r} - \mathbf{r}_j) \,. \tag{8}$$

instead of the density. By including the scattering properties in the density, equation (5) can be written as

$$\frac{\mathrm{d}\sigma}{\mathrm{d}\Omega} = \left\langle \left| \int_{V} \mathrm{d}^{3} r \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}) \rho_{b}(\mathbf{r}) \right|^{2} \right\rangle \,. \tag{9}$$

Thus, the differential cross section is the absolute square of the Fourier transform of the scattering length density. For neutron scattering, this concept is used to obtain a low-resolution description for small-angle scattering and reflectometry. For light scattering the local dielectric constant of the medium plays the rôle of $\rho_b(\mathbf{r})$.

2.2 Structure factor from pair correlation function

The second way to derive the scattering law starts with applying the definition of the absolute square, $|X|^2 = X^*X$ to equation (6):

$$S(\mathbf{Q}) = \frac{1}{N} \left\langle \left(\sum_{j=1}^{N} \exp(-\mathrm{i}\mathbf{Q} \cdot \mathbf{r}_{j}) \right) \left(\sum_{k=1}^{N} \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}_{k}) \right) \right\rangle$$
$$= \frac{1}{N} \sum_{j,k=1}^{N} \left\langle \exp(\mathrm{i}\mathbf{Q} \cdot (\mathbf{r}_{k} - \mathbf{r}_{j})) \right\rangle.$$
(10)

From this expression two characteristic properties of scattering become clear:

- 1. The scattering law arises from particle pairs (j, k).
- 2. Only distances between particles enter the expression, not the individual positions. The scattering law remains invariant under translation of the whole sample.

Analogously to $S(\mathbf{Q})$ the differential scattering cross section can be expressed in terms of particle distances as

$$\frac{\mathrm{d}\sigma}{\mathrm{d}\Omega} = \left\langle \sum_{j,k=1}^{N} b_j^* b_k \exp\left(\mathrm{i}\mathbf{Q}\cdot(\mathbf{r}_k - \mathbf{r}_j)\right) \right\rangle \,. \tag{11}$$

In order to proceed in a similar way as before, we introduce the two-particle density

$$\rho(\mathbf{r}_1)\rho(\mathbf{r}_2) = \sum_{j,k=1}^N \delta(\mathbf{r}_1 - \mathbf{r}_j)\delta(\mathbf{r}_2 - \mathbf{r}_k)$$
(12)

which is the joint probability that particle j is found at \mathbf{r}_1 and particle k at \mathbf{r}_2 . It is important that in general the average of this probability density is not just the product of the average densities:

$$\langle \rho(\mathbf{r}_1)\rho(\mathbf{r}_2)\rangle \neq \langle \rho(\mathbf{r}_1)\rangle \langle \rho(\mathbf{r}_2)\rangle = {\rho_0}^2$$
 (13)

 $(\rho_0 = N/V)$. The reason for this is that usually there is an interaction between particles which enhances or reduces the probability for particles close to each other. E.g. if one imagines particles with a hard core of radius R then $\langle \rho(\mathbf{r}_1)\rho(\mathbf{r}_2)\rangle$ vanishes for all \mathbf{r}_1 and \mathbf{r}_2 which would imply a 'collision' of the particles, $|\mathbf{r}_2 - \mathbf{r}_1| < 2R$. But still, in a translationally invariant system one of the positions can be chosen arbitrarily, especially as the origin, so that

$$\langle \rho(\mathbf{r}_1)\rho(\mathbf{r}_2)\rangle = \langle \rho(\mathbf{0})\rho(\mathbf{r}_2 - \mathbf{r}_1)\rangle = \rho_0 \left\langle \sum_{j,k=1}^N \delta(\mathbf{r}_j - \mathbf{r}_k + \mathbf{r}_2 - \mathbf{r}_1) \right\rangle.$$
 (14)

For a system of identical scatterers the two-particle density (12) can now be used to express the structure factor:

$$S(\mathbf{Q}) = \frac{1}{N} \left\langle \int_{V} \mathrm{d}^{3} r_{1} \int_{V} \mathrm{d}^{3} r_{2} \exp(\mathrm{i}\mathbf{Q} \cdot (\mathbf{r}_{2} - \mathbf{r}_{1}))\rho(\mathbf{r}_{1})\rho(\mathbf{r}_{2}) \right\rangle$$

$$= \frac{1}{\rho_{0}} \int_{V_{d}} \mathrm{d}^{3} r \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}) \left\langle \rho(\mathbf{0})\rho(\mathbf{r}) \right\rangle.$$
(15)

Note that in this last expression r does not have the meaning of an absolute position but that of a vectorial distance and consequently the volume of integration V_d is not the sample volume V but the volume of possible distances within the sample.

In the literature often alternative definitions of the pair correlation function are used (instead of using $\langle \rho(\mathbf{0})\rho(\mathbf{r})\rangle$ directly). The most common definition in the context of soft matter including biological systems is

$$g(\mathbf{r}) = \frac{\langle \rho(\mathbf{0})\rho(\mathbf{r})\rangle}{\rho_0^2} - \frac{\delta(\mathbf{r})}{\rho_0}.$$
 (16)

The normalisation by ρ_0^2 has the effect that for non-interacting particles or at distances where the interaction is weak, $g(\mathbf{r}) = 1$. The subtraction of the delta function removes the singularity

of $\langle \rho(\mathbf{0})\rho(\mathbf{r})\rangle$ at $\mathbf{r} = \mathbf{0}$ due to the j = k terms in (12). With this pair correlation function the structure factor can be written as

$$S(\mathbf{Q}) = 1 + \rho_0 \int_{V_d} \mathrm{d}^3 r \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r})(g(\mathbf{r}) - 1) \,. \tag{17}$$

Here, the 1 + compensates the delta function term subtracted in (16). In addition, one usually writes $g(\mathbf{r}) - 1$ instead of simply $g(\mathbf{r})$ in the Fourier transform. This avoids a delta function term arising in the limit $V_d \to \infty$ at Q = 0. In that limit, this 'trick' only changes the result at Q = 0 which is the (unobservable) forward scattering. Nevertheless, strictly speaking, one loses the scattering contribution by the overall sample shape. But this only affects the very low Q region if the sample has macroscopic dimensions $\gg 2\pi/Q$.

In many physical systems the interaction between particles is not directional with the consequence that $g(\mathbf{r})$ depends only on the distance $r = |\mathbf{r}|$. In this case by symmetry follows that also $S(\mathbf{Q})$ is only a function of $Q = |\mathbf{Q}|$ and equation (17) reduces to a one-dimensional integral:

$$S(Q) = 1 + \frac{4\pi\rho_0}{Q} \int_0^\infty (g(r) - 1)\sin(Qr)rdr.$$
 (18)

3 Scattering from dynamic systems

Here, the more realistic situation will be considered in which the particles of the sample are moving. Their dynamics will be described by *trajectories* $\mathbf{r}_j(t)$ which implies that the result is only valid in classical approximation. This assumption is usually uncritical for biological systems because experiments are done at high temperatures where quantum mechanical effects and 'recoil' (the change of the trajectories by the momentum transfer with the scattered radiation) are unobservable. In most scattering experiments the motion of the particles is observed indirectly via the inelasticity, i.e. the energy transfer $\hbar\omega = E' - E$ of the scattering process. The notable exception is dynamic light scattering where a time correlation function is measured directly.

In analogy to (3) the *double differential cross-section* is defined as the probability density that a neutron is scattered into a solid angle element $d\Omega$ with an energy transfer $\hbar\omega \dots \hbar(\omega + d\omega)$. Within a certain approximation, which is usually valid for biological systems, it can be given as²

$$\frac{\mathrm{d}\sigma}{\mathrm{d}\Omega\mathrm{d}\omega} = \frac{1}{2\pi} \frac{k'}{k} \int_{-\infty}^{\infty} \mathrm{e}^{-\mathrm{i}\omega t} \mathrm{d}t \left\langle \sum_{j,k=1}^{N} b_{j}^{*} b_{k} \exp\left(\mathrm{i}\mathbf{Q}\cdot\left(\mathbf{r}_{k}(t)-\mathbf{r}_{j}(0)\right)\right) \right\rangle \,. \tag{19}$$

It can be seen that the 'indirectness' of scattering leads to another Fourier transform in time $t \to \omega$ (in addition to that in space $r \to Q$). A factor which was not visible before because k' = k for scattering from static systems is the ratio of the wave numbers k'/k. This factor can be easily understood for massive particles as neutrons as the ratio of velocities v'/v which modulates the intensity of the scattered beam.

Another place where $k' \neq k$ has to be taken into account is that Q now does not anymore result from the isosceles construction in Fig. 2 drafted in black but from scattering triangles as those

²A correct derivation of the double differential cross section for inelastic scattering is only possible by quantum mechanical means. This is usually done in textbooks on neutron scattering as refs. 4, 5. It would be beyond the scope of this lecture to present this derivation in detail.

in blue and red. Application of the cosine theorem leads to the following expression for Q in the inelastic situation:

$$Q = \sqrt{k^2 + k'^2 - 2kk'\cos(2\theta)}$$

= $\sqrt{\frac{8\pi^2}{\lambda^2} + \frac{2m_n\omega}{\hbar} - \frac{4\pi}{\lambda}\sqrt{\frac{4\pi^2}{\lambda^2} + \frac{2m_n\omega}{\hbar}\cos(2\theta)}}$ for neutrons. (20)

Especially, it has to be observed now that Q also depends on $\hbar\omega$ implying that Q is not anymore constant for a single scattering angle. This consideration is relevant for neutron scattering but less for x-ray and light scattering where the incident energy usually is much larger than the energy transfer.

Expression (19) is valid for (nuclear, non-polarised) neutron scattering and the most simple when compared to light scattering or x-ray scattering. The expressions for the other scattering probes would include further terms involving the polarisation of the beam and the polarisability of the sample. Because these are not immediately relevant for biological systems, we will continue with this simplest formulation of inelastic scattering.

In order to derive a quantity similar to the structure factor (17), one assumes again a system of N chemically identical particles. But in order to capture the feature of incoherent scattering, present in neutron scattering, it is assumed that these particles do not have identical scattering lengths but individual randomly distributed scattering lengths with the average $\overline{b} = (1/N) \sum_j b_j$ and the variance $\overline{|b|^2} - |\overline{b}|^2 = \overline{|b - \overline{b}|^2} = (1/N) \sum_j |b_j - \overline{b}|^2$. The most obvious reason for the variance of scattering lengths is that chemically identical atoms may be different isotopes. Because the neutron scattering length is a nuclear property it may differ from isotope to isotope. But even in monisotopic systems there may be such a variance due to disorder of the nuclear spin orientations because the scattering length also depends on the combined spin state of the scattered neutron and the scattering nucleus.³

The sum in expression (19) can be decomposed into one over different indices and one over identical indices,

$$\sum_{j,k=1}^{N} b_{j}^{*} b_{k} \mathrm{e}^{\mathrm{i}\mathbf{Q} \cdot (\mathbf{r}_{k}(t) - \mathbf{r}_{j}(0))} = \sum_{j \neq k=1}^{N} b_{j}^{*} b_{k} \mathrm{e}^{\mathrm{i}\mathbf{Q} \cdot (\mathbf{r}_{k}(t) - \mathbf{r}_{j}(0))} + \sum_{j=1}^{N} |b_{j}|^{2} \mathrm{e}^{\mathrm{i}\mathbf{Q} \cdot (\mathbf{r}_{j}(t) - \mathbf{r}_{j}(0))}, \quad (21)$$

which have to be averaged in a different way with respect to the distribution of scattering lengths. In the first term b_j^* and b_k can be averaged separately because the different particle scattering lengths are uncorrelated: $\overline{b^*}\overline{b} = \overline{b}^*\overline{b} = |\overline{b}|^2$. In the second term one has to average *after* taking the absolute square:

$$= \sum_{j\neq k=1}^{N} |\overline{b}|^2 \mathrm{e}^{\mathrm{i}\mathbf{Q}\cdot(\mathbf{r}_k(t)-\mathbf{r}_j(0))} + \sum_{j=1}^{N} \overline{|b|^2} \mathrm{e}^{\mathrm{i}\mathbf{Q}\cdot(\mathbf{r}_j(t)-\mathbf{r}_j(0))} .$$
(22)

 $^{^{3}}$ It is more difficult to see how incoherent scattering may arise in light scattering. There the individual scattering centres may be colloidal particles with an effective *b* depending on the size and the refractive index. The analogue to isotope disorder would be here to have particles with identical size and interaction but different refractive indices. This trick is exploited to obtain incoherent scattering from colloidal systems [16] and thereby gain access to the self-correlation.

In order to avoid the sum over distinct particles, the first sum is complemented by the j = k terms, $|\bar{b}|^2 e^{i\mathbf{Q}\cdot(\mathbf{r}_j(t)-\mathbf{r}_j(0))}$, and to compensate, these terms are subtracted in the second sum:

$$= \sum_{j,k=1}^{N} |\bar{b}|^2 e^{i\mathbf{Q} \cdot (\mathbf{r}_k(t) - \mathbf{r}_j(0))} + \sum_{j=1}^{N} \left(\overline{|b|^2} - |\bar{b}|^2 \right) e^{i\mathbf{Q} \cdot (\mathbf{r}_j(t) - \mathbf{r}_j(0))} .$$
(23)

With this result it is possible to express the double differential cross section as

$$\frac{\partial \sigma}{\partial \Omega \partial \omega} = N \frac{k'}{k} \left(\left| \overline{b} \right|^2 S_{\text{coh}}(\mathbf{Q}, \omega) + \left(\overline{|b|^2} - \left| \overline{b} \right|^2 \right) S_{\text{inc}}(\mathbf{Q}, \omega) \right)$$
(24)

with

$$S_{\rm coh}(\mathbf{Q},\omega) = \frac{1}{2\pi N} \int_{-\infty}^{\infty} e^{-i\omega t} \mathrm{d}t \sum_{j,k=1}^{N} \left\langle e^{i\mathbf{Q}\cdot(\mathbf{r}_k(t)-\mathbf{r}_j(0))} \right\rangle$$
(25)

and

$$S_{\rm inc}(\mathbf{Q},\omega) = \frac{1}{2\pi N} \int_{-\infty}^{\infty} e^{-i\omega t} dt \sum_{j=1}^{N} \left\langle e^{i\mathbf{Q}\cdot(\mathbf{r}_j(t) - \mathbf{r}_j(0))} \right\rangle \,. \tag{26}$$

The quantities defined by (25) and (26) are called *coherent* and *incoherent scattering function* or *dynamic structure factors*.

The prefactors of the scattering functions in expression (24) are often expressed by the scattering cross sections

$$\sigma_{\rm coh} = 4\pi \left| \overline{b} \right|^2, \quad \sigma_{\rm inc} = 4\pi \left(\overline{|b|^2} - \left| \overline{b} \right|^2 \right) \tag{27}$$

which (for the incoherent part in general and for the coherent in the limit $Q \to \infty$) give the scattering into all directions, i.e. the solid angle 4π .

In some cases it is interesting to consider the part of expression (25) before the time-frequency Fourier transform, called *intermediate coherent scattering function*:

$$I_{\rm coh}(\mathbf{Q},t) = \frac{1}{N} \sum_{jk} \left\langle e^{i\mathbf{Q} \cdot (\mathbf{r}_k(t) - \mathbf{r}_j(0))} \right\rangle$$
(28)

 $(g_1(t)$ in dynamic light scattering). Its value for t = 0 expresses the correlation between atoms *at equal times*. A theorem on Fourier transforms tells that this is identical to the integral of the scattering function over all energy transfers:

$$I_{\rm coh}(\mathbf{Q},0) = \frac{1}{N} \sum_{jk} \left\langle e^{i\mathbf{Q}\cdot(\mathbf{r}_k - \mathbf{r}_j)} \right\rangle = S(\mathbf{Q}) = \int_{-\infty}^{\infty} S_{\rm coh}(\mathbf{Q},\omega) \mathrm{d}\omega \,.$$
(29)

The concrete significance of this relation is that a diffraction experiment, which does not discriminate energies and thus implicitly integrates over all $\hbar\omega$, only shows the instantaneous correlation of the atoms, viz the structure of the sample⁴. $S(\mathbf{Q})$ is the structure factor as derived in section 2 for the static situation. The dynamic information is lost in the integration process.

⁴Strictly speaking, this is only an approximation. There are several reasons why the integration in the diffraction experiment is not the 'mathematical' one of (29): (1) On the instrument the integral is taken along a curve of constant 2θ while constant Q would correspond to a straight line in the Q- ω plane. (2) The double differential cross-section (24) contains a factor k'/k which depends on ω via (20). (3) The detector may have an efficiency depending on wavelength which will introduce another ω -dependent weight in the experimental integration. All these effects can be taken into account in the so-called Placzek corrections [17–19].

Similarly the incoherent intermediate scattering function is

$$I_{\rm inc}(\mathbf{Q},t) = \frac{1}{N} \sum_{j=1}^{N} \left\langle e^{i\mathbf{Q} \cdot (\mathbf{r}_j(t) - \mathbf{r}_j(0))} \right\rangle$$
(30)

with

$$I_{\rm inc}(\mathbf{Q},0) = \frac{1}{N} \sum_{j=1}^{N} \left\langle e^{i\mathbf{Q} \cdot (\mathbf{r}_j - \mathbf{r}_j)} \right\rangle = 1 = \int_{-\infty}^{\infty} S_{\rm inc}(\mathbf{Q},\omega) d\omega \,. \tag{31}$$

Note that this result is independent of the actual structure of the sample. Integration of the double-differential cross section (24) over ω shows that also the static scattering contains an incoherent contribution. But because of (31), this term is constant in Q. It constitutes a flat background in addition to the S(Q)-dependent scattering. In some cases (e.g. small-angle scattering) it may be necessary to correct for this, in other cases (e.g. diffraction with polarisation analysis) it may even be helpful to normalise the coherent scattering.

As in the static situation, the scattering law can be traced back to distance distribution functions, the *van Hove correlation functions*, which are time-dependent:

$$G(\mathbf{r},t) = \frac{1}{N} \left\langle \sum_{j,k=1}^{N} \delta(\mathbf{r} - \mathbf{r}_{k}(t) + \mathbf{r}_{j}(0)) \right\rangle, \qquad (32)$$

$$G_{\rm s}(\mathbf{r},t) = \frac{1}{N} \left\langle \sum_{j=1}^{N} \delta(\mathbf{r} - \mathbf{r}_j(t) + \mathbf{r}_j(0)) \right\rangle.$$
(33)

Insertion into

$$I_{\text{[coh]inc]}} = \int_{V_d} G_{[\text{s}]}(\mathbf{r}, t) \exp(\mathrm{i}\mathbf{Q} \cdot \mathbf{r}) \mathrm{d}^3 r$$
(34)

directly proves that the spatial Fourier transforms of the van Hove correlation functions are the intermediate scattering functions.

The two particle version can-as in the static case-be reduced to the microscopic density,

$$\rho(\mathbf{r},t) = \sum_{j=1}^{N} \delta(\mathbf{r} - \mathbf{r}_j(t)).$$
(35)

Again as in the static case $\rho(\mathbf{r})$, $\rho(\mathbf{r}, t)$ can also be seen as a generalisation to the situation of scattering from continua. Its autocorrelation function in space and time is

$$\langle \rho(\mathbf{0},0)\rho(\mathbf{r},t)\rangle$$
. (36)

The 0 is again showing that translational symmetry is assumed. So the correlation function can be replaced by its average over all starting points r_1 in the sample volume:

$$\left\langle \rho(\mathbf{0},0)\rho(\mathbf{r},t)\right\rangle = \frac{1}{V}\int_{V} \mathrm{d}^{3}r_{1}\left\langle \rho(\mathbf{r}_{1},0)\rho(\mathbf{r}_{1}+\mathbf{r},t)\right\rangle.$$
(37)

Insertion of (35) gives

$$\langle \rho(\mathbf{0},0)\rho(\mathbf{r},t)\rangle = \frac{1}{V} \left\langle \sum_{j,k=1}^{N} \int_{V} \mathrm{d}^{3}r_{1}\delta(\mathbf{r}_{1}-\mathbf{r}_{k}(t))\delta(\mathbf{r}_{1}+\mathbf{r}-\mathbf{r}_{j}(t)) \right\rangle$$
(38)

$$= \frac{1}{V} \left\langle \sum_{j,k=1}^{N} \delta(\mathbf{r}_{k}(t) + \mathbf{r} - \mathbf{r}_{j}(t)) \right\rangle$$
(39)

$$G(\mathbf{r},t) = \frac{1}{\rho_0} \langle \rho(\mathbf{0},0)\rho(\mathbf{r},t) \rangle \,. \tag{40}$$

Again setting t = 0 results in the static scattering situation:

$$G(\mathbf{r},0) = \frac{\langle \rho(\mathbf{0},0)\rho(\mathbf{r},0)\rangle}{\rho_0} = \delta(\mathbf{r}) + \rho_0 g(\mathbf{r})$$
(41)

with $g(\mathbf{r})$ from equation (16).

As in the case of static scattering there is an alternative way to derive the scattering function by first Fourier-transforming the density

$$\rho_{\mathbf{Q}}(t) = \int \mathrm{d}^3 r \mathrm{e}^{\mathrm{i}\mathbf{Q}\cdot\mathbf{r}} \rho(\mathbf{r}, t) = \sum_{j=1}^{N} \mathrm{e}^{\mathrm{i}\mathbf{Q}\cdot\mathbf{r}_j(t)}$$
(42)

and then multiplying its conjugated value at t = 0 with that at t:

$$I_{\rm coh}(\mathbf{Q},t) = \frac{1}{N} \left\langle \rho_{\mathbf{Q}}^*(0) \rho_{\mathbf{Q}}(t) \right\rangle \tag{43}$$

and

$$S_{\rm coh}(\mathbf{Q},\omega) = \frac{1}{2\pi N} \int_{-\infty}^{\infty} e^{-i\omega t} \left\langle \rho_{\mathbf{Q}}^*(0)\rho_{\mathbf{Q}}(t) \right\rangle \mathrm{d}t \,. \tag{44}$$

(This is a general consequence of the cross-correlation theorem of Fourier transform.)

Note that a reduction of the self correlation function $G_s(\mathbf{r}, t)$ is *not* possible in the same way because the multiplication $\rho(\mathbf{0}, 0)\rho(\mathbf{r}, t)$ inevitably includes all combinations of particles j, k and not only the terms for identical particles j, j.

The derivation of inelastic scattering presented here, starting from expression (19), was introduced by Vineyard in 1958 [20]. It was discovered already four years later [21] that this is a very rough approximation, among other shortcomings violating *detailed balance*. While equations (25) and (26) with reasonable assumptions about the symmetry of the system imply $S(Q, -\omega) = S(Q, \omega)$, detailed balance requires

$$S(Q, -\omega) = \exp\left(\frac{\hbar\omega}{k_{\rm B}T}\right) S(Q, \omega) \,. \tag{45}$$

In equilibrium, the probability for the scattering system to be in the lower energy state is higher by the factor $\exp(\hbar\omega/k_BT)$. Therefore scattering into a state with higher energy is more probable by the same factor than scattering into the lower energy state. This results from thermodynamics and is independent from the approximation to neglect quantum mechanics. Fortunately, in the case of biological systems most of the time one deals with slow processes at comparatively high temperatures so that $\hbar\omega \ll k_BT$ and the approximations used here are valid nevertheless. A more exact treatment (even on the semi-classical level) would make it necessary to know the momenta of the particles in addition to their trajectories [21]. In case of absence of such information and if it is necessary to access high energy transfers $\hbar\omega$ it is often justified to 'force' detailed balance by an ad hoc factor [22]:

$$S^{\text{corrected}}(Q,\omega) = \exp\left(-\frac{\hbar\omega}{2k_{\text{B}}T}\right)S(Q,\omega).$$
 (46)

4 Scattering probes: photons, neutrons, electrons

The most obvious difference between different scattering probes, i.e. radiation types, is their wavelength and frequency. In the particle picture these correspond to the momentum $(p = h/\lambda)$ and the energy (E = hf) of the particles. For each radiation type there is a relation between both quantities based on the rest mass of the particle, m_0 :

$$\lambda = \frac{h}{\sqrt{2Em_0 + E^2/c^2}} \,. \tag{47}$$

Note that E is the net kinetic energy here, thus $E = (m - m_0) c^2$ relativistically. For the photon, a massless particle, the relation reduces to

$$\lambda = \frac{hc}{E} = \frac{c}{f} \tag{48}$$

with f = E/h being the (common) frequency. So for photons the wavelength decreases inverseproportionally with frequency, $\lambda \propto E^{-1}$. Conversely in the non-relativistic limit ($m_0 \gg E/c^2$) which would be fulfilled for neutrons in general, one obtains:

$$\lambda = \frac{h}{\sqrt{2Em_0}} = \frac{h}{p} \tag{49}$$

with p being the momentum of the particle. Thus for massive (non-relativistical) particles the wavelength only decreases as $\lambda \propto E^{-1/2}$.

Fig. 3 shows the relation between wavelength and energy/frequency for photons, neutrons, and electrons. It can be seen that in the determination of molecular structure, neutrons and x-rays are equivalent because they both have adequate wavelengths. Nevertheless for inelastic scattering experiments x-rays have an energy several orders of magnitude higher at the same wavelength. Therefore the energy transfers which correspond to processes in a time range of nanoseconds and above are much more difficult to resolve by x-ray scattering.

In order to dwell on this argument in more detail, in the following it is explained how the wavelength and energy of the probe beam determine the resolution and range of a scattering method in time and space:

- The minimal distance resolved in space is related directly to the wavelength. This can be seen from Bragg's law, λ = 2d sin θ. Because sin θ ≤ 1 the distance (of Bragg planes) is limited by d ≥ λ/2⁵. Expressed in terms of the scattering vector, equation (4), this corresponds to an upper limit of the scattering vector Q ≤ 4π/λ.
- The maximal distance accessible is also proportional to the wavelength, but here the instrumental resolution is also a determining factor. Again in the picture of Bragg's law, if the smallest technically accessible scattering angle is θ_{\min} then the largest observable Bragg plane distance is $d \leq \lambda/2 \sin \theta_{\min}$. So here the wavelength of the beam does not pose a limit on principle. Nevertheless, even on dedicated small-angle scattering instruments one is not able to reach angles below 0.01° . Therefore, there is a technological limit at about four orders of magnitude above the wavelength of the radiation used in a scattering experiment, $d \leq 10^4 \lambda$.

⁵It has to be noted that in the absence of a crystalline structure this relation is only an estimate. In principle, the whole range of distances r enters a Fourier transform expression as (18) for a given Q. Nevertheless, a relation as $d_{\min} = 2\pi/Q_{\max}$ gives the smallest structure which can be technically resolved also in this case, but with a prefactor which may differ from 2π depending on the definition of what "resolve" means.



Fig. 3: Relation between energy (E) / frequency (f) and wavelength (λ) of different scattering probes. Red: photons, pink: electrons, blue: neutrons. The broadened sections of the curves represent the ranges actually used in scattering experiments on biological systems.

- The minimal time resolvable in dynamics is inversely related to the frequency of the radiation, or the particle energy. Here, the argument is somewhat more complicated because the energy transfer $E = \hbar \omega$, especially in neutron scattering, may be larger than the energy hf. But as equation (20) implies, large energy transfers lead to a strongly offset Qrange. So there is a combined restriction of time- and space range which usually prevents energy transfers to be used which are more than an order of magnitude higher than the incident energy. This leads to a restriction t > 0.2 ps meV/E.
- For the maximal time accessible, as in the case of the maximal distance, the instrumental resolution is a crucial factor. The limit in terms of the minimal resolvable energy transfer is $t < 2 \,\mathrm{ps} \,\mathrm{meV}/\Delta E$ (with a slight dependence of the prefactor on the definition of "accessible time"). But the relative instrumental resolution $\Delta E/E$ differs by orders of magnitude from instrument to instrument. Conventional neutron scattering instruments reach down to 10^{-3} while due to the high flux available inelastic experiments on synchrotrons can reach 10^{-7} . The neutron spin echo technique (see subsection 5.2) can reach times corresponding to $\Delta E/E = 10^{-6}$. A method opening an on-principle unlimited time range is the use of *intensity autocorrelation* on the scattered beam. For this method the resolution is not restricted by the monochromatisation of the incident beam. This is a standard technique for photons, where times up to 10^3 s can be reached. The corresponding to creating autocorrelation are up to 10^3 s can be reached. The corresponding to creating is currently under development reaching roughly the same limit in certain applications.

Fig. 4 shows the spatial dimensions accessible by different elastic scattering methods. It can be seen that for neutrons and x-rays usually small scattering angles have to be used to reach the dimensions of biological objects. On the other hand, light scattering is only suited to very large biological objects as cells and organelles.



Fig. 4: Spatial dimensions accessible by different scattering methods. (Note that these are typical ranges; particular experimental set-ups may extend these.)

Fig. 5 shows in dimensions of space and time the range of different inelastic scattering methods. It can be seen that light scattering covers an enormous range in the dynamics due to the availability of the photon correlation technique. On the other hand it suffers from the long wavelength prohibiting access to shorter length scales. X-rays do cover nanometre sizes but due to their high energy only allow to resolve very fast processes (≤ 1 ps) with filter methods. X-ray photon correlation methods are currently limited to longer times than some microseconds. There remains a large region, roughly the nanometre-nanosecond range, which can only be covered by inelastic neutron scattering.

Another property of a probe beam important for the application is its interaction mechanism with the sample. Photons of long wavelength (light) interact with the dielectric function of the sample. Thus, any sample showing a fluctuation of the index of refraction is suitable for scattering. On the other hand, short-wavelength photons, x-rays, interact directly with the electron density. Therefore, their interaction is stronger with elements of higher atomic number Z (Fig. 6). At roughly the same wavelength neutrons do not show the monotonic dependence on Z and in addition usually the scattering length depends on the isotope. Especially, the scattering lengths of hydrogen and deuterium are completely different. In organic materials this allows to create a contrast without change of the chemical properties.

5 Scattering methods for biological matter investigations

Arguably, the most important scattering technique applied to biological systems is protein crystallography. For that reason, the separate lecture A5 is devoted to this topic. In this lecture only the most important method for non-crystalline samples, small angle (x-ray or neutron) scattering will be presented. For dynamical scattering the neutron spin echo (NSE) technique will be explained. Concrete examples of its application will be shown in lecture B2.



Fig. 5: Ranges in time and space accessible by different scattering methods. PFG-NMR has been added because it provides the intermediate scattering function I(Q,t) as genuine scattering methods do. (Note that these are typical ranges; particular experimental set-ups may extend these.)



Fig. 6: Scattering cross sections for neutrons and x-rays. Red symbols: neutron scattering cross section $\sigma_{\rm coh} + \sigma_{\rm inc}$. The small points indicate individual isotopes and the large open circles averages for the natural composition of the elements. Blue line: x-ray scattering cross section. The unit of the cross section is $1 \text{ barn} = 10^{-28} \text{ m}^2$.

5.1 Small-angle scattering

While single-crystal diffraction allows a detailed study of biological objects (mostly proteins) in the crystallised state, it is often desirable to study their structure in solution or even in their natural environment, which may be a lipid bilayer. This could be either because the objects do not crystallise at all or because it is suspected that the structure in the crystal differs from the 'natural' one. In that kind of scattering experiments [23], due to the random placement and orientation of the objects, the scattering is not concentrated in Bragg peaks but forms a diffuse continuum. Because biological object dimensions ($\gtrsim 1 \, \mathrm{nm}$) usually exceed the wavelength of x-rays ($\lesssim 0.1 \,\mathrm{nm}$) by more than an order of magnitude, x-ray scattering experiments require small angles. Small-angle x-ray scattering (SAXS) is conceptually not different from ordinary Debye-Scherrer diffraction but the small scattering angle often requires some technical 'tricks': (1) All such instruments have to employ an evacuated flight path to avoid scattering from air. (2) The small solid angles of collimation and detection lower the intensities. Therefore either strong sources are necessary (rotating anode tubes for tabletop applications, or a synchrotron) or point collimation has be replaced by slit collimation (Kratky camera [24]). In the latter case the gain in intensity has to be paid for by a distortion of the scattering curves which has to be corrected mathematically. Several more advanced SAXS set-ups are discussed in ref. 25. (3) Because biological materials consist predominantly of elements with lower atomic number than their sample containers, extremely thin-walled containers have to be used (e.g. Mark capillaries).

Neutron scattering experiments on biological matter typically use 'cold neutrons', i.e. neutrons where the kinetic energy is reduced by multiple collisions in a moderator which is kept far below room temperature. A typical temperature would be T = 20 K corresponding to an energy $k_{\rm B}T = 1.7 \,\mathrm{meV}$ and a wavelength $\lambda = 0.7 \,\mathrm{nm}$. Although this wavelength is closer to, e.g., the dimensions of a protein molecule than that of thermal neutrons ($\approx 0.2 \,\mathrm{nm}$), the observation of its internal structure usually requires small scattering angles (SANS) [25, 26]. The standard instrument is a scaled-up pinhole SAXS set-up. Because the width of the neutron beam is usually about a centimetre, 10–100 times larger than an x-ray beam, the whole instrument has to be enlarged by this factor resulting in flight paths up to 40 m. With this conventional instrument, scattering vectors down to $Q = 10^{-2} \,\mathrm{nm}^{-1}$ can be achieved. For lower Q values, focussing elements as mirrors or neutron lenses have to be implemented. The beam size is reduced by apertures to some millimetres and the focussing partially compensates the loss in intensity. The limit of this technology lies at about $10^{-3} \,\mathrm{nm}^{-1}$. For even lower Q, Bonse-Hart set-ups are used as in x-ray scattering (limit $\approx 10^{-4} \,\mathrm{nm}^{-1}$).

The crucial advantage of SANS compared to SAXS, justifying the far higher experimental effort, is the possibility of *contrast variation*. While in x-ray scattering the cross-section is invariably coupled to the atomic number, it can be varied in neutron scattering by isotopic substitution. This allows e.g. to create a scattering contrast between chemically identical molecules. Conversely, one can make structures 'invisible' by 'contrast-matching' certain moeities in a sample. This is done by choosing an isotopic composition which makes their average scattering length densities (see eq. (8)) identical.

The main disadvantage of small-angle scattering methods is that due to the random orientation of the scatterers they only produce $d\sigma/d\Omega$ as a one dimensional function of Q as a result. It is immediately clear that this does not allow an unambiguous reconstruction of a threedimensional structure (as it is possible for single crystal diffraction). Nevertheless, it *is* possible to calculate $d\sigma/d\Omega$ from a given arrangement of atoms [27]. Thus one can verify whether a structure of a protein in its crystalline form is still valid in solution. But in the opposite



Fig. 7: Structural modeling of the arrangement of the three domains in TIA-1. a) Schematics of the procedure. b) SAXS and c) SANS data compared to calculated diffraction patterns. d) χ^2 deviation ranges, left: all 5000 test structures, right: finally selected five structures. Reprinted from ref. 34

direction only a low-resolution structure determination can be done. This is possible e.g. by multipole representation of the shape [28, 29], finite elements [30, 31], or dummy residues [32]. The amount of information in the scattering experiment can be increased by combining SANS and SAXS (different scattering length densities), contrast variation, and selective deuteration of parts of the molecule. Because chemical synthesis is restricted to short peptides for deuteration one often has to resort to biosynthesis in deuterated media [33]

As an example using some of these techniques Fig. 7 shows schematically the structure determination of a RNA-binding protein (RBP), TIA-1 [34]. Although the structures of the three subunits RRM (RNA recognition motif) 1–3 are known, their relative arrangement by flexible linkers in the presence of RNA is unclear. To investigate this, 5000 test structures of the whole complex were generated by molecular dynamics calculations. For all of these the expected SAXS patterns were calculated using CRYSOL [27]. In this way, all but 67 structures could be ruled out. Subsequently, these structures were similarly tested against four SANS patterns differing in the D_2O content of the solvent and deuteration of the subunits. Finally, only five structures remained consistent with the data. From these structures certain general characteristics of the spatial arrangement could be inferred, e.g. that the complex takes an elongated L-shape. (For more detaily please refer to the original publication.)



Fig. 8: Schematic setup of a neutron spin echo spectrometer. x, y, and z denote the general orientation of the coordinate system used in the text with the reservation that z should follow the neutron flight direction.

5.2 Neutron spin echo spectroscopy

While the diffusional motion of biological molecules (mostly proteins) can easily be studied by dynamic light scattering (DLS), the internal dynamics of such molecules happens on a length scale too small to be resolved by light. On the other hand, this internal dynamics is usually slow, so that the inelastic neutron scattering method with the best resolution, i.e. neutron spin echo (NSE) spectroscopy [35] has to be used.

The high resolution is achieved by measuring the *individual* velocities of the incident and scattered neutrons using the Larmor precession of the neutron spin in a magnetic field. The neutron spin vector acts as the hand of an internal clock which is linked to each neutron and connects the result of the velocity measurement to the neutron itself. Thereby the velocities before and after scattering on one and the same neutron can be compared and a direct measurement of the velocity difference becomes possible. The energy resolution is thus decoupled from the monochromatisation of the incident beam. Relative energy resolutions in the order of 10^{-5} can be achieved with an incident neutron spectrum of 20% bandwidth.

The motion of the neutron polarisation P(t)—which is the quantum mechanical expectancy value of the neutron spin—is described by the Bloch equation

$$\frac{\mathrm{d}\mathbf{P}}{\mathrm{d}t} = \frac{\gamma\mu}{\hbar} (\mathbf{P} \times \mathbf{B}) \tag{50}$$

where γ is the gyromagnetic ratio ($\gamma = -3.82$) of the neutron, μ the nuclear magneton and B the magnetic field. Equation (50) is the basis for manipulation of the neutron polarisation by external fields. A simple calculation shows that (50) predicts a precession with the *Larmor* frequency $\gamma \mu B/\hbar$ in a constant magnetic field. Thus, if a neutron of wavelength λ is exposed to a magnetic field B over a length l of its flight path its spin is rotated by

$$\phi = \left(\frac{2\pi |\gamma| \mu \lambda m_{\rm n}}{h^2}\right) Bl \,. \tag{51}$$

The basic setup of an NSE spectrometer is shown in figure 5.2. A velocity selector in the primary neutron beam selects a wavelength interval of 10-20% width. In the primary and secondary flight path of the instrument precession fields *B* and *B'* parallel to the respective path are generated by cylindrical coils. Before entering the first flight path the neutron beam is
Table 1: Evolution of the neutron spin (classical) in the NSE spectrometer. The left column contains the cartesian components of the normalised spin vector \mathbf{s} . The right column shows the neutronic devices the beam passes on its way through the spectrometer. z is defined as always denoting the direction parallel to neutron propagation.

(0,0,1)	
$\pi/2$ flippe	er
(1,0,0) field .	В
$(\cos\phi, -\sin\phi, 0) = \pi \text{ flipped}$	er
$(\cos(-\phi), \sin(-\phi), 0)$ field <i>B</i>	3′
$\pi/2 \text{ flippo}$ $(0, \sin(\phi - \phi'), \cos(\phi - \phi'))$	er

polarised in forward direction⁶. Firstly, a $\pi/2$ flipper rotates the polarisation to the x direction perpendicular to the direction of propagation (z). This is done by exposing the neutrons to a well-defined field for a time defined by their speed and the thickness of a flat coil (Mezei coil). Beginning with this comparatively well-defined initial condition the neutrons start their precession in the field B. After being scattered by the sample the neutrons go through a π flipper. The effect of the π flipper amounts to reverting the precession angle accumulated before scattering. Then the neutrons pass through the second precession field B'. Finally, the neutrons pass through another $\pi/2$ coil which under certain conditions restores their initial polarisation parallel to their flight direction. In order to understand what these conditions are one has to trace the changes of the spin vector as shown in Table 1. In total, the spin is rotated by $\phi - \phi'$ around the x axis when a neutron passes through the spectrometer. This means that the final polarisation is identical to the incident if $\phi = \phi' (+2\pi n)$, especially if $\lambda_i = \lambda_f$ (elastic scattering) and $\int_0^l Bdz = \int_0^{l'} B' dz$ (for homogeneous fields: Bl = B'l') as follows from (51). This condition is called "spin echo" and is independent of the individual velocities of the neutrons because their difference alone determines $\phi - \phi'$.

Leaving spin echo condition, the probability of a single neutron to reach the detector is reduced due to the polarisation analyser by $\cos(\phi' - \phi)$. If we keep the symmetry of the instrument, Bl = B'l', but consider inelastic scattering the precession angle mismatch follows from eq. (51):

$$\phi' - \phi = \left(\frac{2\pi |\gamma| \mu m_{\rm n}}{h^2}\right) Bl(\lambda_{\rm f} - \lambda_{\rm i})$$

$$\approx \underbrace{\frac{|\gamma| \mu m_{\rm n}^2 \lambda^3 Bl}{h^3}}_{=t_{\rm NSE}(B)} \omega$$
(52)

The approximation in the second line is valid for small energy transfers where $\Delta \lambda \approx \hbar \omega / \frac{dE}{d\lambda}$

⁶This is done by a 'polarising supermirror' [36] which only reflects neutrons of that spin—similar to the Nicol prism in optics.

can be used. Because the energy transfer for inelastic scattering is not fixed but distributed as determined by the scattering function $S(\mathbf{Q}, \omega)$ we have to average the factor $\cos(\phi' - \phi)$ weighted by $S(\mathbf{Q}, \omega)$ to get the reduction of count rate at the detector, the effective polarisation

$$P(\mathbf{Q}, t_{\rm NSE}) = \frac{\int_{-\infty}^{\infty} S(\mathbf{Q}, \omega) \cos(\omega t_{\rm NSE}) d\omega}{\int_{-\infty}^{\infty} S(\mathbf{Q}, \omega) d\omega} \,.$$
(53)

Firstly, we note that $S(\mathbf{Q}, \omega)$ in this expression usually is the coherent scattering function with the addition the isotope part of the incoherent scattering. In principle, similar arguments can be used for spin-incoherent scattering because a well-defined fraction (2/3) of neutrons changes its spin. This leads to a "negative echo" because the majority of neutrons invert their polarisation. But because this effect is only partial it is much more difficult to observe. Only recently, NSE spectroscopy could be applied successfully to incoherently scattering samples [37]. Nevertheless, because biological systems contain a large amount of hydrogen, the usual procedure is to replace water by D₂O which eliminates the majority of spin-incoherent scattering.

Secondly, expression (53) reverses the temporal Fourier transform⁷ in equation (25) or (26). Therefore, NSE provides the intermediate scattering function normalised to the structure factor

$$P(\mathbf{Q}, t_{\text{NSE}}) = \frac{I(\mathbf{Q}, t_{\text{NSE}})}{S(\mathbf{Q})}$$
(54)

in the coherent case or just the intermediate scattering function in the incoherent case.

From equation (52) can be seen that the maximum time accessible by NSE increases with the third power of the incident neutron wavelength λ . Therefore, a statement about its limits depends strongly on the required Q value because large Q makes a small λ necessary (eq. (4)). Also for large λ there may be intensity problems because long-wavelength neutrons are not sufficiently present in the source spectrum. For standard applications on the dynamics of biological molecules the limit in Fourier time lies at about 350 ns.

Many variations of the basic NSE concept explained here are used in currently operated NSE spectrometers. The most straightforward, employed at most such instruments, is the use of multidetectors to observe several Q values simultaneously. The technical problem resulting from this measure is that the field integral has to be kept identical for all possible neutron flight paths by elaborate correction coils [38]. For NSE instruments which are mostly operated at low Q, similar focussing techniques as in (static) small-angle scattering can be used, e.g. mirrors [39]. NSE can be combined with the time-of-flight principle in the primary spectrometer. The components of a neutron pulse enter the spectrometer in the sequence from fast to slow. Thus, the NSE spectrometer can be operated sequentially with all incident wavelengths contained in the pulse. In this way a broad range of Q-t combinations is covered simultaneously. For pulsed sources this leads to a clear intensity gain [40]. It is also possible to use the TOF-NSE principle at a continuous source [39], but then the advantage is partially outweighed by the loss of neutrons during the shaping of the pulses by a chopper.

A variant of the NSE technique described here replaces the main precession coils with a static field by smaller radio-frequency driven coils (neutron *resonance* spin echo, NRSE) [41]. It is even possible to modify this set-up in a way that it does not need the encoding of information

⁷The cosine Fourier transform here and the exponential Fourier transform of (25) or (26) are only identical under the conditions used to derive inelastic scattering in section 3. Then $S(Q, \omega)$ is symmetric and I(Q, t) is a real function. The condition $k_{\rm B}T \gg \hbar\omega$ is fulfilled here because soft matter experiments are done around room temperature and NSE is used to detect small energy transfers.

into the neutron spin anymore but uses a modulation of intensity (modulation of intensity by zero effort, MIEZE) [42]. In this way the mentioned problems for spin-incoherent scattering can be avoided. Although the last technique would be beneficial for biological samples it is usually not employed because its disadvantages on the side of intensity and time range outweigh the effort for deuteration.

References

- P. Lindner, T. Zemb (eds.), Neutrons, X-rays and Light: Scattering Methods Applied to Soft Condensed Matter (Elsevier, Amsterdam, 2002).
- [2] R. J. Roe, *Methods of X-ray and neutron scattering in polymer science* (Oxford University Press, Oxford, 2000).
- [3] G. E. Bacon, Neutron Diffraction (Clarendon Press, Oxford, 1975).
- [4] G. L. Squires, *Introduction to the theory of thermal neutron scattering* (Cambridge University Press, Cambridge, 1978).
- [5] S. W. Lovesey, *Theory of Neutron Scattering from Condensed Matter* (Clarendon Press, Oxford, 1984).
- [6] M. Bée, Quasielastic neutron scattering (Adam Hilger, Bristol, 1988).
- [7] T. Brückel, G. Heger, D. Richter, R. Zorn (eds.), *Neutron Scattering* (Forschungszentrum Jülich, 2008), http://hdl.handle.net/2128/3718.
- [8] B. Chu, Laser Light Scattering (Academic Press, New York, 1974).
- [9] B. J. Berne, R. Pecora, Dynamic Light Scattering (Wiley, New York, 1976).
- [10] H. Z. Cummins, E. R. Pike (eds.), *Photon Correlation Spectroscopy and Velocimetry* (Plenum Press, New York, 1977).
- [11] R. Pecora, Dynamic Light Scattering (Plenum Press, New York, 1985).
- [12] W. Brown, Light Scattering: Principles and Development (Oxford University Press, Oxford, 1996).
- [13] R. Hosemann, A. N. Bagchi, *Direct analysis of diffraction by matter* (North-Holland, Amsterdam, 1962).
- [14] L. V. Azaroff, R. Kaplow, N. Kato, R. J. Weiss, A. J. C. Wilson, R. A. Young, X-ray diffraction (McGraw-Hill, Columbus, 1974).
- [15] O. Kratky, O. Glatter, Small Angle X-Ray Scattering (Academic Press, London, 1982).
- [16] W. van Megen, T. C. Mortensen, S. R. Williams, J. Müller, Phys. Rev. E 58, 6073 (1998).
- [17] R. Zorn, D. Richter, *Correlation Functions Measured by Scattering Experiments*, in ref. 7, chapter 4.
- [18] G. Placzek, Phys. Rev. 86, 377 (1952).
- [19] J. L. Yarnell, M. J. Katz, R. G. Wenzel, S. H. Koenig, Phys. Rev. A 7, 2130 (1973).
- [20] G. H. Vineyard, Phys. Rev. 110, 999 (1958).
- [21] R. Aamodt, K. M. Case, M. Rosenbaum, P. F. Zweifel, Phys. Rev. 126, 1165 (1962).

- [22] P. Schofield, Phys. Rev. Lett. 4, 239 (1960).
- [23] D. I. Svergun, M. H. J. Koch, P. A. Timmins, R. P. May, Small Angle X-Ray and Neutron Scattering from Solutions of Biological Macromolecules (Oxford University Press, Oxford, 2013).
- [24] O. Kratky, Experimental Techniques, Slit Collimation, in ref. 15.
- [25] J. S. Pedersen, Instrumentation for Small-Angle X-Ray and Neutron Scattering and Instrumental Smearing Effects, in ref. 1, pp. 127–144.
- [26] H. Frielinghaus, Small-Angle Scattering, in ref. 7, chapter 8.
- [27] D. I. Svergun, C. Barberato, M. H. J. Koch, J. Appl. Cryst. 28, 768 (1995).
- [28] H. B. Stuhrmann, Z. Phys. Chem. Neue Fol. 72, 177 (1970).
- [29] D. I. Svergun, J. Appl. Cryst. 30, 792 (1997).
- [30] D. Franke, D. I. Svergun, J. Appl. Cryst. 42, 342 (2009).
- [31] A. Koutsioubas, S. Jaksch, J. Perez, J. Appl. Cryst. 49, 690 (2016).
- [32] M. V. Petoukhov, D. I. Svergun, J. Appl. Cryst. 36, 540 (2003).
- [33] M. Haertlein, M. Moulin, J. M. Devos, V. Laux, O. Dunne, V. T. Forsyth, Biomolecular Deuteration for Neutron Structural Biology and Dynamics, in Isotope Labeling of Biomolecules - Applications, chapter 5 (Elsevier, Cambridge, 2016).
- [34] M. Sonntag, P. K. A. Jagtap, B. Simon, M.-S. Appavou, A. Geerlof, R. Stehle, F. Gabel, J. Hennig, M. Sattler, Angew. Chem. Int. Ed. 56, 9322 (2017).
- [35] F. Mezei, C. Pappas, T. Gutberlet (eds.), *Neutron Spin Echo Spectroscopy* (Springer, Berlin, 2003).
- [36] F. Mezei, Commun. Phys. 1, 81 (1976).
- [37] M. C. Bellisent-Funel, R. Daniel, D. Durand, M. Ferrand, J. L. Finney, S. Pouget, V. Reat, J. C. Smith, J. Am. Chem. Soc. 120, 7347 (1998).
- [38] M. Ohl, M. Monkenbusch, T. Kozielewski, B. Laatsch, C. Tiemann, D. Richter, Physica B 356, 234 (2005).
- [39] https://www.ill.eu/instruments-support/instruments-groups/ instruments/in15.
- [40] https://neutrons.ornl.gov/nse.
- [41] R. Gähler, R. Golub, Physica B & C 49, 1195 (1988).
- [42] M. Koppe, M. Bleuel, R. Gähler, R. Golub, P. Hank, T. Keller, S. Longeville, U. Rauch, J. Wuttke, Physica B 266, 75 (1999).

A 5 Macromolecular Crystallography

T. E. Schrader Jülich Centre for Neutron Science Forschungszentrum Jülich GmbH

Contents

1	Intr	oduction	2
2	The	Physics of Macromolecular Crystallography	4
	2.1	An X-ray macromolecular crystallography beamline	4
	2.2	A neutron macromolecular crystallography instrument	5
	2.3	Some general aspects of diffraction by a protein crystal	7
	2.4	Model building and refinement	10
3 4	A fi A se	rst case study: Water network around myoglobin	12
Ack	knowle	dgement	
Ref	erence	28	16
Rec	omme	nded Textbooks	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Ideally, one knows each atom position in a cell or even in the whole body of an organism at any given time. With this amount of data one could reconstruct each metabolic pathway. One could understand for example the synthesis of fatty acids or the immune response to an HI-Virus. Unfortunately, it is not easy to retrieve each atom position in such a complex environment as a bacterial cell. One problem is the quantum mechanical uncertainty relation. This is a principal problem which cannot be overcome by any structure resolving technique. Here, only quantum mechanical calculations in combination with structural techniques can help. But mostly quantum mechanics is not the limit. A whole organism or even a single living cell is already too complex in order to obtain structural information on all molecules involved in metabolic processes. Here, one usually constructs model situations "in vitro" meaning: one just dissolves the molecules under investigation in a buffer solution mimicing pH and ionic strength of the cell environment. Small molecules (of less than 1000 g/mol in weight) usually have only few degrees of strutural freedom, so one is more interested in the "macromolecules" which often are polymers like starch (made out of suger units) or DNA/RNA consisting of Nucleotides often matched in base pairs and proteins which are polymers of 20 different aminoacids. Among those macromolecules used by living cells in their metabolism there are two classes which have a strong structure-function relationship and are therefore subjected to many structural studies: DNA or RNA and proteins. But it is already difficult to retrieve high (=atomic) resolution structural information from single molecules in solution. The NMR technique is one example but is often limited in the size of the macromolecule (see lecture A7). When it comes to scattering techniques, only very bright x-ray sources in the future (e. g. XFEL) might enable one to gain enough information out of many single protein molecules in solution, averaging many exposures and orientations[1].

Crystallography provides here the most elegant solution. In a crystal, the macromolecules are all oriented in a very regular way. This is why their scattering intensities add coherently and the contrast is increased. Both DNA/RNA and proteins form crystals. But by far the biggest number of structures recorded on biological macromolecules is on protein crystals using x-rays as a probe.

Apart from water, proteins are the most abundant molecules in living cells. Proteins fulfil numerous functions in the cell, for example enzymatic catalysis which enhances the speed with which molecules like fatty acids are synthesized. Proteins take care for transport and storage of important molecules like oxygen or they are involved in immunology, just to name a few of their functions [2]. In order to perform these functions proteins adopt a unique three dimensional structure with a carefully controlled mixture of flexibility and stiffness. For an understanding of their function knowledge of this three dimensional structure is a prerequisite. Ideally one would like to produce a movie where one can follow the functioning protein in action in slow motion with atomic resolution. In practice, there are techniques available which have a sufficient time resolved infrared spectroscopy. However, there are some promising attempts to do time resolved x-ray scattering to have both strutural resolution at the atomic level and time resolution [3]. But this technique of time resolved x-ray scattering is often limited to certain systems which undergo a photo-initiated process. On the other hand there are methods which provide full atomic resolution but with essentially no time resolution.

With those methods one often stops the functioning process of a protein under investigation in an intermediate state by trapping methods using inhibitor molecules which stop the catalytic process of the protein leaving it trapped it in a certain intermediate state. This lecture focusses on this type of methods since most of the structural information deposited in the protein data bank has been obtained with this "slow" x-ray crystallography technique. Due to its similarity to x-ray crystallography we will also discuss neutron protein crystallography which uses neutrons instead of x-ray photons as a probe.

The protein data bank (www.rcsb.org) is a well known source of structural information on proteins. Data from many different experimental techniques are entered in a standardized format, a .pdb file which essentially contains not only three dimensional coordinates of all observed atoms in a protein but also information on their mean square displacements. The number of stored entries exceeds 120534 as of 31st of October 2017.



Fig. 1: Entries of structual data with respect to the method used as found in the protein data base (www.rcsb.org) up to the 31st of October 2017. The total number of entries for each method is given in brackets. The circle on the left totals to 134656 entries of which 425 entries are taken out to form the second circle on the right. This lecture discussed the method of x-ray (dark blue section, left circle) and neutron (light blue section, right cercle) crystallography with 120534 and 58 entries respectively.

Among them x-ray crystallography has contributed more than 89.5 %. The next in line method with 11921 entries in the data bank is solution NMR spectroscopy. Electron microscopy as a method was used in more than 1776 entries and is the fastest expanding technique at the moment. A whole lecture (A2) in this spring school is dedicated to this technique. The remaining methods count mostly less than 100 entries each including neutron protein crystallography. Since the latter technique is represented by an instrument in the Jülich Centre for Neutron Science (JCNS) and because of its similarity to x-ray crystallography it will be given some space in this lecture. Finally two example case studies are discussed where both techniques give complimentary information.

2 The Physics of Macromolecular Crystallography

The following chapter will show how most of this structural information has been obtained. For both techniques x-ray and neutron protein crystallography a single crystal of the protein of interest is required since the scattering of one protein molecule is very weak. So, in general a crystal has to be grown, especially large ones in case of neutron crystallography since the neutron luminosity of modern sources is much smaller than the x-ray flux reached by synchrotron sources. To grow sufficiently large crystals is a big challenge in the case of many proteins, especially membrane proteins. Here, one has to adjust a large parameter set of protein concentration, pH condition, salt concentration, percipitant concentration and type just to name a few.



Fig. 2: Real space arrangement of myoglobin molecules in a crystal of space group $P2_1$ (on the left) versus diffraction pattern (right) of a myoglobin crystal.

The crystal then serves as a noiseless amplipier of the diffraction signal. But the arrangement of proteins in a crystals brings in another advantage, since the orientational averaging can be avoided, which is always present in the solution phase. Fig. 2: lef hand side shows the regular arrangement of myoglobin molecules in a crystal lattice. The unit cell of the monoclinic lattice (space group P2₁) is indicated by black lines. It bears two myoglobin molecules in one unit cell. The picture on the right shows a diffraction pattern recorded with the instrument BioDiff on a myoglobin crystal. The crystal is rotated by ca. 0.5° while recording one diffraction pattern. In order to map the reciprocal space completely one has to put the crystal in many different orientations into the beam and record diffraction patterns as mentioned above. Fortunately, crystal symmetry helps that some orientations are equivalent to each other and need not be recorded.

2.1 An X-ray macromolecular crystallography beamline

Synchrotron beamlines provide extremely high photon flux for x-ray crystallography. Due to the high demand from the structural biology community, often more than one macromolecular crystallography beamline is operated at a synchrotron. Those beamlines are optimized for special wavelengths and focal diameters. Here as an example the beamline BL14.2 is used to elucidate a typical x-ray protein crystallography experiment.

Fig. 3: left hand side shows a schematic view of the beam path of beamline BL14.2. The beam paths to the other beamlines BL14.1 and BL14.3 have been omitted for clarity. A supraconducting magnet-structure interacts with the electron orbit of the storage ring and creates the so called synchrotron radiation. This radiation is used as a white light X-ray source for the three beamlines. A double crystal monochromator is used to select a very narrow energy band (2 eV at 9 keV) from the broad spectrum of the X-ray source. The mechanics of the double crystal monochromator keeps the out-going beam path constant when changing the wavelength. Focussing mirrors and collimators in the beam path ensure an efficient photon transport from the source to the sample and a small beam size at the sample position of 150 μ m x 100 μ m (hxv, FWHM). The Rayonix MX-225 detector has a pixel size of 37 μ m. Without on chip binning one frame amounts to 6144 x 6144 pixels. The exposure time per frame is typically between 3 to 10 seconds.



Fig. 3: Schematic view of the beamline 14.2 of BESSY beginning with the magnet structure which is used as a white light x-ray photon source. On the right a picture taken from the experimental hutch of beamline 14.1 is shown. The beam enters from the right and the sample goniometer is mounted horizontally from the left. A cryostream sample environment to stabilize the sample temperature is discernible pointing towards the sample from the top right corner. Pictures kindly provided by Dr. Uwe Müller.

Typically, the sample crystals are kept at liquid nitrogen temperatures to avoid radiation damages. To record a full data set takes about 10-30 minutes. The largest diagonal of a typical protein crystal ranges between 10 to 500 μ m

2.2 A neutron macromolecular crystallography instrument

Since x-rays are scattered from the electrons in the crystal and neutrons from the nuclei, hydrogen atoms are hardly seen in x-ray crystallography experiments. Only at very high resolutions of 1 Å or less there is a chance to observe hydrogen atom positions. This resolution is often not within reach because of the crystal quality. Here neutron protein crystallography must be employed to retrieve the hydrogen atom positions. Moreover, neutron scattering can distinguish between different isotopes, especially between hydrogen and

deuterium. Whereas from x-ray crystallography the electron density in the unit cell of the crystal can be calculated, neutron protein crystallography yields the nuclear scattering length density, which is a signed quantity. In fact, the coherent scattering length of hydrogen is negative and the one from deuterium is positive (cf. Fig. 4: left).



Fig. 4: The table on the left lists scattering lengths of selected atoms of biological relevance. On the right there is a comparison of x-ray scattering cross section with scattering lengths from neutron scattering. The circles are scaled to match at the carbon atom.

A major drawback of the method neutron protein crystallography is the required crystal size. Due to the much smaller neutron flux as compared to x-ray flux the crystals required for a neutron crystallography study must be much larger as compared to x-ray crystallography. Here, often crystal diagonals of 1 mm and more have to be reached.

As an example of a neutron diffraction instrument optimised for protein crystallography the instrument BioDiff at the FRM II shall be introduced. It is a collaboration between the Forschungszentrum Jülich (FZJ) and the Forschungs-Neutronenquelle Heinz Maier-Leibnitz (FRM II).



Fig. 5: Schematic view of the BioDiff instrument (left) and a picture taken from a similar view point with the biological shielding removed (right).

Figure 10 shows a schematic view of the instrument from the top and a corresponding picture when the biological shielding has been removed. The neutron beam from the cold source of

the FRM II reactor enters from the right. By Bragg reflection from a pyrolytic graphite crystal (002-reflex) neutrons are taken out from the white neutron spectrum of the neutron guide NL1 and pass a first boron carbide adjustable slit and a velocity selector. The velocity selector acts as a $\lambda/2$ filter. Together with the pyrolytic graphite crystal it forms a monochromator with a $\Delta\lambda\lambda$ of ca. 2.5 %. Behind the velocity selector the beam passes a second variable slit and the main instrument shutter, named γ -shutter. Additionally, a boron carbide neutron shutter is placed directly after the monochromator crystal for a faster shutter operation. Before entering the detector drum of the image plate detector through a Zirconium window a collimator made out of two manually exchangeable boron carbide apertures with fixed diameters between 3 mm and 5 mm shape the beam to fit to the sample size. At present the sample is usually contained in a glass tube (either a thin walled capillary or a NMR-tube for larger crystals). It is fixed to a standard goniometer which is mounted upside-down from the sample stage on top of the instrument. After passing the sample the main neutron beam exits the detector drum through a second Zirconium window and hits finally the beam stop which consists of a cavity of 4 mm thick boron carbide plates surrounded by a 13 cm thick wall of lead shielding bricks. The cylindrical image plate detector is covering roughly half of the total 4π solid angle. It is 200 mm in diameter and 450 mm in height. It can be read out with three different resolutions of 125 μ m, 250 μ m and 500 μ m. As an alternative, one can lower the image plate detector and swing in a neutron sensitive scintillator which is imaged onto a CCD-chip. This CCD-camera set up serves as a second detector. In particular it is used for a fast alignment of the sample crystal with respect to the neutron beam.

With the image plate detector the diffraction pattern shown in Fig. 2: right hand side has been recorded. In fact, a complete crystallographic data set on a myoglobin crystal has been recorded allowing for the calculation of a nuclear scattering length density map. The exposure time was 17 minutes per frame and the crystal was rotated by 0.5° during exposure. 331 frames were recorded before the crystal was manually rotated by ca. 90° in the capillary and another set of 243 frames were recorded. Altogether ca. 9 days of beam time were necessary to record the complete data set. The achieved resolution with sufficient completeness was 1.7 Å. The required time to record this data set was much longer as the 30 minutes from x-ray diffraction.

2.3 Some general aspects of diffraction by a protein crystal

Having recorded a complete data set on a crystal some data treatment is necessary in order to calculate meaningful atom positions. Here only a brief introduction can be given more details can be found in text books [4,5].

Assuming a number of n atoms per unit cell the structure factor of a single unit cell can be written as:

$$F(\mathbf{S}) = \sum_{j=1}^{n} f_j \exp(2\pi i \mathbf{r}_j \mathbf{S})$$
(1)

Here \mathbf{r}_j denote the atom position of atom j and S is the scattering vector perpendicular to the plane which reflects the incident beam. In the prvious lecture A4 the scattering vector \mathbf{q} was defined. It relates to S with the following relation: $\mathbf{q} = 2\pi \mathbf{S}$. In crystallography it is just more convenient to use S instead of \mathbf{q} . f_j can be seen here either as the scattering length of atom j in the neutron scattering case or the atomic scattering factor in case of x-ray diffraction. One can generalize this approach by switching form the summation to an integration to yield:

$$F(\mathbf{S}) = \int_{unixell} \rho(\mathbf{r}) \exp(2\pi i \mathbf{r} \mathbf{S}) d^3 \mathbf{r}$$
(2)

where $\rho(\mathbf{r})$ is the electron density distribution or the scattering length density respectively. Since a crystal consists of AxBxC unit cells, the structure factor of the crystal can be composed as

$$F_{cryst.}(\mathbf{S}) = F(\mathbf{S}) \sum_{u=0}^{A} \exp(2\pi i u \mathbf{a} \mathbf{S}) \sum_{v=0}^{B} \exp(2\pi i v \mathbf{b} \mathbf{S}) \sum_{w=0}^{C} \exp(2\pi i w c \mathbf{S})$$
(3)

The vectors \mathbf{a} , \mathbf{b} and \mathbf{c} denote basis vectors of the unit cell. For an increasing number of unit cells the sums can be represented by delta functions leading to the Laue conditions for the structure factor being non-zero:

$$\mathbf{aS} = h, \mathbf{bS} = k, \mathbf{cS} = l \tag{4}$$

This means that one only gets constructive interference, when the scattering vector is perpendicular to planes in the crystal which can be denoted by the index vector $\mathbf{h} = hkl$. For this reason the diffraction pattern of a single crystals shows distinct peaks, the so called Bragg peaks. The Bragg law can be easily derived from equation 4. Figure 11 shows the Ewald sphere construction. It is a tool to construct the direction of the diffracted beam. The Ewald sphere has its origin at the position of the crystal. Its radius is the reciprocal wavelength used in the scattering experiment. The origin of the reciprocal space lattice is placed at the intersection of the sphere with the incident beam direction. Whenever the orientation of the reciprocal space last is on the Ewald sphere a diffracted beam results in the direction of line running from the centre of the Ewald sphere through that point.

When the crystal is rotated the reciprocal space rotates with it resulting in other lattice points to cause diffracted beams. In practice the incident beam is not strictly monochromatic but has a wavelength distribution which causes the Ewald sphere to be elongated to form a spherical shell of a certain thickness. This increases the number of diffracted beams observed. The beam divergence adds also to its thickness.

So, the positions of the diffracted beams on the detector only depend on the reciprocal lattice. The structure of the protein inside the unit cell is encoded in the amplitude and phase of the structure factor.



Fig. 6: Ewald sphere: On the left the schematic shows how to construct the Ewald sphere. On the right an Ewald sphere (golden colour) construction is shown in three dimensions. The blue square represents a flat two dimensional detector.

To obtain the electron density or the nuclear scattering length density one has to perform the inverse Fourier transformation:

$$\rho(\mathbf{r}) = \frac{1}{V} \sum_{\mathbf{h}} F(\mathbf{h}) \exp(-2\pi i \mathbf{r} \mathbf{h})$$
(5)

Here V is the volume of the unit cell. Unfortunately only the modulus squared of the structure factor is measured as intensity on the detector. The phase information is lost which is known as the phase problem of crystallography.

There are several solutions to the phase problem which are only applicable for the x-ray diffraction case:

- isomorphous replacement: Several crystals of the same crystal structure have to be available for this method. First a crystallographic data set is recorded on an untreated crystal. Then crystals are soaked in at least two different heavy atom salt solutions. In the best case, the different heavy atom ions occupy different regular positions in the unit cell. From these (at least) two crystallographic data sets recorded on the heavy atom treated crystals phase information can be retrieved which is then used to determine the phases of the data set of the untreated crystal.
- anomalous dispersion: Often it is possible to replace one distinct methionine amino acid with an artificial selenomethionine one. The selenium atom has a suitable absorption edge on which anomalous scattering can be performed by tuning the wavelength of the beamline to the anomalous regime. Crystallographic data sets are then recorded at different wavelengths from with

the phase information can be calculated. In some cases this approach can also be adopted for sulfur atoms present in naturally occuring cysteine residues.

 molecular replacement: From the primary structure one can search the protein data base (pdb) for proteins with a similar amino acid sequence. If one finds enough fragments which seem to be sufficiently homologous to the unknown structure one can use those fragments for the calculation of initial phases. In further refinement steps these phases can be improved further. Since the number of unique structures entered in the protein data base is growing this method is increasingly favoured over other methods.

The phase problem of the neutron data sets is solved by using the x-ray structure and the molecular replacement technique.

2.4 Model building and refinement

With the data treatment one has now arrived at a contour map $\rho(\mathbf{r})$ be it either a nuclear scattering length density or an electron density. Now the information on the primary structure of the protein is used and either manually or employing software first the backbone is coarsely fitted into the contour map. Then from this model new amplitudes and phases of the structure factor are calculated using eq. 1. The modulus squared of the structure factor is again compared with the data. One could now think of a least square based fitting procedure to find the optimum arrangement of the protein atoms in the unit cell. In practice however maximum likelihood and simulated annealing molecular dynamics simulations are used because those are superior to the least square approach in terms of overcoming local minima. In these molecular dynamics simulations a lot of stereochemical information is used as restraints for example known bond lengths of CC single bonds or bond angles. The agreement between model and observed contour map is often measured by calculating a so called R-factor:

$$R = \frac{\sum_{\mathbf{h}} \left\| F_{obs}(\mathbf{h}) \right| - \left| F_{calc}(\mathbf{h}) \right\|}{\sum_{\mathbf{h}} \left| F_{obs}(\mathbf{h}) \right|}$$
(6)

The index "obs" denotes the observed structure factors and the index "calc" the calculated structure factors from the model. The value of the R-factor lies in the limits between 0 and 1. A good agreement between model and measured data leads to an R-factor of about 0.2. R-factors of 0.5 and above are indicative for a random agreement between model and data.



Fig. 7: The side chain of the amino acid tryptophan no. 7 of myoglobin measured with neutron diffraction at different resolutions. The contour level of the shown nuclear map is $+1.5\sigma$ (blue) and -1.75σ (red). Exchanged (liable) hydrogen atoms (green) and C- (yellow), N- (blue) and O-atoms (red) appear on a positive contour level. Only not liable hydrogen atoms are seen on the negative contour level.

But even a good R-factor does not guarantee that the model fits to the data. In fact, it is possible in special cases to obtain a reasonably low R-factor when using the amino acid chain in the wrong direction as a model [6]. Here, Brünger et al. [7] have suggested to divide the measured Bragg reflections into two subset one working set denoted by "A" and one test set denoted by "T". With the working set the fitting procedure is performed, whereas the test set only serves to control the model quality by calculating the R_{fipe} factor.

$$R_{free} = \frac{\sum_{\mathbf{h} \in \mathcal{T}} \left\| F_{obs}(\mathbf{h}) \right| - \left| F_{calc}(\mathbf{h}) \right\|}{\sum_{\mathbf{h} \in \mathcal{T}} \left| F_{obs}(\mathbf{h}) \right|}$$
(7)

The test set usually consists of 5 to 10 % of all structure factors, uniformly distributed over the resolution range.

The R-factor and R_{free} factor should not differ too much from each other. In general, it is good practice to always look at the resulting model and its fit to the calculated map after each refinement step. Ramachandran plots can also be used to judge whether the amino acid backbone adopts a reasonable fold. With decreasing resolution (cf. Fig. 12) of the data it becomes more and more difficult to find the right orientations of side chains or even errors in the registry of the protein backbone can occur, whereby for example one amino acid is left out.



Fig. 8: Water network in the contact region between two myoglobin molecules in the crystal. In grey colour on the left amino acids 51 to 52 of one myoglobin molecule are shown. On the right amino acids 58 to 60 (from top to bottom) are depicted in green. In the centre of picture a sulfate ion SO_4^- is seen with the sulfur atom shown in yellow, the oxygen atoms shown in red. The deuterium atoms of the water molecules are coloured grey. The x-ray map is shown in magenta at a contour level of +2.7s. The nuclear map is shown at a contour level of -1.75s in red and at +2.3s in blue. Data taken from ref. [9]. (The picture is similar to the one shown in [8] or [10] since it is based on the same data.)

3 A first case study: Water network around myoglobin

This example is chosen since it nicely shows the interplay between x-ray and neutron crystallography. Myoglobin has been used quite frequently as an example throughout this lecture. Its function is to take over the oxygen molecules from the blood heme proteins in the red blood cells and to transport and store it within the muscle cells. Therefore its binding affinity to oxygene must be stronger than that of the hemoglobin. In order to perform the transport tasks myoglobin has to be highly soluble and movable within the context of a muscle cell. Let alone therefore it is interesting to study the water network surrounding of myoglobin. Since it is a joint neutron and x-ray diffraction study the crystal under investigation has been soaked in D_2O in order to exchange all liable hydrogen atoms with deuterium atoms.

Fig. 8: shows a contact region between two myoglobin molecules. In the centre a sulfate ion is observed. Here, in the nuclear map the central sulfur atom is hardly seen because of its small scattering length. The oxygen atoms of the sulfate ion and of the water molecules are readily observed by both techniques. The deuterium atoms of the water molecules are discernible only in the nuclear map. When the water molecule is fixed by hydrogen bonds, all three atoms can be observed. These water molecules exhibit a triangular shape in the nuclear map. In case

A5.13

atom is distributed over a large volume and is therefore averaged out. Those water molecules are denoted as "short ellipsoidal" by Chatake et al. [8]. The long ellipsoidal water molecules are fixed at both deuterium atoms but only the oxygen can rotate freely around the DD-axis. In the fourth case only the oxygen atom is observed. In this case the orientation of the water molecule is not fixed, only the oxygen atom is held in place.

This classification of water molecules helps to judge the flexibility of the water network around proteins. It can also be used to classify observed water molecules and their hydrogen bonding pattern in molecular crystals in general.

4 A second case study: Fighting Meticillin resistant bacteria

This example also shows the nice interplay between x-ray and neutron crystallography. It is chosen because of the relevance it has to our everyday life. The protein involved here is the β lactamase. It is produced by bacteria and partly secreted to their sourrounding environment in order to split a certain bond in the four-membered ring (= β -lactam motif) of a series of very common antibiotics. This is one of the mechanisms which render such bacteria resistant against this type of antibiotics. The antibiotic drug called meticillin also bears such a four membered ring. Meticillin resistant *Staphylococcus aureus* is assumed to possess this mechanism of resistance against β -lactam antibiotics. This is why it would be interesting to investigate the catalytic cylce of this protein in order to block it. This would result in the bacteria being affected by the antibiotica again.

The experiment was designed in a clever way. The catalytic cycle of the β -lactamase consists of an acylation step where a covalent aduct between protein and substrate (here the antibiotics molecule) is formed (Fig. 9: top part). Here, the CO double bond is split into a single bond and a tetrahedral intermediate between substrate and protein is produced. The second part is a deacylation step where the four-membered ring (the so-called β -lactam motif) is split and the product is released. The product cannot act as an antibiotic any more due to its different chemical struture.





Fig. 9: Catalytic cycle of a β-lactamase enzyme. (The Figure is taken from ref. [11]. This research was originally published in Journal of Biological Chemistry, Stephen J. Tomanicek et al. "Neutron and X-ray Crystal Structures of a Perdeuterated Enzyme Inhibitor Complex Reveal the Catalytic Proton Network of the Toho-1 β-Lactamase for the Acylation Reaction." © the American Society for Biochemistry and Molecular Biology.)

Due to its relevance to clinics the enzyme family of β -lactamases have been subjected to countless studies mostly by x-ray crystallography but the main question yet not addressed was the nature of the active base which takes the excess proton in the acylation step. It was not clear which part, which side chain of the protein takes over this role. But this was the key information for improving drugs to block this enzyme. Since a base is best detected when it takes a proton in the acylation step, the problem consisted of stopping the catalytic cycle in the acylation step and hunting for a proton which was not there in the ground state of the protein when no substrate was bound.



Fig. 10: X-ray (right) and neutron (left) data on the BZB aduct in the catalytic cycle of a βlactamase enzyme. (The Figure is taken from ref. [11]. This research was originally published in Journal of Biological Chemistry, Stephen J. Tomanicek et al. "Neutron and X-ray Crystal Structures of a Perdeuterated Enzyme Inhibitor Complex Reveal the Catalytic Proton Network of the Toho-1 β-Lactamase for the Acylation Reaction." © the American Society for Biochemistry and Molecular Biology.)

Here, a different ligand was used as the natural antibiotics: benzothiophene-2-boronic Acid (BZB). This lignand was known to stop the protein in the acylation step. In addition to that, this ligand was also known to mimic this covalent intermediate between substrate and protein. In addition to that, neutron protein crystallography was used to detect the proton unambigously. For the latter technique the protein was expressed in deuterated media such that a fully deuterated protein resulted. The BZB was also synthesized using deuterium instead of hydrogen atoms. In fact, even a special Boron isotope was used which shows less neutron absorption than the natural abundant boron which is known as a good neutron absorber. Fig. 10: left side shows the x-ray structure and Fig. 10: right side shows the neutron structure. The proton is clearly seen in the neutron structure and therefore the side chain Glutamate 166 could be identified as the base in the acylation step. With this knowledge better antibiotic drugs can be designed which bind this sidechain firmly blocking the protein's catalytic cycle in its first step.

Acknowledgement

The author would like to thank Christian Felder for software support. The author is especially grateful to Andreas Ostermann for supplying some of the Figures.

References

[1] F. Maia, T. Ekeberg, N. Timneanu, D. van der Spoel, and J. Hajdu, Physical Review E **80**, 031905 (2009).

[2] L. Stryer, *Biochemie* (Spektrum Akad. Verlag, Heidelberg Berlin New York, 1991), 4. edn.

[3] P. Fromme and J. C. H. Spence, Curr. Opin. Struct. Biol. 21, 509 (2011).

[4] J. Drenth, *Principles of Protein X-Ray Crystallography* (Springer Science+Business Media, LLC, New York, 2007).

[5] N. Niimura and A. Podjarny, *Neutron Protein Crystallography - Hydrogen, Protons, and Hydration in Bio-macromolecules* (Oxford University Press, Oxford, New York, 2011), IUCr Monographs on Crystallography, 25.

[6] G. J. Kleywegt and T. A. Jones, Structure 3, 535 (1995).

[7] A. T. Brunger, Nature 355, 472 (1992).

[8] T. Chatake, A. Ostermann, K. Kurihara, F. G. Parak, and N. Niimura, Proteins-Structure Function and Genetics **50**, 516 (2003).

[9] A. Ostermann, I. Tanaka, N. Engler, N. Niimura, and F. G. Parak, Biophys. Chem. 95, 183 (2002).

[10] T. Chatake, A. Ostermann, K. Kurihara, F. G. Parak, N. Mizuno, G. Voordouw, Y. Higuchi, I. Tanaka, and N. Niimura, Journal of Synchrotron Radiation **11**, 72 (2004).

[11] S. J. Tomanicek, R. F. Standaert, K. L. Weiss, A. Ostermann, T. E. Schrader, J. D. Ng, and L. Coates, J. Biol. Chem. **288**, 4715 (2013).

Recommended Textbooks

on X-ray crystallography:

J. Drenth, *Principles of Protein X-Ray Crystallography* (Springer Science+Business Media, LLC, New York, 2007)

on neutron protein crystallography:

N. Niimura, and A. Podjarny, *Neutron Protein Crystallography - Hydrogen, Protons, and Hydration in Bio-macromolecules* (Oxford University Press, Oxford, New York, 2011)

A 6 Neutron Imaging

S. Förster JCNS-1/ICS-1 Forschungszentrum Jülich GmbH

Contents

1	Introduction	.2
2	History	.2
3	Principles of Imaging	.2
4	Neutron Beams and Neutron Detection	.4
5	Data Evaluation and Tomography	.6
6	Application	.8
Refer	ences	.9

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Neutron imaging is a method of making an image with neutrons. It is a non-destructive technique that reveals the interior structure of a material. The image results from the neutron attenuation properties of the imaged object. Different materials attenuate neutrons differently. Neutron imaging can yield two-dimensional (radiography) or three-dimensional structural information (tomography).

In a typical neutron radiography experiment, the object being investigated is positioned in a neutron beam. A 2D position-sensitive detector placed on one side of the object records the radiation transmitted through the object. The shadow image thus produced gives information about the internal structure of the object.

Neutrons are unique probes for imaging because they specifically interact with the atomic nuclei. The interaction is weak for most elements and strong only for a few elements, the most prominent being hydrogen. Neutrons sense differences between isotopes of a given element and between elements with similar atomic numbers. Neutron sources at large-scale research facilities such as research reactors or spallation sources are necessary to produce sufficiently intense beams of free neutrons.

Neutron imaging is complementary to other non-destructive imaging methods, in particular X-ray imaging. X-rays are scattered and absorbed by the electron cloud of an atom and thus are less sensitive than neutrons to light elements such as hydrogen, lithium, boron, carbon, or nitrogen. The high penetration depth of neutrons allows to study large samples of some materials, with volumes of up to 100 cm3. [1]

2 History

Following the discovery of neutrons by Chadwick in 1932, the first demonstration of neutron radiography was made by Kallmann and Kuhn in the late 1930s and 1940s, initially using sources based on radioisotopes and later on accelerators. [2,3] With the advent of reactor sources in the late 1940s the experimental techniques were refined and optimized, and first applications were explored. The impetus for developing neutron radiography was the need for a technique that could be employed in applications that precluded the use of X-rays, e.g. in investigations on reactor fuel assemblies.

In the early days of neutron imaging the detection was done with film, using gadolinium as a converter, which was not quantitative. In the 1990's CCD cameras became available such that grey values could be measured and imaging became a quantitative method. With the rapid development of digital detectors with better spatial and temporal resolution, neutron imaging has much improved and diversified and is now used in research for a wide range of applications. Several facilities in Europe offer access to instruments for radiography and tomography, e.g. NEUTROGRAPH at the Institut Laue-Langevin (ILL) Grenoble, CONRAD at the Hahn-Meitner Institut (HM) Berlin, ANTARES at the Maier-Leibniz-Zentrum (MLZ) Berlin, and NEUTRA and ICON at the Paul-Scherrer-Institut (PSI) Villigen.

3 Principles of Imaging

Transmission imaging is a technique that uses incident penetrating radiation to investigate the internal structures and material compositions of optically opaque objects. The types of radiation range from particles (neutrons, electrons) to electromagnetic waves (X-rays, γ -rays), with specific selection depending upon the requirements of the application. The underlying

principle is the same for all: as the radiation passes through the object, it is attenuated by an amount that depends both on the thickness of the sample along the path taken by the radiation, and on the materials present along that path.



Fig. 1: Scheme of the experimental setup for neutron imaging equipped with a rotation stage for tomography. The collimator defines the neutron beam that passes through the sample to produce a two-dimensional image on the detector.

Radiography involves the two-dimensional detection of the transmitted beam intensity in a plane perpendicular to the beam propagation (Fig. 1). It produces a two-dimensional grey-scale image, which may be considered as a measure of the spatially varying integrated attenuation properties of the object, and holds information about sample thickness (including the presence of cracks and voids) and the chemical composition.

The transmission, T, is the ratio of the transmitted beam intensity, I, to the incident beam intensity, I_0

$$T = \frac{I}{I_0}$$

along a given propagation path s. The transmission is given by the basic law of radiation attenuation in matter

$$I = I_0 e^{\int \alpha ds}$$

where α is the local linear attenuation coefficient (in units of cm⁻¹). The attenuation coefficient is a material property and is given by

$$\alpha = \sigma \frac{\rho N_A}{M}$$

where σ is the total interaction cross-section (in units of cm²), ρ is the material density, N_A is Avogadro's number, and M is the molar mass of the material. Mass attenuation coefficients (in units of cm²/g) for the most common elements are given in Fig. 2.



Fig. 2: Mass attenuation coefficients for thermal neutrons (25 meV) and X-rays (100 keV) for elements as function of their atomic number. For neutrons there can be large variations of the attenuation coefficients for adjacent elements. After ref. [4].

The interaction cross-section for radiation passing through matter can generally be divided into two contributions, both of which depend on the radiation energy: absorption and scattering. Absorption is the retention of incident radiation without further transmission, reflection or scattering. In scattering events, local inhomogeneities cause the dispersal of radiation into a range of directions. Whereas absorption is an ideal attenuation process, scattering attenuates by altering the propagation direction of the radiation. As there is the possibility for scattered radiation to be redirected onto the detector at any arbitrary position, it can result in unavoidable background noise and the blurring of features and boundaries. The total neutron cross-section is given by

$$\sigma = \sigma_{\rm scattering} + \frac{v_{\rm thermal}}{v} \sigma_{\rm absorption}$$

where $\sigma_{\text{scattering}}$ is the sum of coherent and incoherent scattering effects, *v* the neutron velocity, and v_{thermal} and $\sigma_{\text{absorption}}$ are the velocity (~2200 m/s) and absorption cross-section for thermal neutrons, respectively. It follows that for cold neutrons having smaller velocities the absorption cross-section increases. Many metals (including aluminum and lead) have very small neutron cross-sections and are thus highly transparent. Neutron scattering is particularly strong for hydrogen. Several elements have absorption resonances, where the absorption cross-sections increase by several orders of magnitude for a narrow range of energies, lying within the thermal spectrum. These elements (including boron, gadolinium, cadmium, lithium) are thus ideal for use as tracer particles or in radiation shielding.

4 Neutron Beams and Neutron Detection

Radiation beams applied for imaging purposes are characterized in terms of four properties: the flux distribution, the geometry, the divergence, and the energy spectrum. These properties may be fine-tuned by introducing various pieces of equipment along the flight path. [3]

Beam size and flux distribution. Neutron beams used for neutron imaging can have dimensions as large as 220x220 mm² (e.g. Neutrograph at ILL, Grenoble). Larger objects can be visualized by aligning and splicing exposures at multiple positions. Ideally, the beam has a homogeneous flux distribution.

Beam geometry. The beam geometry describes the spatial distribution of the radiation leaving the source. The two most commonly employed beam geometries are the cone beam and the parallel beam. Cone beams are produced naturally as radiation emanates isotropically from a point source. Parallel beams can be approximated by taking a relatively small portion of a cone beam far from the source. The advantage of parallel beams is that tomographic imaging only requires rotation about the sample over angles $0...\pi$, since any projection P_{θ} obtained at a rotation angle θ will be an exact mirror image of the projection $P_{\theta+\pi}$.



Fig. 3. Collimation of a neutron beam using an aperture of width *D*. Neutron beams passing through the same position on the sample, but coming from different directions will arrive on different positions on the detector, defining the resolution d.

Collimation. Neutrons arriving from the source will be travel traveling in many different directions. To produce a good image, neutrons beams need to have a fairly uniform direction, which is accomplished using a collimator. It has an aperture, an opening that will allow neutrons to pass through. In the collimator neutron absorption materials (e.g. boron) absorb neutrons that are not traveling the length of the collimator in the desired direction. A shorter collimation system or larger aperture will produce a more intense neutron beam but the neutrons will be traveling at a wider range of angles, while a longer collimator or a smaller aperture will produce more uniformity in the direction of travel of the neutrons, but significantly fewer neutrons will pass and a longer exposure time will result. As collimation determines the resolution, collimators have become the key factor in neutron imaging instruments.

Beam divergence. The divergence of the beam determines the extent to which it spreads out as it propagates. It has a direct bearing on the spatial resolution in the resulting image. For imaging instruments the most useful figure of merit of the divergence is the ratio of the beam's effective length L to its diameter D (Fig. 3). From this value it is possible to gauge the maximum achievable spatial resolution, d, using simple geometry

$$\frac{L}{D} = \frac{l}{d}$$

where l is the distance between the sample and the detector. A higher L/D-value equates to a more parallel beam and returns sharper images.

In principle, L/D could be increased arbitrarily by increasing L and/or decreasing D. However, a small beam diameter would decrease the field of view, and an increasing the beam length, apart from potential impracticalities, would be detrimental to the available flux, which is inversely proportional to the square of the distance from the source. The range of L/D at neutron imaging facilities range from L/D = 100 - 800. If the detector is placed at l=1cm behind the object, the achievable resolution is in the range of $d = 100 - 12.5 \,\mu\text{m}$.

Energy spectrum. The intensity at lower energies, where neutrons are approaching thermodynamic equilibrium with the moderator, is described by a broad Maxwell-Boltzmann distribution, peaking at 25 meV (300 K). Thermal neutrons are well suited for imaging purposes. Cold neutrons are in some cases even more so because they have higher interaction cross-sections, increasing the sensitivity to the materials, and thus yield better contrast and higher detection efficiency.

Radiography and tomography require two-dimensional position-sensitive detection. The neutrons are detected in two steps: First, the neutron is captured by a strong neutron absorber, e.g. ³He, ¹⁰B, ⁶Li, or natural Gd. This capture is accompanied by the emission of charged particles or X-rays. This secondary radiation travels for a distance and then interacts with a scintillating material, thus producing visible light which can be detected by conventional position-sensitive detectors, for instance CCD cameras. The mean free path of the secondary radiation in the scintillating material (usually ZnS), gives rise to an uncertainty and limits spatial resolution to currently about 20 μ m.



Fig. 4: *Typical detection system for neutron imaging consisting of a scintillator screen which converts neutrons into photons, a mirror and a CCD-camera*

The scintillator is a homogeneous plastic sheet containing a conversion material, e.g. ⁶LiF, ¹⁰B, ¹¹³Cd, or ¹⁵⁵Gd/¹⁵⁷Gd, and a phosphor, e.g. ZnS:Cu,Al,Au. The conversion materials absorb neutrons to form energetic particles (α ,e⁻) or radiation (γ), which cause ionization in the phosphor. Upon ionic relaxation, phosphorescent photons in the visible wavelength range are emitted, which are detected by a CCD camera. A lens can be used to focus the light onto the CCD chip, which allows the field of view to be enlarged.

5 Data Evaluation and Tomography

High spatial resolution and contrast are two important objectives of radiographic imaging. High quality imaging thus relies on the precision and accuracy with which the beam intensity can be measured, as limited by systematic and statistical errors and the degree to which they can be corrected or minimized.

Image normalization. Raw data always require two corrections: an offset related to fluctuations in the electronics of the CCD camera, and a normalization related to inhomogeneities in the beam intensity and in the detection efficiency of the scintillator and the CCD chip. The magnitude of the dark noise offset is found by measuring an image with no incident photons, i.e. an image where the neutron beam shutters and the internal camera shutter are closed. Normalization is to the unperturbed or open beam, i.e. the intensity of the beam with nothing in the field of view. The beam attenuation, A, is then given by

$$A = \int \alpha ds = -ln \frac{I}{I_0} = -ln \left(\frac{i \cdot b}{o \cdot b}\right)$$

where i is the sample image intensity, o is the open beam image intensity, and b is the dark image intensity. If the sample is placed in a sample holder or located in a different material, this is normalized as well and subtracted.

Time resolution. A single radiographic image can require as little as 1 ms exposure, although 50 - 250 ms is more usual. As an example, a ~15 µm thick polyethylene sheet is detectable with an exposure time of ~1.5 s.

Contrast. In some instances it is possible to greatly improve the contrast by selectively altering the material present in a region of interest – either by introducing additional material or by completely replacing the material with a new one. This contrast agent is usually chosen to be either a strong a weak absorber (depending on the other materials present in the sample) such that the region of interest will contrast well with neighbouring regions. For neutrons, good examples of the former case are boron, water-soluble gadolinium salts, liquid gadolinium, and helium-3. The latter case may be particularly desirable for example with water, which may be exchanged with heavy water; water is a strong scatterer, while heavy water is not.

Tomography. Whereas radiography results in a two-dimensional image of the sample studied, tomography allows to visualize samples in three dimensions. For neutron tomography, the sample is rotated stepwise around a fixed (mostly vertical) axis, whereby for each rotation angle a projection image is recorded by a position sensitive detector. A 3D representation of the volume of the object can be reconstructed using a mathematical algorithm.

"Computed axial tomography" (CT) was first demonstrated in 1972 by Hounsfield. It is based on a mathematical framework formulated by Radon in 1917, and provides a system for reconstructing the three-dimensional properties of a sample from a set of many radiographs taken at a series of angles.



Fig. 5: Sequence of data normalization and transformation for the tomographic reconstruction of a plant. From ref. [3]

Today, one of the fastest, most efficient, and most widely-used mathematical algorithms is filtered back-projection. All projections are first rearranged into a new set of images such that the *n*-th pixel row of each projection is now stored sequentially in an individual image, the so-called sinogram (Fig. 4). A sinogram contains all the attenuation information for all angles for a row of pixels. An implementation of the inverse two-dimensional Radon transform is then used to calculate a cross-sectional slice from each sinogram. Each reconstructed slice lies in a plane perpendicular to the axis of rotation around the sample, and is, by definition, exactly one pixel thick. Collecting all slices into an image stack represents the three-dimensional attenuation properties of the object; each pixel in every slice is a finite volume, a so-called voxel, whose value represents the attenuation properties of the corresponding volume in the sample. This image stack can be rendered and manipulated with software to highlight specific volumes, surfaces, planes or contours, or to view the selected features from any perspective.

6 Application

Neutron imaging is used routinely to highlight light materials such as hydrogeneous substances with high contrast in engine parts or in hydrogen storage tanks and fuel cells. It allows to visualize the movement of fluids, such as oil or water, in large metal objects or sandstone. It is also used to analyze archaeological artefacts or pieces of art.

Stroboscopic imaging can be used to investigate fast, periodic events. Fig. 6 shows neutron radiography images from the inside of a running motorcycle engine.



Fig. 6: Radiography images of a motorcycle engine at 1200 rpm recorded with an exposure time of 0.1 milliseconds. From ref. [5]

Neutron imaging is frequently used to investigate water penetration or take-up in sandstone or plants (Fig. 7).



Fig. 7: Neutron radiography image of a tomato seedling taking up heavy water (D_2O) from the roots. From ref. [6,7].

References

[1] H. Kallman, E. Kuhn, Research 1 (1947), 254.

[2] L. Santodonato, J. Bilheux, B. Bailey, H. Bilheux, Introduction to Neutron Imaging, Oak Ridge National Laboratory, 2014

[3] M. N. Dawson, Applications of Neutron Radiography & Tomography, Ph.D. thesis, University of Leeds, 2008

[4] M. Strobl et al., J. Phys. D: Appl. Phys. 42 (2009) 243001.

[5] Neutron Imaging - How neutrons create pictures, Paul Scherrer Institute, 2007

[6] U. Matsushima *et al.*, Application potential of cold neutron radiography in plant science research, *J. Appl. Botany and Food Quality* 82 (2008), 90-98.

[7] N. Kardjilov, Oxford School on Neutron Imaging, 2017

A 7 Protein NMR Spectroscopy in Solution

P. Neudecker

Physikalische Biologie, Heinrich-Heine-Universität Düsseldorf ICS-6 (Structural Biochemistry) Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Intr	oduction	2
2	Fun	damentals of Protein NMR Spectroscopy	2
	2.1	Nuclear magnetic resonance	2
	2.2	One- and multi-dimensional pulsed FT NMR spectroscopy	5
	2.3	Local magnetic fields and interactions	8
	2.4	Chemical shift	8
	2.5	Indirect (scalar) and direct (dipolar) coupling	9
3	Pro	tein Dynamics	11
	3.1	Time scales and NMR spectral parameters	11
	3.2	Examples	13
Ref	erence	·S	15

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

On the atomic level the biological function of a protein is usually determined by the position of key chemical groups in its static three-dimensional structure and by the dynamics of the three-dimensional conformation along the reaction coordinate. Although nuclear magnetic resonance (NMR) spectroscopy has been known since 1946 [1] [2], the spectral complexity of biomacromolecules required the development of multi-dimensional techniques with sufficient resolution and of computers with sufficient computing power to allow the first determination of the high-resolution three-dimensional structure of a protein in 1985 [3]. As of January 2, 2018, 10,599 out of a total of 126,802 protein and/or nucleic acid coordinate entries in the Protein Data Bank (PDB; https://www.wwpdb.org/) have been determined by NMR spectroscopy. While X-ray crystallography continues to be the most widely used method for experimental structural biology and electron microscopy (EM) is increasing in popularity, NMR spectroscopy offers several distinct advantages: First, it is directly applicable to proteins in aqueous solution, which is the natural environment for most proteins (note that NMR spectroscopy of proteins in the solid state is also possible but beyond the scope of this article). Second, the various NMR spectral parameters are highly sensitive to dynamics on almost all biologically relevant time scales [4], which makes NMR spectroscopy an ideal tool to study biochemical processes as diverse as protein folding, ligand binding, allosteric regulation, enzyme catalysis etc. with a unique combination of high spatial and temporal resolution under near-physiological solution conditions.

To take full advantage of the arsenal of modern NMR spectroscopic techniques it is typically necessary to produce several hundred microliters of highly pure (> 95%) aqueous sample solutions of proteins enriched in rare isotopes (¹³C, ¹⁵N, and sometimes ²H) at relatively high protein concentrations (up to about 1 mM). In practice, this is most commonly achieved by high-yield biosynthetic protein expression in bacteria or yeast grown in suitably isotopeenriched media, followed by two to three orthogonal steps of chromatographic purification and screening of buffer conditions for optimal stability and solubility of the protein of interest. Due to the practical limitations imposed by increasing spectral overlap and deteriorating relaxation properties with the size and hydrodynamic radius of the molecule under study, high-resolution structure determination by solution NMR spectroscopy is best suited for proteins with a molecular weight of about 30 kDa or less. These limitations notwithstanding, solution NMR spectroscopy can still yield valuable mechanistic insights into the action of supra-molecular machines with molecular weights of several hundred kilodaltons [5]. In this lecture, the fundamentals of protein NMR spectroscopy in solution are introduced and illustrated based on practical examples, with a special emphasis on recent developments to study conformational dynamics in protein folding and function. For further study, several excellent textbooks detailing the fundamentals of NMR structural biology are available by now, e. g. [6] [7] [8] [9] [10].

2 Fundamentals of Protein NMR Spectroscopy

2.1 Nuclear magnetic resonance

In addition to rest mass and electric charge most elementary particles possess another fundamental property, their intrinsic angular momentum, which is called spin. The spin – like
all angular momenta – is a vectorial quantity, which exists only in quantized units (quanta) according to the laws of quantum mechanics. The magnitude of the spin vector is given by

$$\left|\vec{I}\right| = \hbar \sqrt{I(I+1)} \tag{1}$$

where $\hbar = h/(2\pi)$ is the Planck constant divided by 2π and *I* is a spin quantum number that is characteristic for a particular elementary particle. The spin quantum number *I* of a particle can be half-integer (so-called "fermions") or integer (so-called "bosons"). The spin component in one of the three directions in space (by convention chosen to be the *z* axis) can only adopt the following values:

$$I_z = m_l \hbar \tag{2}$$

where $m_I = -I$, -I+1, ..., 0, ..., I-1, I (2I+1 values in total) is the magnetic quantum number. Matter is composed of quarks and electrons, which are fermions with spin I = 1/2. In this case, the spin has two eigenstates characterized by the magnetic quantum numbers $m_I = +1/2$ (spin up, often denoted as $|\uparrow\rangle$ or $|\alpha\rangle$) and $m_I = -1/2$ (spin down, $|\downarrow\rangle$, $|\beta\rangle$).

In classical electrodynamics an electrically charged particle on a circular orbit (e. g. electrons in a solenoid) generates a magnetic dipole field, i. e. the angular momentum of an electrically charged particle is associated with a magnetic (dipole) moment. The potential energy of a magnetic dipole moment $\vec{\mu}$ in a magnetic field \vec{B} is

$$E_{\text{not}} = -\vec{\mu} \bullet \vec{B} \tag{3}$$

The magnetic field exerts a torque on the magnetic moment:

$$\vec{N} = \vec{\mu} \times \vec{B} \tag{4}$$

In quantum mechanics the spin of an electrically charged particle is also associated with a magnetic (dipole) moment:

$$\vec{\mu} = \gamma \vec{I} \tag{5}$$

 γ is called the gyromagnetic ratio of the particle. Both the proton and the neutron have a spin of I = 1/2 with an associated magnetic moment because they both consist of three electrically charged quarks. As a result, the atomic nuclei of most isotopes also possess a spin with an associated magnetic dipole moment. A notable exception are nuclei with an even number of protons and an even number of neutrons that are all paired to yield a net nuclear spin of 0, most importantly the most abundant carbon isotope ¹²C, which is unfortunately silent in NMR spectroscopy. For reasons of symmetry, nuclei with I = 1/2 do not show a quadrupole moment, which results in much more favorable relaxation properties and hence narrower linewidth (see below) due to the absence of the strong interaction between the nuclear quadrupole moment and the electric field gradient at the nucleus. Spin-1/2 nuclei such as ¹H, ¹³C, ¹⁵N, and ³¹P are therefore particularly well suited for high-resolution NMR spectroscopy of biomacromolecules. The NMR properties of the most relevant stable isotopes for biomolecular NMR are compiled in Tab. 1. Note that there is no stable oxygen isotope suitable for high-resolution NMR spectroscopy.

For NMR spectroscopy, the sample tube is inserted into a probe head located in the center of the bore of a superconducting magnet that generates a strong static magnetic field $\vec{B}_0 = B_0 \vec{e}_z$, which is extremely homogeneous over the active sample volume, extremely stable over time, and which defines the longitudinal (z) axis. Due to (2) and (3), the energy levels of the two eigenstates for I = 1/2 will then split up:

$$E_{pot} = -\vec{\mu} \bullet \vec{B} = -\gamma \vec{I} \bullet \vec{B}_0 = -\gamma B_0 I_z = -m_I \gamma \hbar B_0 = \mp \frac{1}{2} \gamma \hbar B_0$$
(6)

Isotope	Ι	$\frac{\gamma_I}{[10^7 \text{ rad}/(\text{T s})]}$	$ \omega_0 /(2\pi)$ at $B_0 = 14.09$ T [MHz]	Natural Abundance	<i>Υι/</i> Υ <i>Η</i>
$^{1}\mathrm{H}$	1/2	26.752	600.0	99.99%	1 (<i>per def.</i>)
² H	1	4.107	92.1	0.01%	0.153506088
¹² C	0	-	-	98.9%	-
¹³ C	1/2	6.728	150.9	1.1%	0.251449530
¹⁴ N	1	1.934	43.4	99.6%	
¹⁵ N	1/2	-2.713	60.8	0.4%	0.101329118
¹⁶ O	0	-	-	99.76%	-
¹⁷ O	5/2	-3.628	81.4	0.04%	
¹⁸ O	0	-	-	0.20%	-
¹⁹ F	1/2	25.18	564.7	100.0%	
³¹ P	1/2	10.839	243.1	100.0%	0.404808636

Tab. 1: *NMR* properties of selected stable isotopes [9] and IUPAC-IUB recommended *chemical shift ratios* γ_{I}/γ_{H} [11].

The resulting energy difference between the spin states,

$$\Delta E = \gamma \hbar B_0 = \hbar \left| \omega_0 \right| \tag{7}$$

is proportional to the static magnetic field B_0 . Nuclear magnetic resonance is therefore observed at a frequency of

$$\omega_0 = -\gamma B_0 \tag{8}$$

the so-called Larmor frequency. For commercially available superconducting NMR magnets of presently up to about 23 T this frequency is in the radiofrequency (RF) range (see Tab. 1).

According to the Boltzmann distribution,

$$\frac{N_{\beta}}{N_{\alpha}} = e^{-\frac{\Delta E}{k_{B}T}} \tag{9}$$

there is a population difference between the two spin states in equilibrium and the vector sum over all the magnetic moments in the sample results in a net longitudinal macroscopic magnetization $\vec{M}_0 = M_0 \vec{e}_z$. Unfortunately, this macroscopic equilibrium magnetization represents only about 0.01% of all the spins in the sample because the energy difference (7) is about four orders of magnitude smaller than the thermal energy k_BT at room temperature, which greatly reduces the sensitivity of solution NMR spectroscopy and results in the aforementioned requirement of high sample concentrations. Unlike the longitudinal component (2), the components of the spin in the transverse (x, y) plane are not eigenstates and transverse macroscopic magnetization can only be generated by creating phase coherence of the spins. For reasons of symmetry, all transverse orientations are equally probable in equilibrium and consequently there is no phase coherence and no net transverse magnetization. As a result, the sample shows a macroscopic magnetization parallel to the static magnetic field in equilibrium.

Upon perturbation the magnetization will relax back to equilibrium in an exponential fashion with characteristic time constants. For proteins, the longitudinal relaxation times T_1

characterizing the buildup of the longitudinal Boltzmann magnetization are typically on a time scale of seconds, thereby limiting the repetition rate with which NMR experiments can be repeated in order to obtain better signal to noise ratio. The transverse relaxation times T_2 characterizing the loss of any phase coherence in the transverse plane are typically on a time scale of tens to a few hundreds of milliseconds. As a result, transverse phase coherence is readily created by coherently manipulating a thermodynamic ensemble of spins with microsecond radiofrequency pulses applied at the Larmor frequency (8) with magnetic field components in the transverse direction, which are generated by transverse transmitter/receiver coils in suitable tuned RF circuits in the probe head. It is this inherently coherent nature which makes pulsed NMR experiments can no longer be fully understood by considering spins individually but require thermodynamic quantum mechanical analyis of an ensemble of coherent spins using the density operator formalism or shorthand notations thereof such as the product operator formalism [9], which is beyond the scope of this article.

2.2 One- and multi-dimensional pulsed FT NMR spectroscopy

A simplified description is the so-called vector model of NMR spectroscopy. Because the nuclear spin is an angular momentum on which the magnetic field exerts a torque according to equation (4), the Bloch equations governing the motion of nuclear magnetic moments are analogous to the equations of motion of a spinning top:

$$\frac{d}{dt}\vec{\mu}(t) = \gamma \frac{d}{dt}\vec{I}(t) = \gamma \vec{N} = \gamma \vec{\mu}(t) \times \vec{B} = -\gamma \vec{B} \times \vec{\mu}(t)$$
(10)

Or, for macroscopic magnetization:

$$\frac{d}{dt}\vec{M}(t) = -\gamma \vec{B} \times \vec{M}(t) \tag{11}$$

In analogy to a spinning top, what these equations describe is that any magnetization that is not parallel to the magnetic field will undergo a circular precession about the direction of the magnetic field with an angular frequency of $\omega = -\gamma B$. In the absence of any RF pulses, $\vec{B} = \vec{B}_0 = B_0 \vec{e}_z$ and hence any transverse magnetization precesses about the z axis with the Larmor frequency $\omega_0 = -\gamma B_0$. By contrast, the magnetic field of a linearly polarized transverse RF pulse, $\vec{B}(t) = \vec{B}_1(t) = 2B_1\vec{e}_x \cos(\omega_{RF}t)$, is explicitly time dependent. This time dependence is removed by decomposing the linearly polarized RF wave into two counterrotating circularly polarized RF waves followed by transformation into a coordinate frame that rotates with the radio frequency $\omega_{RF} \approx \omega_0$ about the z axis. In the rotating frame (rf), one of the circularly polarized components becomes stationary, $\overline{B_1^{rf}}(t) = B_1 \overline{e_x^{rf}}$, and it can be shown that the effect of the counter-rotating component can be neglected [9]. If B_1 is sufficiently strong that off-resonance effects can be neglected, its effect in the rotating frame is therefore a precession about the transverse direction along which the pulse is applied (the phase of the pulse, in this case x) with angular frequency $\omega_1 = -\gamma B_1$, and a pulse of duration pw rotates the magnetization about the pulse axis by an angle of $-\gamma B_{1\times}pw$. For example, a pulse of $pw = 10.0 \,\mu\text{s}$ and phase +x applied at a field strength of $\gamma B_1/(2\pi) = 25.0 \,\text{kHz}$ will rotate z magnetization by -90° , so that it will end up along the y axis. Similarly, a 180° pulse is obtained by doubling the duration to $pw = 20.0 \ \mu s$ (or, alternatively, by doubling the field strength while keeping pw constant) and will rotate z magnetization to -z. In other words, as long as the spins maintain sufficient phase coherence (limited by the transverse relaxation time T_2) we can freely rotate the magnetization along a trajectory of our choosing by designing a suitable sequence of RF pulses; during the delays between the RF pulses any transverse magnetization present at this point in the sequence precesses about the *z* axis with the Larmor frequency ω_0 in the lab frame and with the difference frequency $\Omega = \omega_0 - \omega_{RF}$ in the rotating frame.

Instead of being limited to continuous wave (cw) NMR spectroscopy, in which either the external magnetic field B_0 or the radio frequency is swept to obtain the nuclear magnetic resonance spectrum, the long lifetime of phase coherences of nuclear spins thus allows us to obtain the entire NMR spectrum in a single pulsed Fourier transform (FT) NMR experiment. To this end, the longitudinal equilibrium magnetization is rotated into the transverse plane using a 90° pulse. The ensuing precession of the transverse magnetization about the z axis with the Larmor frequency ω_0 induces an oscillating radiofrequency voltage in the transmitter/receiver coils in the probe head, which is amplified, digitally sampled using phasesensitive detectors (PSDs) and analog-to-digital converters (ADCs), and stored in a computer. For historic reasons this oscillating signal as a function of time is called the free induction decay (FID). The NMR spectrum is readily obtained by spectral decomposition of the FID via Fourier transformation (FT). In addition to the multiplex advantage of being able to record the entire spectrum in a single measurement that is faster than the (obsolete) cw detection of NMR, pulsed FT NMR also offers significantly more control and enables us to custom-tailor multi-pulse sequences that allow efficient solvent and artifact suppression, isolation of a single spectral parameter of interest from the complications arising from all the other spectral parameters in dedicated experiments, and recording of multi-dimensional NMR experiments in order to resolve the severe spectral overlap for large biomacromolecules.



Fig. 1: 1D ¹H NMR spectrum of 1.0 mM [U-¹⁵N] Fyn SH3 domain mutant A39V/N53P/V55L in 50 mM sodium phosphate, 0.2 mM EDTA, 0.05% (w/v) NaN₃, 10% D₂O (pH 7.0) recorded at 800 MHz and a temperature of 20.0°C. The ¹H₂O resonance at 4.82 ppm was suppressed by excitation sculpting [12] with water flip-back [13].



Fig. 2: [¹H,¹⁵N] HSQC spectrum of 1.0 mM [U-¹⁵N] Fyn SH3 domain mutant A39V/N53P/V55L in 50 mM sodium phosphate, 0.2 mM EDTA, 0.05% (w/v) NaN₃, 10% D₂O (pH 7.0) recorded at 800 MHz and a temperature of 20.0°C. Amide proton resonances are labeled according to their residue numbers [14] (Thr2 corresponds to Thr84 of full-length Gallus gallus Fyn, residues -3 to 1 are from a cloning artifact). The Hε resonances of Arg13, Arg40, and Arg60 (blue) are aliased in the ¹⁵N dimension.

Even small protein domains such as SH3 domains contain several hundreds of protons, whose resonances cannot normally be resolved in a one-dimensional (1D) NMR spectrum (Fig. 1). As suggested by J. Jeener in 1971, two-dimensional (2D) NMR spectra can be recorded using pulse sequences consisting of the following four building blocks:

- (a) A preparation period consisting of a single 90° pulse or a multi-pulse sequence that generates transverse magnetization of the first spin of interest,
- (b) an evolution time consisting of a delay, t_1 , during which the first spin of interest is allowed to precess,
- (c) a mixing pulse sequence that transfers transverse magnetization from the first spin of interest to the second spin of interest, and
- (d) a detection period, t_2 , during which the precession (FID) of the second spin of interest is recorded.

The FID of the second spin is recorded directly in such a sequence and the corresponding 1D spectrum is obtained by FT along t_2 . However, as described in more detail in the textbook by

J. Keeler [6], the precession of the first spin during the evolution time t_1 determines the initial conditions for the mixing sequence and thereby modulates the amplitude of the NMR spectrum of the second spin. A second dimension is obtained indirectly by recording many such 1D spectra while systematically varying the evolution period t_1 , followed by spectral decomposition of the precession frequency of the first spin that is encoded in the amplitude modulation of the set of 1D spectra by FT along t_1 . Because the NMR spectrum of the backbone amide groups is particularly well dispersed in folded proteins, the 2D [¹H,¹⁵N] heteronuclear single quantum coherence (HSQC) spectrum, in which the ¹⁵N resonance frequency is measured in the indirect dimension and the ¹H resonance frequency of the attached amide proton in the direct dimension, is commonly used as a "fingerprint" spectrum (Fig. 2). If further resolution is desired, the strategy outlined above is readily extended to record 3D or 4D experiments by introducing additional evolution & mixing periods.

2.3 Local magnetic fields and interactions

If the nuclear spins were experiencing only the static external magnetic field B_0 , all ¹H (and similarly ¹³C or ¹⁵N) nuclei would resonate at exactly the same frequency, the Larmor frequency given by equation (8). Figs. 1 and 2 clearly show that this is not the case; there are additional interactions with other spins and with the electrons in the molecule, which give rise to locally varying magnetic fields and which are a treasure trove of information about molecular structure and dynamics.

2.4 Chemical shift

The electron cloud of the molecule under study causes a tiny (parts per million, ppm) diamagnetic screening of the external magnetic field at the location of the nucleus:

$$B = B_0 \left(1 - \sigma \right) \tag{12}$$

This screening reflects the chemical environment and is associated with a shift of the resonance frequency,

$$\omega = -\gamma B_0 (1 - \sigma) \tag{13}$$

the so-called chemical shift. Note that the chemical shift is also dependent on orientation and generally has to be described by a tensorial quantity; in isotropic solution, however, the chemical shift anisotropy (CSA) is averaged to 0 by rotational diffusion and only the isotropic component σ of the tensor remains, which is a scalar quantity. As seen in Fig. 1, the spectral width of a typical protein ¹H NMR spectrum is about 10 ppm = 10^{-5} of the external magnetic field. To resolve two ¹H nuclei with a chemical shift difference of 0.01 ppm = 10^{-8} it is therefore imperative that B_0 is homogeneous to eight decimals over the active sample volume and remains constant to eight decimals over the duration of the measurement. This requirement is achieved via two correction mechanisms built into commercially available NMR spectrometers, the shim coil system to create spatial homogeneity and the fieldfrequency lock operating on the ²H channel to compensate for the field drift; note that operation of the lock channel necessitates the use of a partially deuterated solvent, for example 10% D₂O. Calculation of the chemical shift σ from the measured resonance frequency ω according to (13) also requires knowledge of the absolute value of B_0 to an accuracy of eight decimals. Because only NMR itself offers this level of accuracy, the magnetic field is determined from the resonance frequency ω_{ref} measured for a reference compound and the chemical shift is conventionally reported as

$$\delta = \frac{\omega - \omega_{ref}}{\omega_{ref}} \approx \sigma_{ref} - \sigma \tag{14}$$

The IUPAC-recommended reference standard for protein ¹H NMR spectroscopy is the methyl group resonance of 2,2-dimethylsilapentane-5-sulfonic acid (DSS), ω_{DSS} ; reference frequencies for ²H, ¹³C, ¹⁵N, ³¹P are obtained from the measured ω_{DSS} by multiplication with the indirect chemical shift referencing ratios γ_{I}/γ_{H} reported in Tab. 1 [11]. The chemical shift scale δ defined by (14) offers the key advantage of being independent of B_0 , thus allowing direct comparison of spectra recorded on different NMR spectrometers.

Protein chemical shifts are highly sensitive to a wide variety of chemical parameters, including all levels of protein structure - primary, secondary, tertiary, and quaternary structure [15]. Identification of regular secondary structure elements such as α -helices and β -strands from chemical shifts is relatively straightfoward using approaches such as the chemical shift index (CSI) [15] or TALOS-N [16]. Determination of high-resolution tertiary structures from chemical data is significantly more challenging; in recent years, different computational strategies such as CS-ROSETTA [17] and CamShift [18] were developed for this purpose. Chemical shifts are also sensitive reporters of local backbone dynamics [19]. Furthermore, NMR chemical shifts are highly sensitive to ligand binding and other changes in chemical environment and are therefore routinely used to detect and quantify protein-protein and protein-ligand interactions via titration experiments, often even allowing determination of the complex structure by molecular docking guided by the observed chemical shift perturbations [20].

2.5 Indirect (scalar) and direct (dipolar) coupling

Via a combination of the hyperfine coupling (the magnetic interaction between the nuclear and the electronic spins) and the exchange interaction (the electric repulsion between the electrons in a molecular orbital) the magnetic field experienced by a nuclear spin depends on the spin state (up or down) of neighboring spins connected by common molecular orbitals, that is, one or very few (in practice, usually not more than three) chemical bonds. This indirect interaction between two nuclear spins is described by a Heisenberg-type Hamiltonian,

$$H_J = \frac{2\pi}{\hbar} J \vec{I}_1 \bullet \vec{I}_2 \tag{15}$$

and therefore also called scalar coupling. Scalar couplings cause a splitting of the NMR resonances into multiplets. The separation of the multiplet components is given by the scalar coupling constant J, which is independent of B_0 and molecular orientation and reported in units of Hz. As a consequence of the scalar coupling constant of ${}^{1}J_{HN} \approx -93$ Hz between the backbone amide ${}^{1}H$ and the directly bonded ${}^{15}N$, all backbone amide resonances (resonating in the range from about 7 ppm to about 10 ppm) in Fig. 1 are actually doublets with a separation of about 93 Hz. Scalar coupling constant between the amide proton and the alpha proton depends on the protein backbone torsion angle Φ according to the Karplus equation,

$${}^{3}J_{HNH\alpha} = 6.4 \text{Hz} \times \cos^{2}(\Phi - 60^{\circ}) - 1.4 \text{Hz} \times \cos(\Phi - 60^{\circ}) + 1.9 \text{Hz}$$
 (16)

and can therefore be used to determine Φ experimentally [9]. Perhaps more importantly, the scalar coupling can be used for the coherence transfer during the mixing periods in multidimensional NMR experiments, thereby correlating the chemical shifts of spins that are connected by chemical bonds. For example, for each backbone amide ¹⁵N the [¹H,¹⁵N] HSQC spectrum shown in Fig. 2 reveals the chemical shifts of the attached backbone amide ¹H. These chemical shift correlations via scalar couplings – most commonly obtained from 3D or 4D triple resonance experiments on ¹³C/¹⁵N-enriched samples - are exploited to obtain sequence-specific chemical shift assignments [9].

In addition to the indirect interaction, spins also directly experience the magnetic field associated with the nuclear magnetic dipole moment of any spin that is located within a distance of approximately 5 Å. This dipole-dipole interaction between two nuclear spins A and B depends on the orientation of the internuclear vector relative to the static external magnetic field B_0 and is inversely proportional to the cube of their distance, r_{AB}^{3} . The direct interaction is completely anisotropic and averages to zero in isotropic solution. However, in second order the dipole-dipole coupling causes longitudinal cross-relaxation between spins that are close to each other in space, that is, longitudinal magnetization can be transferred from one spin to the other one during the mixing period of a suitably designed multidimensional NMR experiment (NOE SpectroscopY or NOESY). This second order effect is called the Nuclear Overhauser Enhancement (NOE) and it is inversely proportional to r_{AB}^{6} . In practice, ¹H-¹H NOEs are usually observed up to a distance of about 5 Å. The identification of many hundred ¹H-¹H NOE distance restraints from multi-dimensional NOESY spectra is the primary source of experimental information for conventional NMR protein structure determination. Reconstructing a protein conformation that satisfies both, the covalent geometry of the polypeptide chain and the hundreds of experimental restraints is a complex non-linear global optimization problem most commonly solved by a simulated annealing molecular dynamics simulation protocol [9]. As an example, Fig. 3 shows the ensemble of structures representing the 159-residue major cherry allergen Pru av 1 calculated from a total of 2438 experimental restraints, including 2299 NOE distance restraints [21].



Fig. 3: Backbone overlay of the 22 accepted structures of the cherry allergen Pru av 1 [21].

3 Protein Dynamics

3.1 Time scales and NMR spectral parameters

Proteins in aqueous solution at room temperature are not static but dynamic on a wide range of time scales (Fig. 4). Typically, bond vibrations are on the femtosecond time scale, side chain rotations on the picosecond time scale, local backbone dynamics and overall rotational diffusion on the nanosecond time scale, domain and loop motions on the nanosecond to millisecond time scale. Protein folding can be as fast as a few microseconds for small systems, but there is no evolutionary pressure to fold faster than the ribosomes can synthesize the polypeptide chain, so most proteins fold on the millisecond to minute time scale. The microsecond to minute time scale is also particularly relevant for many other biochemical processes such als ligand binding and enzyme catalysis.



Fig. 4: Typical time scales of protein motions (above the logarithmic axis of time scales) and NMR spectral parameters sensitive to these motions (below). Modified from [4] [22].

As shown in Fig. 4, protein dynamics on virtually all these time scales affect one or more NMR spectral parameters, which can hence be used to characterize these motions [4] [9]. The longitudinal and transverse relaxation times, T_1 and T_2 , respectively, together with the heteronuclear NOE between the amide ¹H and the amide ¹⁵N are determined by molecular motions on the pico- to nanosecond time scale of the Larmor frequency. Quantitative NMR relaxation experiments are therefore routinely performed to identify disordered segments of the polypeptide chain in a protein as well as to measure the correlation times of the overall rotational diffusion of the molecule, which reports on the hydrodynamic radius and hence the oligomerization state.

By contrast, if a protein exists in an equilibrium between two conformations A and B with different chemical shifts and the exchange happens on the micro- to millisecond time scale of the chemical shift difference (in units of rad/s) then this so-called intermediate chemical

exchange causes exchange line broadening of the NMR resonances [6], because (in a manner not dissimilar to the blurred perception of the spokes of a turning wheel) the resonances of the individual conformations can no longer be completely resolved due to the exchange. This line broadening can be quantified in CPMG relaxation dispersion experiments, in which the effective transverse relation rate, $R_{eff} = 1/T_{eff}$, is measured as a function of the pulsing frequency v_{CPMG} with which a train of 180° pulses is applied during a transverse relaxation interval of constant duration [4] [9]. If the 180° CPMG pulse train is applied with a sufficiently high CPMG frequency, the exchange line broadening is completely quenched and the intrinsic transverse relaxation time is recovered, $T_{eff} = T_2$. The relaxation dispersion profile $R_{eff}(v_{CPMG})$, which is usually recorded at two or more different magnetic fields B_0 , depends in a non-linear fashion on the kinetics of the exchange process (the exchange rates), the thermodynamics of the equilibrium (populations of the states according to the Boltzmann distribution), and the chemical shift difference of each spin between the two different states. Non-linear least squares fitting of suitable exchange models to the experimental relaxation dispersion profiles therefore allows quantification of the kinetics and thermodynamics of the conformational exchange and reconstruction of the chemical shifts of the exchanging conformations. Note that NMR relaxation dispersion spectroscopy is highly sensitive and allows identification of minor conformations populated to as little as 0.5% in equilibrium, which is so low that the resonances from such minor conformations are practically invisible in the NMR spectra themselves and are only accessible via quantification of the line broadening they cause to the resonances of the dominant conformation.

CPMG Relaxation Dispersion Spectroscopy



Fig. 5: Principle of CPMG relaxation dispersion spectroscopy.

Protein motions on the second time scale can be quantified using exchange spectroscopy (EXSY, zz exchange) [4] [9], including chemical exchange saturation transfer (CEST) [23]. A complementary technique that is particularly useful for studying protein folding is amide proton H/D exchange spectroscopy [24]. Dynamic processes on the minute time scale and

slower are, of course, directly accessible by observing the changes in the NMR spectra in real time [22].

3.2 Examples



Fig. 6: Kinetics, profile of the free energy, ΔG , and evolution of the three-dimensional structure along the folding pathway of the Fyn SH3 domain mutant A39V/N53P/V55L at 35°C as determined from NMR RD experiments. The SH3 domain folds from a random-coil like unfolded state U with no detectable structure, first by forming a 3-stranded β -sheet in the early transition state, via a metastable 4-stranded intermediate I (PDB 2L2P) into the native 5-stranded β -sandwich conformation F, which constitutes the directly observable ground state and was thus amenable to X-ray crystallographic and "conventional" NMR-spectroscopic studies (PDB 1SHF, 1NYG, 3CQT). The partially folded, aggregation-prone intermediate is stabilized by a non-native hydrophobic core, which has to be broken up and rehydrated in the late transition state to allow formation of the native core [14] [25] [26].

As mentioned above, CPMG relaxation dispersion (RD) spectroscopy allows the quantitative characterization of conformational exchange processes on the millisecond time-scale. Application of this technique to the Fyn SH3 domain mutant A39V/N53P/V55L (whose NMR spectra are displayed in Figs. 1 and 2) allowed us to delineate and accurately measure the

kinetics and thermodynamics of its three-state folding pathway [14]. The kinetics obtained from such experiments also served as a basis for NMR RD-based Φ -value analysis to probe the structure of rate-limiting transition states [25]. The chemical shifts reconstructed from CPMG RD experimentes even allowed us to determine the atomic-resolution structure of an "invisible" low-populated folding intermediate [26] using the CamShift structure calculation protocol [18] mentioned above. Together these studies provided a rare glimpse of a complex protein folding pathway in equilibrium under native conditions in atomic detail (Fig. 6).

¹HN and ¹⁵N Carr–Purcell–Meiboom–Gill (CPMG) relaxation dispersion experiments at various temperatures revealed that the third WW domain (WW3*) of the human ubiquitin ligase Nedd4-1 exists in an equilibrium between the natively folded peptide binding-competent state and a random coil-like denatured state [27]. The thermodynamics of the folding equilibrium was determined by fitting a thermal denaturation profile monitored by circular dichroism (CD) spectroscopy in combination with the CPMG data, leading to the conclusion that the unfolded state is populated to about 20% at 37°C. The kinetics of binding to the PY motif of the epithelial Na⁺ channel α subunit (α -hENaC) was determined by NMR lineshape analysis. These results show that the binding of the hNedd4–1 WW3* domain to α -hENaC is coupled to the folding equilibrium (Fig. 7). Our ongoing studies reveal that these coupled folding and binding equilibria are retained in the context of neighboring hNedd4-1 WW domains and likely also exist for the neighboring WW domains themselves [28].



Fig. 7: Nedd4-1 WW3* folding equilibrium (unfolded state, U, natively folded state, F) determined by CPMG coupled with the binding equilibrium of the PY motif ligand, L, determined by NMR lineshape analysis at 5°C. Modified from [27].

References

- [1] Purcell et al, Phys. Rev. 69, 37 (1946)
- [2] Bloch et al, Phys. Rev. 69, 127 (1946)
- [3] Williamson et al, J. Mol. Biol. 182, 295 (1985)
- [4] A. G. Palmer, Chem. Rev. 104, 3623 (2004)
- [5] A. M. Ruschak and L. E. Kay, J. Biomol. NMR 46, 75 (2010)
- [6] J. Keeler, Understanding NMR Spectroscopy (2nd ed., Wiley, Chichester, 2010)
- [7] M. H. Levitt, Spin Dynamics (2nd ed., Wiley, Chichester, 2008)
- [8] G. S. Rule and T. K. Hitchens, Fundamentals of Protein NMR Spectroscopy (Springer, Dordrecht, 2006)
- [9] Cavanagh et al, Protein NMR Spectroscopy (2nd ed., Academic Press, Burlington, 2007)
- [10] Q. Teng, Structural Biology (2nd ed., Springer, New York, 2013)
- [11] Markley et al, Pure Appl. Chem. 70, 117 (1998)
- [12] T.-L. Hwang and A. J. Shaka, J. Magn. Reson. A112, 275 (1995)
- [13] S. Grzesiek and A. Bax, J. Am. Chem. Soc. 115, 12593 (1993)
- [14] Neudecker et al, J. Mol. Biol. 363, 958 (2006)
- [15] D. S. Wishart and D. A. Case, Meth. Enzymol. 338, 3 (2001)
- [16] Y. Shen and A. Bax, J. Biomol. NMR 56, 227 (2013)
- [17] Shen et al, Proc. Natl. Acad. Sci. (USA) 105, 4685 (2008)
- [18] Robustelli et al, Structure 18, 923 (2010)
- [19] M. V. Berjanskii and D. S. Wishart, J. Am. Chem. Soc. 127, 14970 (2005)
- [20] Dominguez et al, J. Am. Chem. Soc. 125, 1731 (2003)
- [21] Neudecker at al, J. Biol. Chem. 276, 22756 (2001)
- [22] M. Zeeb and J. Balbach, Methods 34, 65 (2004)
- [23] Vallurupalli et al, J. Am. Chem. Soc. 134, 8148 (2012)
- [24] Krishna et al, Methods 34, 51 (2004)
- [25] Neudecker et al, Proc. Natl. Acad. Sci. (USA) 104, 15717 (2007)
- [26] Neudecker et al, Science 336, 362 (2012)
- [27] Panwalkar et al, Biochemistry 2016, 659 (2016)
- [28] Panwalkar et al, FEBS Lett. 591, 1573 (2017)

A 8 Computational Image Analysis

P. Kollmannsberger

Center for Computational and Theoretical Biology University of Würzburg

Contents

1	Intr	oduction	2								
	1.1	Motivation	2								
	1.2	Image data	2								
	1.3	Workflows	3								
	1.4	Common Tools	3								
2	Filte	Filtering and Denoising									
	2.1	Image Noise	4								
	2.2	Convolution	4								
	2.3	Nonlinear Filters	6								
	2.4	Feature Detection	7								
3	Ima	ge Segmentation	7								
	3.1	Thresholding and Watershed	7								
	3.2	Clustering and Superpixels	8								
	3.3	Morphology and Measurements	9								
	3.4	Machine Learning	0								
4	Spec	cial Topics 1	2								
	4.1	Motion Tracking	2								
	4.2	Stitching and Registration	2								
	4.3	Colocalization and FRET	3								
	4.4	Superresolution Analysis	3								

Lecture Notes of the $49^{\rm th}$ IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

1.1 Motivation

Modern imaging methods enable us to capture biological processes in unprecedented detail and with high temporal and spatial resolution. Examples are the macromolecular and subcellular ultrastructure in electron microscopy, the arrangement of proteins and cellular structures in superresolution microscopy, or the dynamics of cells and tissues using lightsheet imaging. High-resolution detectors and automated image acquisition result in enourmous amounts of data that push the limits of data transfer and storage infrastructure. Besides this technical problem, the more important challenge is to make full use of the potential of such "big" imaging data to test hypotheses and to develop new theories and models. Manual or semi-automated image analysis workflows quickly become a bottleneck, since they do not scale well with the amount and complexity of the data. Given that nearly all biological and biophysical research relies on modern microscopy methods, image analysis is a highly relevant part of the skill set in these fields. This lecture gives an overview of typical image analysis workflows, presenting the most important techniques and their typical usage scenarios. Rather than going into the mathematical detail of the different algorithms and procedures, the emphasis is on providing guidance on which tools and workflows are applied and why.

1.2 Image data

Images are spatial maps of intensity and color information, arranged in a regular coordinate grid of "pixels". It is important to note that pixels are not extended objects like squares or cubes but merely represent the amount of photons sampled at this location. A classical gray-level image is a two-dimensional grid I(x, y), whereas a RGB color image contains three such gray-level images representing the three colors, I(x, y, c). An image volume I(x, y, z) has an additional spatial dimension, whereas a movie contains an additional time or frame index, I(x, y, t). Different colors or spectral channels can also be represented as an image dimension, e.g. $I(x, y, \lambda)$. From a computational viewpoint, images can have an arbitrary number of dimensions. Many common image processing algorithms are however defined or implemented specifically for 2D or 3D images.

The intensity resolution of an image is defined as the amount of information that is used to enconde a single pixel. For example, an 8-bit gray-level image can distinguish $2^8 = 256$ different intensity levels, and a 16-bit image correspondingly $2^{16} = 65536$ different intensity levels. The total size of an image results from the number of pixels in the different image dimensions multiplied by the number of bytes used to encode a single pixel.

When storing image data for later analysis, it is important to also keep the information on how and when the image was recorded, and what the different image dimensions represent. Such meta-data includes for example the calibration of an image (conversion from pixels to metric length units), information on the type of illumination and microscope objective lens used to image the sample, up to details on experimental conditions and sample type. Many file formats exist that can store metadata alongside the raw pixel data. A common standard is the TIFF format. Specifically for microscopy image data, a common standard for TIFF metadata was defined by the OMERO consortium [1] to facilitate data exchange between different microscope brands and research groups. in the image. Adaptation of brightness and contrast in an image corresponds to changing the mapping between the actual intensity values in each pixel and the representation on the screen (brightness, color etc.). Sometimes, the foreground and background of an image are visible as two distinct distributions in the histogram, which can be used to automatically determine a threshold to separate both.

1.3 Workflows

Independent of the specific scientific question and type of image, there are certain stereotypic steps in image analysis that occur regularly. By chaining these components into analysis work-flows, many complex analysis task can be assembled from simpler steps. It is helpful to think in terms of such "workflows" when designing and reporting data analysis procedures in order to make it easier for others to understand and reproduce the different steps of such an analysis. While the specific components of such a workflow can vary greatly depending on the question that is to be addressed, stereotypic workflow in computational image analysis is the following:

- 1. Filtering and preprocessing (remove imaging artifacts and noise)
- 2. Applying a threshold (separate relevant signal from background)
- 3. Segmentation (identify and label structures of interest)
- 4. *Quantification* (measure properties of the segmented objects)

1.4 Common Tools

While a large pool of problem-specific software exists for computational image analysis, there are a number of open source as well as commercial frameworks and tools with a general scope that will be shortly described here. While there often is a division into graphical user interfaces (GUIs) and programming languages, the border between both is somewhat blurry, since there are examples where both paradigms can be used.

In general, the approach of encoding workflows in the form of macros, scripts or source code has a number of advantages when sharing and describing research data analysis. It is easier to reproduce and document if there is code rather than a more or less detailed description of which buttons were clicked in a graphical program. It is also more flexible than predefined functions in standard software. This can be important for scientific data analysis which often requires new ways of analyzing and interpreting data beyond pre-existing paths. While many powerful image processing and analysis algorithms are available in Python and C/C++, there are also many tools for non-programmers that provide more easy access and are attractive for scientists since they follow the open source paradigm:

- *ImageJ* is a versatile and extensible platform-independent tool covering most aspects of image analysis. It has a flexible plugin architecture and a built-in interpreter for many common scripting languages. *Fiji* is a distribution of ImageJ that includes a large collection of plugins [2].
- *Icy* is an open community platform for bioimage informatics that provides the software resources to visualize, annotate and quantify bioimaging data [3].

- *ilastik*, the interactive learning and segmentation toolkit, is a simple, user-friendly tool for interactive image classification, segmentation and analysis using modern machine learning [4].
- *CellProfiler* is an open-source software designed to enable researchers without training in computer vision or programming to quantitatively measure cell phenotypes from many images automatically [5].
- *KNIME* is an open data analytics platform that includes a growing number of image analysis functions and is optimized for scalability to high-performance clusters [6].

2 Filtering and Denoising

2.1 Image Noise

In principle, the representation of a real object in an image is always compromised by different types of noise and imaging artifacts (Fig. 1). The primary objective of image preprocessing and denoising is to restore the original ideal (noise-free) image as good as possible before extracting quantitative information. In order to perform denoising, it is important to have a good model of the type of noise present in an image.

One very generic and common type of noise is Gaussian white noise, which is caused e.g. by thermal fluctuations. Here, the intensity in each pixel varies independently according to a Gaussian distribution with $\mu = 0$ and width σ . This noise is added to the actual pixel intensities. Gaussian noise is an example of high-frequency noise, since the noise intensity varies at the smallest possible scale (largest frequency) present in the image, i.e. the pixel scale. Similarly, the stochastic noise due to low numbers of photons or electrons detected in a single pixel can be modeled as uncorrelated, high-frequency noise with a Poisson distribution.

The opposite would be low-frequency noise that varies on the scale of the entire image, for example due to uneven illumination that results in a continuous change of brightness from one end of the image to the other. Another type of noise is spot noise or spike noise, where only a certain percentage of pixels deviates by a large amount from the noise-free intensity. This can be caused by artifacts of analog-to-digital conversion.

Other types of noise or image defects are caused by the optical properties of the imaging setup. For example, all microscope objectives suffer from chromatic aberration: light of different wavelengths emanating from the same point in the specimen ends at slightly different locations on the image sensor. As a result, images taken in different color channels are slightly offset with respect to each other. The same can occur if the sample is moving slightly between acquisitions. Such image defects also need to be corrected prior to further quantitative analysis. Finally, the optical transfer function (point spread function or PSF) of any real-world imaging setup is imperfect, leading e.g. to blurring. If the PSF is known, for example by imaging a point-like object and recording the resulting intensity distribution, it can be used to deconvolve the image, i.e. to invert the transfer function of the optical setup. Due to noise, this is usually not possible and requires approximation and regularization.

2.2 Convolution

The most straightforward approach to minimize white, uncorrelated noise is to acquire several images in a row and simply average the intensity in each pixel over several images. With in-



Low Frequency Noise High-Pass Filter: Subtract Background

White Random Noise Low-Pass Filter: <u>Gaussian Blur</u>

Salt and Pepper Noise Nonlinear Filter: <u>Median Filter</u>

Fig. 1: Examples for different types of noise. Left: low-frequency noise, e.g. due to uneven illumination, can be filtered by a high-pass filter. Middle: High-frequency white noise is e.g. caused by thermal fluctuations in the camera or electronics and can be suppressed by a low-pass filter such as Gaussian blur. Right: Spot noise or salt-and-pepper noise can arise from conversion artefacts and is best suppressed by a nonlinear median filter.

creasing number of images N, the signal-to-noise ratio increases with \sqrt{N} according to the central limit theorem, since the noise-free intensity does not change between acquisitions. Unfortunately, averaging over many images is often not possible due to limited acquisition time. As an alternative, one could average over *different* pixels in the *same* image. For example, assuming that neighbouring pixels have similar intensities, a 3×3 local means filter would replace the noisy pixel intensity I(x, y) by the mean of the noisy intensities of all pixels in a 3×3 window around (x, y). This works well if the intensity does not vary much within that window, but will fail if there is a strong edge at this location in the image, i.e. a transition from dark to bright region or vice versa. Such edges will be blurred as a result. This example illustrates how a local means filter acts as a low-pass filter: high frequencies (noise and edges) are damped, while low frequencies (smooth changes in intensity) remain unaffected.

The local means filter described above can be formulated in terms of a *convolution* :

$$I'(x,y) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} I(x+i,y+j) \cdot H(i,j) = I * H$$
(1)

The noisy image I(x, y) is convolved with a 3×3 window (convolution kernel) H(i, j) with entries 1/9. This way, any local filter where the filtered pixel intensity is a function of the pixel intensities in a window around it can be represented as a convolution. Examples for typical filters are Gaussian blur (Gaussian weighted means), or Sobel and Laplacian (edge enhancement) (Fig. 2). The kernel entries can be negative, but they should always sum to one in order to maintain total intensity levels. If the kernel is symmetric, the corresponding convolution can be separated into individual convolutions along the different image dimensions, which greatly increases computational efficiency in the case of high-dimensional images.



Fig. 2: Examples for convolution based linear filters. The convolution kernel is a 3×3 matrix by which the neighbouring pixel intensities at each location are multiplied and then summed up. Here, two examples are shown for a high-pass filter (Laplacian kernel) to enhance high image frequencies (edges), and a low-pass filter (mean kernel) that suppresses high frequencies by locally averaging the pixel intensities.

One interesting property of convolution is that, according to the convolution theorem, it corresponds to multiplication in the Fourier domain:

$$I * H = \mathcal{F}^{-1} \{ \mathcal{F} \{ I \} \cdot \mathcal{F} \{ H \} \}$$

$$\tag{2}$$

In other words, if both the image I and the convolution kernel H are Fourier transformed into the frequency domain, the filtered image can be obtained by multiplying both and transforming the result back to real space.

This also illustrates the property of many filters to act as high pass or low pass filters. If an image is filtered with a Gaussian blur filter, high frequencies are filtered out. This is equivalent with scaling the image down, since very fine detail (encoded in the high frequencies) is lost. This property is sometimes used in image processing algorithms to construct scale-invariant operations by acting on scale space representations, i.e. a collection of images consisting of the original image filtered with Gaussian kernels of increasing size.

Convolutions are an example for an operation that can easily be parallelized, since the same operation (convolution or multiplication with a constant kernel) is applied to a large set of different data (image subsets). For this reason, they are well suited to be carried out on graphics processing units (GPUs), which are highly optimized for this type of calculation. Linear image filters (convolutions) can be accelerated by two orders of magnitude if they are carried out on consumer-level graphics cards.

2.3 Nonlinear Filters

Convolution-based local linear filters as described above are only one example for image filters. In general, filters need neither be local nor linear. As a simple example, we could replace the mean of the intensities in the local window by the median of intensities. There is no way to express this operation in terms of a convolution, hence this is a nonlinear local filter. It works very well to remove spike noise (salt and pepper) where a small number of pixels contains a large amount of noise.

Many filters that efficiently remove noise while retaining detail such as edges and texture are nonlinear. For example, a bilateral filter is a generalization of a Gaussian filter where the contribution of each pixel in the neighbourhood to the mean is not only weighted by its distance to the center as in the regular Gaussian blur filter, but also by its intensity. Similarly, the anisotropic filter performs averaging along image edges to maintain detail. The non-local means filter averages over similar image patches across the entire image and is therefore a non-local, non-linear filter.

2.4 Feature Detection

Image filters are not only used for denoising, but also to enhance specific image features. For example, edges can be enhanced using a Sobel filter kernel, or combinations of linear filters can be used to enhance blobs. The result of applying a filter to an image contains for each pixel information about its local context (e.g. if it is located near an edge). In this sense, the intensities stored in a filtered image can be used as features that characterize each pixel in terms of its local context in the image. If several such filter responses or feature maps are calculated on the image, the result is a "feature vector" for each pixel. This "feature space" can be used to perform supervised or unsupervised clustering of similar pixels into objects, which is the goal of image segmentation.

3 Image Segmentation

3.1 Thresholding and Watershed

In order to measure properties of objects in an image and to obtain quantitative results, the grid of pixels needs to be converted into a list of objects that are present in the image together with their properties. This step essentially corresponds to assigning class labels to each pixel. For example, if an image contains a cell with a nucleus and a membrane, four labels could be used to segment the image: nucleus, cytoplasm, membrane and background. The resulting image has the same dimensions as the original image, but each pixel intensity is replaced by the class label of that pixel.

The most basic example of image segmentation is applying an intensity threshold to an image (Fig. 3). All pixels in the image are labelled either as foreground (1) or background (0), according to whether their intensity is above or below a certain threshold value. The correct value of the threshold can either be set manually, or obtained automatically by a variety of methods. For example, the commonly used Otsu method tries to minimize intra-class variance [7]. In general, the threshold value can either be global for the entire image, or locally adaptive.

More generally, there are more than two different class labels in an image, one for each distinct object (e.g. cell nucleus). Such a segmentation can be obtained by first thresholding the original image, and then looking for connected regions of pixels that are either foreground or background. This is called region labeling or connected component analysis. Each connected region receives a distinct region label by which the corresponding object in the image can be



Fig. 3: An intensity threshold is used to separate foreground from background. In the histogram of all pixel intensities shown in the middle, two distributions are visible for bright foreground and dark background pixels. The ideal threshold in this case is in between the two distributions. If all pixels with intensity above the treshold are set to 1 and all other pixels to 0, the result is a mask of the foreground object (right).

identified. The pixel image with labels as pixel intensities could now be replaced by a list of objects with the pixel coordinates as object properties without losing information (Fig. 4). Often, it is not sufficient to label connected regions in order to identify all distinct objects. Sometimes, two objects will receive the same label because they touch each other, leading to a single connected component. A typical example would be two slightly overlapping cell nuclei or synaptic vesicles [8]. Sometimes such merged objects can be split using the *watershed algorithm*. The watershed algorithm starts from a number of seed pixels. These seeds are grown step by step into larger regions, bounded by the background. If two distinct regions that originated from different seed pixels touch, a new border is introduced at this position. The final result is equivalent to the original image, but with new borders separating the previously connected objects. The key step is setting the right seed points. This can be done manually, or automatically e.g. by selecting local extrema of the gradient image or of the distance map. The name "watershed" comes from the analogy of viewing each seed point as the bottom of a waterbasin. The basins are flooded, and watersheds are introduced whenever the water from

3.2 Clustering and Superpixels

one basin starts to flow into an adjacent basin.

Intensity-based thresholding followed by connected region labeling is often not sufficient to segment objects in more complex images. In this case, other pixel properties besides intensity and connectivity can be used to determine which pixels are "close" to each other and belong to the same object. As mentioned above, local convolution-based image filters collect information about the local context of each pixel, such as edge and texture information. This information can be used to determine which pixels are similar in terms of their context and belong to the same connected region, e.g. cell nucleus or cytoplasm. The pixel features (pixel intensities in the filtered images, e.g. after applying a gradient filter) correspond to a vector for each pixel, similar to its spatial coordinates (x, y, z) but in the "feature space". When the position of each pixel in this high-dimensional feature space is known, clustering can be applied to group similar pixels into objects.

In large images, clustering of individual pixels can become slow and inefficient. Splitting under-



Fig. 4: Typical image segmentation workflow: The original image is converted into a binary mask by applying an intensity threshold. The three leafs in the middle are a single connected object. To separate them, the distance map of the binary image is used to obtain seed points for a watershed segmentation. In the watershed image, the connected object is split into three separate objects. In the final step, all connected regions are labelled, and their properties (e.g. size) are calculated.

segmented regions (when two or more objects are merged into a single region) is difficult unless watershed can be applied. On the other hand, agglomeration of oversegmented regions (when an object is split into several regions) can be done in similar ways as pixel-based clustering. Clustering such oversegmented *superpixels* is more efficient than working on the "ultimate" segmentation (individual pixels), since pixels are already grouped into regions that belong together and there are by far fewer objects that need to be clustered. The features of superpixels are not filter-based as for pixels, but describe e.g. their shape, contrast or mean intensity, or properties of the interface to adjacent superpixels.

Many algorithms have been developed to oversegment images into superpixels, often adapted specifically to certain problems or image types. One generally applicable algorithm (SLIC)[9] has found widespread application in recent years, as it is largely independent of the specific image type. Starting from a regular grid of seed points, pixels are clustered by iteratively applying *k-means* clustering using the spatial coordinates (proximity) and the color information.

3.3 Morphology and Measurements

The final result of segmentation is an image where the original pixel intensities have been replaced by an integer number that denotes the label of the object to which a pixel belongs. The same information could be stored as a list of objects, each of which has a list of pixel positions attached to it. To reflect the relationships between objects, a tree instead of a list could also be used. To annotate this list or tree of objects with the properties of these objects, measurements can be carried out on the segmented label image. Examples would be simple properties such as total size, shape properties such as ellipticity, or properties related to the original pixel intensity,



Fig. 5: Elementary morphological operations (from left to right): erosion, dilatation, and their combinations morphological opening (dilation of the eroded image) and closing (erosion of the dilated image). The original image is dark blue, and the resulting image light blue. The structuring element here is a circle (grey). Source: Wikipedia

such as mean intensity or contrast. Those can be obtained by applying the corresponding label of each object as a mask to the original image.

For further processing of binary images, there are a number of morphological operations, equivalent to filters for intensity images (Fig. 5). The two elementary morphological operations, erosion and dilation, also are based on a filter window or kernel, called "structuring element". Erosion maintains all pixels where the structuring element, when centered on that pixel, entirely fits into the foreground object. This results in erosion of the edges of the object defined by the shape of the structuring element. Complementary to that, dilation adds pixels at the object edges corresponding to the structuring element for all pixels where the center of the structuring element is inside the object. Other morphological operations can be obtained by combining erosion and dilation: morphological closing is defined as the erosion of the dilated image, whereas morphological opening is the dilation of the eroded image.

Many other operations can be carried out on binary images to obtain different properties of the objects. For example, eroding the image iteratively until only isolated pixels remain is called ultimate erosion. Removing isolated foreground or background objects corresponds to filling holes. Computing the medial axis or skeleton of a binary image is a useful transformation: it maintains the topological properties of the objects and provides information on connectivity, substructure and shape, which we use e.g. to quantify cellular networks [10]. Another important transformation is the Euclidean distance transform of a binary image, which assings to each foreground pixel its Euclidean distance to the next background pixel or vice versa. The resulting distance map can e.g. be used to obtain seed points for the watershed algorithm, or as a starting point for various skeletonization and thinning algorithms.

3.4 Machine Learning

A typical computational image analysis workflow contains many of the abovementioned steps, from filtering to thresholding to segmentation and morphological processing, until object properties can be measured. In each step, there are many possible choices for different filters or algorithms, each of which has several parameters that need to be carefully tuned to obtain the desired result. In all but the most simple cases, this abundance of choices and parameters requires a lot of expertise and can be overwhelming for non-experts. An alternative approach is to

use machine learning: define a general path towards the solution together with some examples (training data), and use a trainable algorithm (classifier) to adjust the free parameters of the workflow such that the desired solution is obtained. Examples for common trainable classifier algorithms are k-nearest neighbour, support vector machine, variations of random forests, and artificial neuronal networks. The readily trained algorithm can then be tested on unseen data for which the solution is known to evaluate its performance. After repeating these steps until there is good performance on the test set, the algorithm can be applied, or the learned parameters can be included in a workflow.

For image segmentation, the training examples are pixels that have manually been assigned the correct label (e.g. cell nucleus, membrane etc.), and the workflow is typically a set of image filters (features) that enhance different structures in an image, such as edges and textures. The third component is a trainable algorithm that predicts the correct label for a pixel based on its feature vector (the list of filter responses for that pixel). In contrast to unsupervised clustering as described above for pixels and superpixels, a trainable algorithm learns from the training data and improves if the training set is increased. On the other hand, the training data set needs to reflect the diversity of the data on which the trained classifier will be applied, since it can only use the information in the training set. A typical problem that can arise is therefore overfitting of the training data: the classifier learns to perfectly classify the training examples, but will not be able to generalize well to previously unseen data. The same type of trainable classifiers can also be applied to superpixels rather than pixels, but the type of features is again different.

While trainable image segmentation is a significant improvement compared to fixed "targeted" segmentation pipelines, there still is the problem of finding or optimizing the right filters or features to separate the different classes of pixels or objects. This can be solved by turning the filter kernels themselves into trainable parameters in a convolutional neuronal network (CNN), also known as "deep learning". CNNs are artifical neuronal networks with special types of layers (Fig. 6). They contain convolutional layers where the neurons of a subsequent layer correspond to a filtered (convoluted) image of the previous layer - hence the name convolutional network. The weights of the layer correspond to the entries of the filter kernel of the convolution, and the activations of the first layer are simply the pixel intensities of the image to be classified. The other type of layer in CNNs are pooling layers that simply scale down the previous layer or image. The last layer is typically not convolutional, but fully connected (all input neurons are connected to all output neurons). The number of output neurons corresponds to the number of classes to be predicted. The term "deep learning" comes from the fact that CNNs typically are deep, i.e. have a much larger number of layers than "classical" artifical neuronal networks such as multilayer perceptrons. While CNNs were initially used to classify whole images (e.g. for digit recognition), it is also possible to use them for pixel-level segmentation by predicting the probabilities of individual pixels to belong to a certain class.

CNNs are trained by presenting the network with an input image, predicting the outcome class (forward pass), computing the error relative to the desired output class, and adapting the weights (filter parameters) by a small amount in the direction that reduces the error by calculating the gradient (backward pass). During training, the convolutional layers turn into filters that enhance the characteristic structures in the training images, analogous to feature filters in trainable segmentation. The resulting networks are typically better adapted to describe the segmentation problem compared to standard filters or hand-crafted features. This leads to a better prediction accuracy, making CNNs the current state-of-the-art of image segmentation. The trade-offs of this superior performance are very high computational requirements during training and the need for a large amount of high-quality hand-labelled training data. Also at the moment, the



Fig. 6: Schematic of a convolutional neuronal network (CNN). The feature maps result from applying linear image filters (convolutions) to the previous layer. The first layer of the network is the input image. Convolution layers alternate with pooling layers (subsampling). The last layer is fully connected and consists of one output neuron for each predicted class. The predictions can be classifications of entire images or of individual pixels. Source: Wikipedia

process of training and fine-tuning CNNs requires detailed expertise, which currently limits the widespread use of deep learning by non-experts. Due to the ongoing development in this field, however, improvements on these limitations can be expected.

4 Special Topics

4.1 Motion Tracking

A common task is to not only detect or measure objects in individual images, but to track their motion or deformation over time across several frames of a movie. Examples woulds be cell migration, single molecule tracking, or live-cell microrheology [11]. There is a large number of techniques and algorithms for solving different variants of this problem. In general, there are two steps involved in a tracking workflow. In the first step, objects to be tracked are identified in the individual frames of the movie. Depending on the type of problem, this can involve simple spot detection (e.g. for tracking single fluorescent molecules) or a complex segmentation workflow. In the second step, the detected objects from different frames are linked together to generate tracks. This means that the same object has to be identified in subsequent frames. If there are only a few particles per frame that do not move much between timepoints, the solution is to simply link to the closest particle in the consecutive frame. In the case of high particle density and larger displacements between frames, a more complex approach has to be taken. Object features have to be calculated to identify the same object in different frames, and the linking problem has to be solved by scoring different combinations of link assignments across all particles using additional plausibility criteria. A more general approach towards tracking does not identify objects or features, but tries to find identical image subwindows in consecutive frames, e.g. by using cross correlation.

4.2 Stitching and Registration

Image stitching is the problem of combining several smaller overlapping images of a region into a single large image. This occurs for example if a larger area of a sample is imaged by consecutively taking images of adjacent regions and moving the sample. The images have to overlap by a certain amount (typically at least a few percent of the total size) in order to obtain enough information to match the images together. The aim is then to minimize the mismatch between the overlapping regions of two images by translating and rotating the images, or applying more complex non-rigid transformations. Volume stitching is the same problem using 3D images. The mismatch can be determined by using pixel-wise differences, or more accurately by using intrinsic image features or fiducial markers that are present in both images and minimizing their distance. Image registration is a related problem, where the individual frames of a 3D stack or movie need to be aligned in order to eliminate drift or misalignment. Using one image as reference, the other images are again translated, rotated or transformed in order to minimize some mismatch criterion with respect to the reference image.

4.3 Colocalization and FRET

A common question that arises in biological imaging is whether two fluorescently stained molecules (e.g. a ligand and its receptor) colocalize, i.e. are in the "same" position within the limits of the image resolution and henceforth interact with each other. A simple qualitative way is to overlay images of both molecules imaged in two different channels using two different colors (e.g. red and green), and to look for overlapping (yellow) regions. This offers however only limited information and is not quantitative. A more quantitative approach towards colocalization analysis is to plot the intensities in both channels for each pixel as a scatterplot, and do a statistical correlation analysis. If the two molecules are present in the same pixel with high probability, the two channels will show a strong correlation. If the two molecules mostly occur in different locations, there will not be a strong correlation in the pixel intensities of both channels. Care has to be taken to avoid common pitfalls during imaging that can strongly influence the correlation analysis and lead to wrong conclusions. For example, crosstalk between the fluorescent channels has to be corrected for, and over- and underexposure of the images has to be avoided.

Interaction between fluorescently stained molecules, or different domains of a protein, can be determined by measuring the Foerster resonant energy transfer (FRET). If two fluorescent molecules (fluorophores) that form a FRET pair are within close vicinity (typically less than 10 nm), the excitation of the higher-energy fluorophore (donor) can result in non-radiating energy transfer to the lower-energy fluorophore (acceptor) and subsequent emission of a corresponding lower-energy photon (see Chap. A3 for more detail). By detecting and analyzing the relative emission of the donor and acceptor fluorophores, the FRET efficiency and hence the average distance between the two fluorophores can be detected. Due to the abovementioned limitations and sensitivity to errors for two-channel measurements, it is more reliable to instead measure the fluorescence lifetime of the donor channel using photon counting detectors, since this lifetime also depends on the amount of FRET.

4.4 Superresolution Analysis

A number of new techniques have recently been developed that made it possible to "break" the resolution limit of conventional fluorescence microscopy and increase the image resolution beyond the diffraction limit (see Chap. A1 for details). Many of these techniques (e.g. dSTORM, PALM, STED) share a common principle: they are based on single molecule localization microscopy. The idea behind this is that only a small fraction of the fluorescence

emitting molecules in a sample are excited at a time by quenching and selective or stochastic activation. The resulting images only contain a few isolated points, which are then localized individually. By collecting many such images, the full information in the image can be reconstructed. Since a single isolated point source can be localized to a precision beyond the diffraction limit, the final image contains a superresolved version of the same image obtained by classical fluorescence microscopy.

Superresolution microscopy by single-molecule localization poses two specific challenges to image analysis. The first challenge is the localization problem: obtaining the coordinates of the individual fluorescence-emitting molecules with high precision and high throughput. The common solution is to fit the intensity with a Gaussian distribution taking into account the effect of noise, using highly efficient parallel algorithms. The second challenge is to deal with the resulting list of point coordinates in order to extract the relevant information and to leverage the full potential of superresolution microscopy. This part of the problem goes beyond classical image analysis since the data are not in the form of images any more, but lists of coordinates. On the other hand, this offers new possibilites for analysis, for example to apply clustering methods in the coordinate space. Clustering is an important step to identify structures in the superresolution image from the individual point coordinates. A common clustering algorithm that is applied to point cloud data is "density-based spatial clustering of applications with noise" (DBSCAN) [12].

References

- [1] Goldberg, I., C. Allan, J.-M. Burel, D. Creager, A. Falconi, H. Hochheiser, J. Johnston, J. Mellen, P.K. Sorger, and J.R. Swedlow. (2005). "The Open Microscopy Environment (OME) Data Model and XML File: Open Tools for Informatics and Quantitative Analysis in Biological Imaging." *Genome Biol.* 6:R47. (http://www.openmicroscopy.org/)
- [2] Schindelin, J., Arganda-Carreras, I., Frise, E. et al. (2012). "Fiji: an open-source platform for biological-image analysis, *Nature methods 9*(7): 676-682. (http://fiji.sc)
- [3] de Chaumont, F. et al. (2012). "Icy: an open bioimage informatics platform for extended reproducible research", *Nature methods*, 9(7), 690-696. (http://icy.bioimageanalysis.org)
- [4] Sommer, C; Straehle, C; Koethe, U; Hamprecht, FA (2011). "ilastik: Interactive Learning and Segmentation Toolkit". *IEEE International Symposium on Biomedical Imaging:* 23033. (http://ilastik.org)
- [5] Carpenter AE, Jones TR, Lamprecht MR, Clarke C, Kang IH, Friman O, Guertin DA, Chang JH, Lindquist RA, Moffat J, Golland P, Sabatini DM (2006). "CellProfiler: image analysis software for identifying and quantifying cell phenotypes." *Genome Biology* 7:R100. (http://cellprofiler.org)
- [6] Berthold M.R. et al. (2008). "KNIME: The Konstanz Information Miner." In: Preisach C., Burkhardt H., Schmidt-Thieme L., Decker R. (eds) Data Analysis, Machine Learning and Applications. Studies in Classification, Data Analysis, and Knowledge Organization. Springer, Berlin, Heidelberg. (https://www.knime.com/)
- [7] Nobuyuki Otsu (1979). "A threshold selection method from gray-level histograms". *IEEE Trans. Sys., Man., Cyber. 9 (1): 6266.*
- [8] KV Kaltdorf, K Schulze, F Helmprobst, P Kollmannsberger, T Dandekar, C Stigloher

(2017). "FIJI Macro 3D ART VeSElecT: 3D Automated Reconstruction Tool for Vesicle Structures of Electron Tomograms." *PLOS Computational Biology 13 (1), e1005317.*

- [9] R Achanta, A Shaji, K Smith, A Lucchi, P Fua, and S Süsstrunk (2012), "SLIC Superpixels Compared to State-of-the-art Superpixel Methods", *IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 34, num. 11, p. 2274 - 2282.*
- [10] P Kollmannsberger, M Kerschnitzki, F Repp, W Wagermaier, R Weinkamer, P Fratzl (2017). "The small world of osteocytes: connectomics of the lacuno-canalicular network in bone." *New Journal of Physics 19 (7), 073019.*
- [11] P Kollmannsberger, CT Mierke, B Fabry (2011). "Nonlinear viscoelasticity of adherent cells is controlled by cytoskeletal tension." *Soft Matter* 7 (7), *3127-3132*.
- [12] M Ester, HP Kriegel, J Sander, and X Xu. (1996). "A density-based algorithm for discovering clusters in large spatial databases with noise." In: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD'96), Evangelos Simoudis, Jiawei Han, and Usama Fayyad (Eds.). AAAI Press 226-231.

A 9 Integrative Structural Biology and Hybrid Modeling

G.F. Schröder Structural Biochemistry Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Introduction	.2
2	Refinement of a Single Model	.3
3	De novo Modeling at Intermediate Resolution	.4
4	Uncertainty, Error and Dynamics	.4
5	Maximum Entropy and Bayes	.6
6	Comparing Models with Raw Data	.7
7	Discussion	.7
Refer	ences	.8

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.



Amount of Experimental Information

Fig. 1: The motivation for hybrid modeling is that the more information is used to build a model the more accurate it will be. Both experiment and simulation should be considered information, both improve the accuracy of a structural model.

1 Introduction

Hybrid modeling in structural biology describes the combination of computational modeling with experimental data to determine macromolecular structures (cf. Fig. 1). It has become particularly imortant to determine structures with low-resolution or sparse data, where the data alone would not suffice to build molecular models. Hybrid modeling is often used synonymously with integrative modeling, which emphasizes more the simultaneous use of different types of experimental information in the structure determination process.

Even though hybrid modeling in structural biology is a highly modern topic and a very active field of research, its birth can be traced to the method proposed by Jack and Levitt in 1978 [1] who introduced a hybrid energy function to optimize the energy of a protein at the same time as its fit to X-ray diffraction data. The fit of the protein model to the data has been defined *ad hoc* as an energy term E_{Data} , which is added with a weight *w* to the molecular mechanics energy E_{MM} of the protein

$$E_{\text{Hybrid}} = E_{\text{MM}} + wE_{\text{Data}}$$

 E_{Data} in its most basic form simply describes the deviation of experimental observables from those calculated from the model, e.g.

$$E_{\text{Data}} = \sum_{hkl} \left[F_{\text{obs}}(hkl) - F_{\text{calc}}(hkl) \right]^2$$

Minimization of this combined hybrid energy function, E_{Hybrid} , yields a refined structure that fulfills both the restraints imposed by the experimental data, as well as the stereo-chemical

restraints, which represent information we have on protein structures in general. Such refinement is instrumental in the interpretation of the data, e.g. in the case of crystallographic data, to improve phase information. Later this hybrid approach made it also possible to determine protein structures using restraints derived from NMR experiments. These early developments lead 30 years later to the notion that including as much information as possible is the best way of building models that are as accurate as possible [2]. A tour de force in such exhaustive integrative modeling was the determination of the molecular architecture of the nuclear pore complex [3, 4].

This chapter focuses on the use of intermediate to low-resolution X-ray diffraction and singleparticle cryo-electron microscopy (cryo-EM) data to determine (pseudo-) atomic models of protein structures. In particular cryo-EM have made tremendous progress in the past few years, which spurred the development of several new computational model-building techniques. The classical approach of building a single model that best fits the data is still prevalent even though the uncertainty in the modeling process could be captured more appropriately by generating an entire ensemble of models. However, the determination of model ensembles poses significant challenges, as will be discussed further below.

2 Refinement of a Single Model

At intermediate to low-resolution (4 - 8 Å) the parameter-to-observation ratio is too low to completely determine the atomic coordinates from the data alone. This problem can be solved by either reducing the number of parameters (e.g. by allowing only torsional degrees of freedom) or by adding information. The additional information could come from very different sources: For example the knowledge about similar structures can be used to guide the refinement; approaches such as the Deformable Elastic Network [5, 6], jelly-body, or reference model[7-9] restraints have been implemented in crystallographic refinement programs.

Another source of additional information are simulation and structure prediction techniques, which use molecular mechanics force fields. For example electrostatic interactions had disappeared completely from standard crystallographic refinement procedures, but it has recently been shown to improve the refinement [10]; even at high resolution when accounting for polarizability and anisotropic structure factors [11, 12].

More exhaustive sampling of protein conformations using all-atom explicit solvent MD simulations allow for larger conformational changes and can lead to significant improvement of the refinement with increased radius of convergence and better phases [13] as compared to standard crystallographic refinement. Similarly, the combination of energy-guided remodeling by the program Rosetta with the Autobuild and refinement tools of the program Phenix shows significant improvement in the refinement [14, 15]. This clearly demonstrates that force field/energy functions can provide valuable information that help to build better models or even solve structures that could otherwise not have been solved.

The same strategies have been followed in the refinement of (high-resolution) X-ray protein structures against (lower resolution) cryo-EM density maps, which is also referred to as flexible fitting. Good reviews on flexible fitting and modeling with cryo-EM density maps have been recently published[16-19]. When fitting structures against low-resolution data, over-fitting is a major concern. The standard tool to detect over-fitting is cross-validation, where a portion of the data, which needs to be independent from the rest, is not used for the

model fitting but only for model validation. Two cross-validation methods for the refinement of models into cryo-EM densities have been recently presented; one splits the set of particle images into two parts [20], and the other splits the density into a high- and low-resolution part[21].

3 De novo Modeling at Intermediate Resolution

In X-ray crystallography at intermediate to low resolution initially the phase information is often weak, resulting in erroneous and fragmented electron density. Therefore, the complete model cannot be built at once, but needs to be improved and extended iteratively while the phase information (and with it the electron density) improves. Furthermore, sequence assignment to the electron density is very difficult if initially only small model fragments can be built and side-chain densities are not clearly visible.

In contrast to X-ray crystallography, cryo-EM yields also phase information and therefore allows for computing a three-dimensional density map directly from the particle images. This makes it possible to build the entire model at once by using global optimization techniques. At intermediate to low resolution $(4 - 8\text{\AA})$, models can be built by detecting secondary structure elements, which can then be connected by modeling the missing loop segments [22-24].

Cryo-EM has made amazing progress in increasing the resolution in the past few years, mostly due to a new electron detector technology [25]. Resolutions below 4 Å have been reported even for particles with no or low symmetry [26-29] and resolutions of even better than 3 Å seem to be in reach. At such resolutions it is possible to build protein structures *de novo* from the density [30, 31].

Traditionally, the first step in building protein models is to trace the backbone in the density. The Pathwalker approach [32] treats the tracing procedure as a Traveling Salesman Problem, where each region in the density map needs to be visited exactly once. This provides a powerful additional restraint on finding the correct topology. Another promising although computationally expensive global optimization approach is the ACMI program [33, 34], which combines a local matching procedure and a global constraint procedure in a probabilistic framework.

4 Uncertainty, Error and Dynamics

Every model has errors. When building a model the uncertainty in the atomic coordinates needs to be accounted for. A good overview of uncertainty in integrative structural modeling is given in Ref. [35]. The uncertainty can be expressed by an ensemble of models as usual in NMR, or by B-factors in crystallography. However, how much of this uncertainty is just error or actually conformational variance is often unclear and cannot easily be decided.

It has been suggested that the extent of atomic motions is significantly underestimated by B-factor refinement [36], especially when distinct alternate conformations are sampled [37]. It becomes increasingly clear that even with high-resolution X-ray diffraction data (and even more so with low-resolution data), a single model might often not be sufficient to describe the experimental data, due to conformational heterogeneity and dynamics [38-41].



Fig. 2: Showing different scenarios of combining experimental data (green) with predicted ensembles (orange). Two different interpretations of the distribution of coordinates (structural ensemble) are considered.

a) Using Bayes formalism to combine prior knowledge (predicted ensemble) with the probability distribution of an experimental observable (likelihood to observe measured data for a given set of coordinates), which also encodes the error of the measurement. The posterior is obtained by multiplying the prior and the likelihood. Note that the width of the posterior distribution is smaller since the uncertainty is decreased by the measurement.

b) Schematic plot of a predicted ensemble that is minimally biased by experimental measurement of the ensemble average (green dashed line). Here the biased ensemble (blue) has the same average value as the measured ensemble average, but is otherwise as similar as possible to the predicted ensemble.

c) Predicted and measured ensembles both yield a distribution of coordinates. The best estimate for the distribution is the (possibly weighted) average of both distributions (blue). No errors are considered. Note that the distribution does not necessarily become narrower. In the extreme case, if prediction and experiment both yield the same correct distribution, combining this knowledge should yield that exact same distribution.

d) Same as in *c)* but here errors of measured and predicted model ensembles are represented by an ensemble of coordinate distributions. Note that while the error of the combined distribution (blue) becomes smaller, its width does not.

Interestingly, careful analysis of electron density by estimating its noise level reveals hidden conformational motions, low-occupancy alternative conformations and ligands [42-44].

Two interpretations of a model ensemble need to be distinguished (cf. Fig. 2): In the first interpretation, the ensemble represents the uncertainty and the width of the ensemble reflects the fact that the amount of restraints (structural information) that the experiment provides is limited. In that case the aim is to find the *broadest* ensemble that is in agreement with data, according to the principle of maximum entropy.

In the second interpretation, it is assumed that the data provide information on the dynamics of the molecule and that the ensemble therefore represents true conformational variance, which means that the refinement of an (restrained) ensemble increases the number of parameters. In that case the aim is to find the *smallest* ensemble that describes the data well (in accordance to Occam's razor) to keep the parameter-to-observations ratio low.

Following the second interpretation, the Sparse Ensemble Selection (SES) method selects the smallest (sparsest) non-uniformly weighted representative ensemble that explains the experimental data to within a desired error. SES does not require any prior information such as a force field; the only restraint is sparsity [45].

Several ensemble refinement methods have been developed for NMR structure determination [46-53] and also for X-ray crystallography [54-59]. While the modeling of ensembles yields a picture of the structural heterogeneity, the combination of NMR with X-ray diffraction data adds information about the timescale of atomic motions and shows that the same picosecond motions observed in solution also occur in the crystalline state at room temperature [60].

5 Maximum Entropy and Bayes

Oftentimes it is easier to measure ensemble averages than the probability distribution of conformational states. For example, the 3D reconstruction from single-particle cryo-EM images is an ensemble average, while information about the conformational distribution is hidden in the collection of very noisy particle images. The determination of the structural distribution is usually an under-determined problem and there exist therefore many different ensembles that lead to the same ensemble average.

To help defining the ensemble, an estimate of the conformational distribution can be obtained from force field based molecular simulations. Recently, Pitera and Chodera presented an approach to bias molecular simulations by experimental data according to the maximum entropy principle, such that the simulated ensemble is minimally biased by the experimental observations [61].

Later it was proven that maximum entropy ensembles are obtained by restrained-ensemble simulations [62-64]. The recently proposed method of experiment directed simulation (EDS) yields such a biasing potential more efficiently [65]. Olsson *et al.* showed how an expectation maximization algorithm yields a minimally biased native state ensemble with NMR data [66].

Errors of experimental data are best included using Bayes formalism, as this allows for using exact error and noise models. The application of Bayes formalism to NMR structure determination has been described in the influential work by Rieping *et al.* [67, 68]. Recently,
the formalism has been applied to modeling structures with single particle cryo-EM images [69] and X-ray single-particle diffraction images [70].

6 Comparing Models with Raw Data

For the interpretation of the experiment, it is in general best to compare the model with the raw experimental data to fully exploit all information provided by the experiment. However, oftentimes the raw data sets are huge and difficult to handle. For example, in cryo-EM the number of single-particle images is typically on the order of 10,000 to several 100,000, such that refinement of a model against the individual particle images is computationally expensive.

Using class-averages instead reduces the number of images to work with and allows for building the assembly structure of macromolecular complexes. Velazquez *et al.* [71] score candidate models of the assembly by the similarity of model projections to class-averages. A Monte-Carlo sampling of this scoring function then yields assembly structures. Also the refinement of models against class-averages is possible [72]. Class-averages, however, have the limitation that any variance is lost and sample heterogeneity is averaged out. Cossio and Hummer [69] showed that correct model conformations can be detected by comparing model projections against the raw particle images; even the correct model ensemble can be determined from a set of images containing a mixture of difference conformations.

Also in X-ray diffraction images there is more information than just the Bragg reflection peaks. Diffuse scattering outside the Bragg reflection reports on correlations between the motions of atoms, which could provide valuable restraints on protein motion [73].

7 Discussion

Model building starts to become a bottleneck in cryo-EM, in particular since typical molecular systems studied by cryo-EM are large and model building is time consuming. Building atomic models at 4.5 Å or worse is still a considerable challenge and often cannot be done reliably let alone in an automated way.

Ensembles predicted from molecular simulations provide important additional information. However, combining experimental data with predicted ensembles requires accounting for uncertainties not only of the data but also of the prediction arising from errors in the force fields as well as insufficient conformational sampling.

How the error of the force field parameters (e.g. partial charges of atoms) can be determined and how these errors propagate to the error of the predicted ensemble, in particular the error on ensemble averages or populations of conformational states is still an open question. But it needs to be resolved to avoid an overly optimistic influence of the predicted information on hybrid modeling.

A general quality measure of hybrid models, depending on how much information (both predicted and experimental) was used to build it, is still missing. For this the information content of any experimental data set [74] needs to be quantified. Such a measure is necessary to compare the quality of hybrid models obtained with different types and combinations of experimental data and predictions.

References

- 1. Jack, A. and M. Levitt, Acta Crystallogr A, 1978. 34: p. 931-935.
- 2. Russel, D., et al., PLoS Biology, 2012. 10: p. e1001244.
- 3. Alber, F., et al., Nature, 2007. 450: p. 683-694.
- 4. Alber, F., et al., Nature, 2007. 450: p. 695-701.
- 5. Schröder, G.F., A.T. Brunger, and M. Levitt, Structure, 2007. 15: p. 1630-41.
- 6. Schröder, G.F., M. Levitt, and A.T. Brunger, Nature, 2010. 464: p. 1218-1222.
- 7. Murshudov, G.N., et al., Acta Crystallogr D, 2011. 67: p. 355-367.
- 8. Nicholls, R.a., F. Long, and G.N. Murshudov, Acta Crystallogr D, 2012. 68: p. 404-17.
- 9. Headd, J.J., et al., Acta Crystallogr D, 2012. 68: p. 381-90.
- 10. Fenn, T.D., et al., Structure, 2011. 19: p. 523-33.
- 11. Schnieders, M.J., et al., Acta Crystallogr D, 2009. 65(9): p. 952-965.
- 12. Fenn, T.D. and M.J. Schnieders, Acta Crystallogr D, 2011. 67(11): p. 957-965.
- 13. McGreevy, R., et al., Acta Crystallogr D: Biol Crystallogr, 2014. 70: p. 2344-2355.
- 14. DiMaio, F., et al., Nature, 2011. 473: p. 540-543.
- 15. DiMaio, F., et al., Nature Methods, 2013. 10: p. 1102-4.
- 16. Esquivel-Rodríguez, J. and D. Kihara, J Struct Biol, 2013. 184: p. 93-102.
- 17. Villa, E. and K. Lasker, Curr Opin Struct Biol, 2014. 25: p. 118-25.
- 18. Lindert, S., P.L. Stewart, and J. Meiler, Curr Opin Struct Biol, 2009. 19: p. 218-25.
- 19. Lander, G.C., H.R. Saibil, and E. Nogales, Curr Opin Struct Biol, 2012. 22: p. 627-35.
- 20. DiMaio, F., et al., Protein Science, 2013. 22: p. 865-8.
- 21. Falkner, B. and G.F. Schröder, Proc Natl Acad Sci USA, 2013. 110: p. 8930-5.
- 22. Lindert, S., et al., Structure, 2012. 20: p. 464-78.
- 23. Baker, M.L., et al., Biopolymers, 2012. 97(9): p. 655-668.
- 24. Baker, M.L., T. Ju, and W. Chiu, Structure, 2007. 15: p. 7-19.
- 25. Bammes, B.E., et al., J Struct Biol, 2012. 177: p. 589-601.
- 26. Lu, A., et al., Cell, 2014. 156: p. 1193-206.
- 27. Bai, X.-C., et al., eLife, 2013. 2: p. e00461.
- 28. Li, X., et al., Nature Methods, 2013. 10: p. 584-90.
- 29. Liao, M., et al., Nature, 2013. 504: p. 107-12.
- 30. Wang, Z., et al., Nat Commun, 2014. 5: p. 4808.
- 31. Baker, M.L., et al., J Struct Biol, 2011. 174: p. 360-373.
- 32. Baker, M.R., et al., Structure, 2012. 20: p. 450-63.
- 33. DiMaio, F., J. Shavlik, and G.N. Phillips, Bioinformatics, 2006. 22: p. e81-e89.
- 34. Soni, A. and J. Shavlik, J Bioinform Comput Biol, 2012. 10(1): p. 1240009.
- Schneidman-Duhovny, D., R. Pellarin, and A. Sali, Curr Opin Struct Biol, 2014. 28C: p. 96-104.
- 36. Kuzmanic, A., N.S. Pannu, and B. Zagrovic, Nat Commun, 2014. 5: p. 3220.
- 37. Janowski, P.A., et al., J Am Chem Soc, 2013. 135: p. 7938-48.
- 38. Forneris, F., B.T. Burnley, and P. Gros, Acta Crystallogr D, 2014. 70: p. 733-43.
- 39. Burnley, B.T., et al., eLife, 2012. 1: p. e00311.
- 40. Furnham, N., et al., Nat Struct Mol Biol, 2006. 13: p. 184-185.
- 41. Altman, R.B. and O. Jardetzky, J Biochem, 1986. 100: p. 1403-1423.
- 42. Fraser, J.S., et al.,. Proc Natl Acad Sci USA, 2011. 108: p. 16247-52.
- 43. Lang, P.T., et al., Proc Natl Acad Sci USA, 2014. 111: p. 237-42.
- 44. van den Bedem, H., et al., Nature Methods, 2013. 10: p. 896-902.
- 45. Berlin, K., et al., J Am Chem Soc, 2013. 135: p. 16595-609.

- 46. Fennen, J., A. Torda, and W.F. van Gunsteren, J Biomol NMR, 1995. 6(2): p. 163-170.
- 47. Bürgi, R., J.W. Pitera, and W.F. van Gunsteren, J Biomol NMR, 2001. 19(4): p. 305-320.
- 48. Kim, Y. and J.H. Prestegard, Biochemistry, 1989. 28(22): p. 8792-8797.
- 49. Lindorff-Larsen, K., et al., Nature, 2005. 433.
- Lindorff-Larsen, K., R.B. Best, and M. Vendruscolo, J Biomol NMR, 2005. 32: p. 273-280.
- 51. Richter, B., et al., J Biomol NMR, 2007. 37: p. 117-135.
- 52. Lange, O.F., et al.,. Science, 2008. 320: p. 1471-1475.
- 53. Linge, J.P., et al., Proteins, 2003. 50: p. 496-506.
- 54. Kuriyan, J., et al., Proteins, 1991. 10(4): p. 340-358.
- 55. Levin, E.J., et al.,. Structure, 2007. 15(9): p. 1040-1052.
- 56. Burnley, B.T., P.V. Afonine, and P. Gros, eLife, 2012: p. e00311.
- 57. Forneris, F., B.T. Burnley, and P. Gros, Acta Crystallogr D, 2014. 70(3): p. 733-743.
- 58. Kohn, J.E., et al., PLoS Comp Biol, 2010. 6: p. 1-5.
- 59. Burling, F.T. and A.T. Brunger, Israel Journal of Chemistry, 1994. 34: p. 165-175.
- 60. Fenwick, R.B., et al.,. Proc Natl Acad Sci USA, 2014. 111: p. E445-54.
- 61. Pitera, J.W. and J.D. Chodera, J Chem Theo Comp, 2012. 8: p. 3445-3451.
- 62. Roux, B. and J. Weare, J Chem Phys, 2013. 138: p. 084107.
- 63. Cavalli, A., C. Camilloni, and M. Vendruscolo, J Chem Phys, 2013. 138(9): p. 094112.
- Boomsma, W., J. Ferkinghoff-Borg, and K. Lindorff-Larsen, PLoS Comp Biol, 2014. 10: p. e1003406.
- 65. White, A.D. and G.A. Voth, J Chem Theo Comp, 2014. 10: p. 3023-3030.
- 66. Olsson, S., et al., J Chem Theo Comp, 2014. 10: p. 3484-3491.
- 67. Rieping, W., M. Habeck, and M. Nilges, Science, 2005. 309: p. 303-306.
- 68. Nilges, M., et al., Structure, 2008. 16(9): p. 1305-1312.
- 69. Cossio, P. and G. Hummer, J Struct Biol, 2013. 184: p. 427-37.
- 70. Walczak, M. and H. Grubmüller, Phys Rev E, 2014. 90: p. 022714.
- 71. Velázquez-Muriel, J., et al., Proc Natl Acad Sci USA, 2012. 109: p. 18821-6.
- 72. Zhang, J., P. Minary, and M. Levitt, Proc Natl Acad Sci USA, 2012. 109: p. 9845-50.
- 73. Wall, M.E., et al., Structure, 2014. 22: p. 182-4.
- 74. Berman, M. and P. Van Eerdewegh, American Journal of Physiology, 1983. 245: p. R620-R623.

A 10 Simulation Methods in a Nutshell

M. Ripoll

Theoretical Soft Matter and Biophysics Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Why	y do we need computer simulations?	2
2	Mic	roscopic simulation methods	3
	2.1	Molecular Dynamics simulations	3
	2.2	Interaction potentials	3
	2.3	System observables	6
	2.4	Non-equilibrium simulations	7
	2.5	Monte Carlo Simulations	8
3	Hyd	rodynamic simulation methods	9
	3.1	Brownian dynamics (BD)	11
	3.2	Smoothed particle hydrodynamics (SPH)	11
	3.3	Lattice Boltzmann (LB)	11
	3.4	Dissipative particle dynamics (DPD)	12
	3.5	Multiparticle collision dynamics (MPC)	14
4	Con	clusions	15

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Why do we need computer simulations?

Physics, Biology, and Science in general have fruitfully developed for centuries without the help of computer simulations. The description of material properties was obtained by combining experiments and analytical approaches. The experimental observation of a system needs always of an apparatus that can be the eye, a clock, a microscope, or a nuclear magnetic resonance device. These apparatus are able to access different but limited amounts of information like average density, or the position of particle in a limited size range. To interpret and understand the experimental observations of a particular system is necessary the construction of a model, what requires the assumption of ideal conditions, and simplified interactions. The behavior of such model is normally described in terms of equations which are soluble only in a very few exceptional cases. For example, to solve the motion of more than two interacting particles, even with simple Newtonian mechanics, resulted essentially impossible before the appearance of computers. The standard procedure is then the performance of more or less dramatic approximations. Therefore, for an analytical theory to explain satisfactorily experimental results is required that, enough precision of apparatus, good choice of the system model, and validity of the performed approximations.

With the appearance of computers (in the early 1950's) a new perspective started to be employed to understand physical systems. The equations that describe a model with simplified interactions are solved now for a few particles or for millions of them. The algorithms necessary to perform the simulations also have limitations, but they are in general enormously less restrictive than the analytical approximations performed in theory. Comparing results from simulations and analytical theories serves first to test the performed approximations. Comparing results from simulations and experiments serves to test the simplified interactions assumed for designing the model (see Fig. 1). Furthermore, computer simulations can be used to predict material properties, by anticipating systems or conditions that are not always easy or cheap to perform experimentally. The information provided by a simulation is in general also much more detailed than the experimental one, since very precise information of the system constituents is generated. This dual role of the simulations, that allows to bridge the distance between the models and the theory, and between the models and the experiments, makes that they are sometimes referred *theoretical simulations* by experimentalists or as *computer experiments* by theoreticians.

Metropolis, Rosenbluth, Rosenbluth, Marshall, Teller and Teller, introduced Monte Carlo in 1953 [1] to investigate fluid properties. Alder and Wainwright introduced Molecular Dynamics in 1957 [2] to study the phase transitions of hard sphere systems, and in 1964 Rahnman carried out the first simulations using a realistic potential for liquid Argon [3]. The first simulation of a protein was done by Levitt and Warshel in 1975 [4]. Since then, the simulation techniques have substantially evolved. Particular problems have developed their own specialized techniques, like combination with quantum mechanical methods to find application at the quantum level [5], or hybrid simulation with mesoscopic techniques in order to bridge the gap with larger scales as those in soft matter systems.

In this contribution, some general concepts, common to all simulation methods will be discussed, and we will introduce a glimpse over the main simulation methods from the microscopic and the mesoscopic viewpoints. For further insight in these and other simulation methods we recommend, several specialized books and journal reviews [6, 7, 8, 9], that provide very good detailed descriptions, and 'recipes' of a large list of simulations aspects.





Fig. 1: Relations between experiment, theory and simulation when analyzing a system.

2 Microscopic simulation methods

2.1 Molecular Dynamics simulations

Molecular Dynamics (MD) is a technique to compute steady and transport properties of manybody systems, and it was together with Monte Carlo, the first simulation method to be proposed. The basic idea of the approach is that the material properties are described by the Newton's equation of motion of each component atom, molecule, monomer, or any other relevant subunit. More precisely, N particles of masses m_i with (i = 1, ..., N) are considered to satisfy the equations of motion,

$$m_i \frac{d^2}{dt^2} \mathbf{r}_i = \mathbf{F}_i,\tag{1}$$

where \mathbf{r}_i are the particle positions, and \mathbf{F}_i are the resulting forces exerted on particle *i*. In the case of conservative systems, the forces are obtained from the potential energy with $\mathbf{F}_i = -\nabla U(\mathbf{r}^N)$, with $\mathbf{r}^N = (\mathbf{r}_1, \dots, \mathbf{r}_N)$. Through the numerical integration of these equations of motion, a natural time scale is built in, and the phase space is deterministically sampled. MD generates information at a microscopic level, since position and velocities of all the component atoms are known as a function of time. Macroscopic observables like pressure, or heat capacity, are obtained from the microscopic information via statistical mechanics.

2.2 Interaction potentials

The potential energy of N interacting particles can be divided in terms of the number of particles simultaneously involved in the calculation of the interaction, this is

$$U(\mathbf{r}^{N}) = \sum_{i} u_{1}(\mathbf{r}_{i}) + \frac{1}{2} \sum_{i} \sum_{j} u_{2}(\mathbf{r}_{i}, \mathbf{r}_{j}) + \frac{1}{4} \sum_{i} \sum_{j} \sum_{k} u_{3}(\mathbf{r}_{i}, \mathbf{r}_{j}), \mathbf{r}_{k}) + \cdots$$
(2)

The first contribution corresponds to the effect of external fields applied on the system, like gravity or electric fields. The second term is the pair potential and includes the force that the presence of one particle induces in a second one. The third term, report forces induced due to

the presence of triples, and similarly subsequent terms can be added for increasing the number of particles in the interaction. The calculation of u_1 is then and operation of order N, and u_2 of order N^2 since all pairs of particles have to be identified. Even more costly will be to determine u_3 , since there are $\mathcal{O}(N^3)$ triples on a system. Furthermore, high order contributions are typically small, and the pair potential is chosen such that gives and effective potential resulting from all the high order terms. A real two-body potential would reproduce experimental data, without and explicit dependence on system parameters like density or temperature. The normally employed effective pair potentials will though depend on system parameters.

It is in the choice of the interaction potentials where the system under study will be mostly determined. In this way, the interactions are very different when simulating an ideal gas, or a system of charged particles; when simulating flexible polymer, or a molecule with a particular defined configuration like a water molecule; or when simulating an spherical vesicle or a red blood cell. Moreover, the interaction will determine the level of description. Each simulated particle can correspond with a particular atom and its precise location inside a molecule, like DNA; or each particle can represent a group of atoms, such that the averaged properties will be to some extent similar those of the real molecule.

The interactions are divided in bonded and non-bonded. The first ones are all the intramolecular interactions, while in the later, interactions such as electrostatic or van der Waals are considered. Here some of the most commonly employed potentials are introduced (see Fig. 2).

Bond potential To model a bond between two particles, it is most common to employ the harmonic potential,

$$U^{H}(r_{ij}) = k^{H} \left(r_{ij} - b \right)^{2} \tag{3}$$

where b is the reference *bond length*, and k^{H} the spring constant. The harmonic potential is basically a Taylor approximation to more sophisticated potentials around the reference bond length.

Angle potential Between the two bonds of three consecutively linked particles an angle θ can be defined. This angle can be fixed to have a certain value θ_0 equal or different from zero,

$$U^{A}(\mathbf{r}^{N}) = k^{b} \left(\theta - \theta_{0}\right)^{2}.$$
(4)

The bending constant k^b determines how strongly are the deviations from the reference angle. To fix the relative position of four consecutively linked particles can be done with torsion potentials. The dihedral angle potential, for example is used to constrain the rotation of one particular bond around the plane defined with another two bonds.

Van der Waals potential Two particles not bounded experience frequently a combination of repulsive and attractive interactions. The strong repulsive interaction at short distances are due to excluded volume effects, while the soft attraction at larger distances is a consequence of the correlation between the electron clouds surrounding the atoms ('van der Waals or 'London' dispersion). The most common potential to model the van der Waals interactions is known as the *Lennard-Jones potential*,

$$U^{LJ}(r) = 4\epsilon \left[\left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right], \qquad r_c = 2.5\sigma.$$
(5)

where the two potential parameters are ϵ , the minimum potential depth, and σ the collision diameter. For effectively reasons, the cut-off radius r_c introduces a distance from which the potential is approximated to zero.

To study systems of purely repulsive particles, a frequent approach is to consider only the repulsive part of this potential by truncating it precisely at the minimum and increasing its value to zero at that value,

$$U^{WCA}(r) = 4\epsilon \left[\left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right] + \epsilon, \qquad r < r_c = 2^{1/6}\sigma.$$
(6)

This is the repulsive Lennard Jones potential also know as the *Weeks-Chandler-Andersen potential* (WCA).

Electrostatic interaction The electrostatic interaction between two particles of charge q_i , and q_j is given by the Coulomb potential,

$$U^{C}(r_{ij}) = \frac{1}{4\pi\varepsilon_{0}\varepsilon_{R}} \frac{q_{i}q_{j}}{r_{ij}}$$

$$\tag{7}$$

where ε_0 is the permittivity of vacuum, and ε_R is the dielectric constant.



Fig. 2: Illustration of the most basic interparticle interactions

By combining the above potentials, the number of possible macromolecules and structures that can be modeled is unlimited. Some basic examples are shown in Fig. 3. Colloidal dispersions can be modelled by particles where the main interaction are van der Waals or Coulomb. By linking a determined number of particles with harmonic springs, the properties of polymers can be studied. If the linked particles have also a bending potential that keeps the structure elongated and stiff, we will be analyzing the properties of rod-like colloids. If the string of particles is consider simultaneously with a surrounding solvent and a few particles have an attractive interaction with the solvent and the rest a repulsive attraction with the solvent, we will have a basic lipid molecule , that with proper concentration values will be able to self-assemble.



Fig. 3: Illustration of the construction of macromolecules by combining various interaction potentials.

2.3 System observables

The information generated by computer simulations is very detailed and lays at the microscopic level. To translate it into macroscopic terms concepts of Statistical Mechanics are to be employed [10, 11]. The number of properties that can be characterized and quantified is therefore very large. It will depend then on the particularities of the chosen system, and on the interests of the study. The system temperature T is for example defined through E_k , the system kinetic energy

$$E_K = \left\langle \frac{1}{2} \sum_{i}^{N} m_i v_i^2 \right\rangle = \frac{d}{2} N k_B T \tag{8}$$

where v_i is the particle velocity and d the system dimension. The pressure can be calculated from the virial equation,

$$P = \rho k_B T + \frac{1}{dV} \left\langle \sum_{i < j} \mathbf{F}_{ij} \cdot \mathbf{r}_{ij} \right\rangle, \tag{9}$$

where V is the system volume. Note that in the two previous expressions, averaged quantities are employed. These averages are defined when a quantity A can be independently measured N_R times. The average value of A is then,

$$\langle A \rangle = \frac{1}{N_R} \sum_{k=1}^{N_R} A_k.$$
⁽¹⁰⁾

In systems where time is defined, the independent measurements can be made at separated enough time points, or by repeating the simulation with different initial values. And finally, in order to estimate the accuracy of the measurement, the standard deviation of the average is calculated as

$$\delta(A) = \sqrt{\frac{1}{N_R - 1} \sum_{k=1}^{N_R} \left(A_k - \langle A \rangle\right)^2} \tag{11}$$

Apart from thermodynamic properties, static information in real or reciprocal space permits to obtain structural properties which are of fundamental when characterizing a system. The pair correlation function g(r) for example, measures the probability to find a particle at a distance r from another particle. If this function shows or not some periodic behavior will allow to determine if the system is a gas, a liquid, or a crystal. Dynamical properties are also of great interest.

In MD the time evolution is intrinsic such that time dependent quantities can be described. As an example, the mean squared displacement $\langle (\Delta r(t))^2 \rangle = \langle \sum_i (\mathbf{r}_i(t) - \mathbf{r}_i(0))^2 \rangle$, shows the transition from a short time ballistic regime in which $\langle (\Delta r(t))^2 \rangle \sim t^2$ to a long time diffusive regime, $\langle (\Delta r(t))^2 \rangle \sim t$, what allows to determine the self-diffusion coefficient D_s

$$\left\langle (\Delta r(t))^2 \right\rangle = 2dD_s t. \tag{12}$$

Reduced units Quantities like time or temperature obtained from simulation results are expressed in terms of the model parameters. In order to relate simulation results with experimental data is convenient to map the model parameters into physical units. The comparison of simulation results with the Lennard Jones potential in Eq. (5) and experimental results with Argon provides one of the most established choices for simulation units, the so called *reduced units*. These consist in choosing σ as a unit of length and relate to $\sigma = 3.40510^{-10}$ m, ϵ as a unit of energy with $\epsilon/k_B = 119.8$ K, and the reference particle mass m_i as the unit of mass with $m_i = 0.03994$ kg/mol. Some related quantities are specified in Table 1.

Quantity	Reduced unit	Simulation value	Physical value
temperature	$T^* = k_B T / \epsilon$	$T^{*} = 1$	$T=119.8~{\rm K}$
time	$t^* = t\sqrt{\epsilon/m_i\sigma^2}$	$t^{*} = 0.01$	$t = 4.32210^{-10} \text{ s}$
pressure	$P^* = P\sigma^3/\epsilon$	$P^{*} = 1$	$P=41.9~\mathrm{MPa}$
density	$\rho^* = \rho \sigma^3$	$\rho^* = 0.4$	$\rho=960~{\rm kg/m^3}$

Table 1: Most relevant reduced units and their translation into physical units

The main importance of the reduced units is that many other parameter combinations translate into the same reduced units, such that their results should be equivalent. This is the *law of corresponding states*.

2.4 Non-equilibrium simulations

The systems under study can be made progressively more complex by increasing the number of particles involved, by distinguishing more than one type of particles, or by introducing many different interactions between particles. Extra degrees of complexity can be introduced by considering an applied external field like gravity or an electric field as stated in Eq. 2. In all these cases though the system will be studied in equilibrium conditions. In the last decades, the explosion of the computer power, and the improvement of the algorithms is making accessible to study by means of simulations the properties of complex systems in various non-equilibrium conditions.

Confinement When simulating a fluid confined between walls two types of boundary conditions can be identified by characterizing the velocity of the fluid in the proximity of the walls. The fluid may 'slide' at the wall and have therefore a different velocity. This characterize the so called *slip boundary conditions*. The *stick boundary conditions* though impose that the fluid has the same velocity as the walls in their neighbourhood.

A perfectly flat wall with slip boundary conditions can be implemented by considering a repulsive Lennard Jones interaction along the wall direction. A rough wall, with boundary conditions close to stick can be included by considering a random distribution of particles with fixed position (similar to infinity mass) that interact with he fluid for example with WCA interactions. Stick can also be obtained with bounce-back, here when a particle arrives to the wall sees its velocity inverted, such the trajectory is reverted at least in part.

Poisuelle flow If walls with stick boundary conditions are implemented together with an applied force0 parallel to the walls, a parabolic velocity profile will develop (see Fig. 4). These Poisuelle or capillary flows are an essential ingredient in many *microfluidic* systems where the geometry of the confinement can be made arbitrarily complex [12, 13].



Fig. 4: Illustration of a fluid confined between walls in one direction and periodic boundary conditions in the other. In the case of stick boundary conditions, and in the presence of an applied field, a capillary flow will be present.

Further non-equilibrium conditions Various other relevant non-equilibrium conditions can be considered, such as shear flow, external temperature gradients, applied external electric currents, or magnetic fields. When the system is in shear flow it displays a linearly varying velocity with position. This is a very interesting situation form the scientific viewpoint since it is possible to study the dynamical behaviour in a non-equilibrium steady state [12, 13]. The number of technological applications is also large ranging from food production to blood flow. Shear flow can be implemented in a MD code mainly following two strategies. The first one is to confine the system between parallel walls which are moving at different velocities. The other one is a modification of the periodic boundary conditions introduced by Lees and Edwards [14]. Effects induced by temperature inhomogeneities are present in a large number of technical and biological systems [15, 16]. To include such temperature gradients, an artificial energy transfer is typically imposed from a cold to a hot slab producing that these two slabs have indeed different average kinetic energies [17, 18]. The system in between the slabs will experience a physical linear temperature variation that will have opposite signs in the two simulation box halves.

2.5 Monte Carlo Simulations

Monte Carlo (MC) simulations focus in providing and efficient stochastic sampling of the configurational and conformational phase space or parts of it [1]. The time scale is therefore not naturally built and the objective of these simulations is to obtain good approximations for statistical quantities such an expectation values, probabilities, correlation functions, or densities of states. Note that MD simulations can become extremely slow when applied to complex systems on microscopic or mesoscopic scales, and that the understanding of many interesting questions does not require to consider the intrinsic dynamics of the system. In some of these systems, MC constitutes a really powerful tool. The main idea of the method consists in obtaining the thermodynamic properties of a system by stochastically sampling of the particle configurations proportional to their Boltzmann weight. The Boltzmann factor of a system with N particles with positions \mathbf{r}^N and temperature T is

$$p(\mathbf{r}^{N}) = \frac{e^{-E(\mathbf{r}^{N})/k_{B}T}}{\int d\mathbf{r}^{N} e^{-E(\mathbf{r}^{N})/k_{B}T}} = \frac{e^{-E(\mathbf{r}^{N})/k_{B}T}}{Q(N,V,T)},$$
(13)

where Q(N, V, T) is the partition function in the canonical ensemble. The averages of an observable $A(\mathbf{r}^N)$ which depends on the configuration \mathbf{r}^N can be calculated then as,

$$\langle A \rangle = \int d\mathbf{r}^N A(\mathbf{r}^N) p(\mathbf{r}^N).$$
(14)

The basic idea of MC is then that instead of calculating a weighted average over all (unweighted) configurations one can also calculate an un-weighted average over weighted configurations.

To illustrate the Monte Carlo method we first discuss the simplest Monte Carlo approach in the calculation of integrals. This method is called random or simple Monte Carlo sampling. Suppose one has an integral of the form $I = \int_a^b f(x) dx$. Using random sampling, the value of this integral is determined by evaluating f(x) at M randomly chosen values of x in the interval $\{a...,b\}$ and averaging over the corresponding f(x). The integral is obtained by $I = (b-a) \langle f(x) \rangle$. This method works for low dimensional integrals, and if the integrand behaves properly, i.e. is a smooth, slowly varying function of the variables.

However, for cases in which the integrand is a rapidly varying function of the configurations \mathbf{r}^N the previous approach is very inefficient. Moreover, the integrand is zero for most configurations and sharply peaked around the average value of the energy. This behaviour of the integrand can be understood because if we were to generate configurations randomly, then most of these configurations would contain overlapping pairs of particles for which the energy is large (infinite). These configurations would have a zero Boltzmann factor. Using the simple (random sampling) Monte Carlo scheme we would thus spend a lot of time in calculating zero contributions to the integral.

For these kind of situations it is necessary to employ importance sampling. In this method, the integrand is sampled often when the Boltzmann factor is large and less frequent when the Boltzmann factor is small. Efficient schemes are then developed to generate sequences of configurations proportional to their Boltzmann factor. Using these configurations an average is calculated in which the weighing is contained in the selection of configurations r^N . A new configuration is created from a preceding one according to a specified transition probability.

Various much more sophisticated sampling methods have been and are being developed and they are employed depending on the model details and relevant questions. Some examples of these methods are the Metropolis algorithm, Rosenbluth sampling, configurational bias sampling, parallel tempering, umbrella sampling, grand canonical MC, Gibbs ensemble MC or the Wang-Landau method.

3 Hydrodynamic simulation methods

The above microscopic techniques are appropriate for systems in which all the relevant component or subunits are of the same order of magnitude. However, this is not the case in most soft matter or biological systems where various components with sizes ranging from $10nm-10\mu$ m such as colloids, polymers, membranes, or vesicles, are mixed and coexist with a solvent whose component particles are molecules of approximately 0.2nm. To perform simulations in which the mesoscopic component and the surrounding solvent is accounted in detail is therefore far away from our computational possibilities. The aim to 'bridge the length and time scales gap', and the increasing availability of computing power, has stimulated the development of several mesoscale simulation techniques in recent years. This has been done fundamentally either from the 'top-down' approach which consists on discretizations of the continuum equations, or from the 'bottom-up' approach which consist in coarse grained descriptions of the fluid where the microscopic scale is strongly simplified, but relevant effects are still taken into account.

The most relevant effect of the solvent is to offer a resistance to any particle motion which will be proportional to the solvent viscosity η . Given the small size of the diluted particles this friction force will be typically accompanied by the thermal fluctuations, by which the particle will perform a random motion, due to the stochastic interaction with the solvent particles. In the cases, where these two contributions are the only relevant ones, a force \mathbf{F}_i applied to the solute particle will result in a linear increase of the particle velocity

$$\mathbf{v}_i = \frac{D_0}{k_B T} \mathbf{F}_i,\tag{15}$$

where k_B is the Boltzmann constant, T the system temperature, and D_0 the diffusion coefficient. D_0 is a constant that depends on the solute and solvent properties, and it is typically inversely proportional to the solvent viscosity η . But very frequently this description will be not sufficient dues to the presence of the solvent induced *hydrodynamic interactions*. In the presence of such hydrodynamic situations, the velocity increase of a particle due to an applied force is not linear anymore but depends on the position and velocities of all neighboring particles,

$$\mathbf{v}_i = \frac{1}{k_B T} \sum_j \mathbb{D}_{ij}(\mathbf{r}^N) \mathbf{F}_j,\tag{16}$$

where $\mathbb{D}_{ij}(\mathbf{r}^N)$ is the hydrodynamic mobility tensor or *Oseen tensor*. Hydrodynamic interactions are therefore dynamical many-body forces characterized by long range interactions that in three dimensions decay like the inverse particle distance, $\sim r^{-1}$. The proper description of such hydrodynamic interactions is fundamental in many systems and processes like polymer or colloid suspensions, microphase separation or microfluidics devices.



When one particle in solution moves, the fluid around is perturbed such that the movement of all neighboring particles is affected.

Fig. 5: Diagram of hydrodynamic interactions.

Following, most of the existing mesoscopic hydrodynamic simulation techniques are briefly mentioned as they have historically appeared. A closer description of the basic implementation details is given for only one of the most competitive methods, multiparticle collision dynamics (MPC). For further explanations, I refer the interested reader to more extended reviews [7, 19, 20, 21] and to the original literature.

3.1 Brownian dynamics (BD)

In BD the solvent particles are omitted from the simulation [22], and their effects upon the solute is represented by a combination of random forces and frictional terms. The BD method computes the space trajectories of a collection of particles that individually obey Langevin equations in a field of force, which reproduces the diffusive behavior of for instance colloidal dispersions. The main drawback of the method is that momentum is not a locally conserved quantity, which means that the behavior is not hydrodynamic but diffusive. Hydrodynamic interactions between particles can be incorporated through a tensorial dependence in the Langevin equations. The most common choice for such a dependence is the hydrodynamic Oseen tensor. This method has reproduced the behavior of polymer and colloidal dispersions, but the two main problems are that detailed time dependent information about the solvent is lost, and that the size of the considered tensor increases with the the number of particles, what limits the size of the systems which can be considered.

3.2 Smoothed particle hydrodynamics (SPH)

This technique was developed in the 1970s in the context of astrophysical flow problems [23, 24] and more recently it has been applied to the study of fluid dynamic problems like viscous [25] and thermal flows [26] in simple geometries. This new application is frequently named smoothed particle applied mechanics (SPAM). The essence is to discretize the macroscopic partial differential equations, such as the Navier-Stokes equations for Newtonian fluids or the elasticity equation for solids. The discretization takes place in an irregular and Lagrangian moving grid in such a way that the nodes can be interpreted as soft particles. In fact, the technique allows one to solve partial differential equations with molecular dynamic simulation codes. SPH has not been applied to study hydrodynamic problems where fluctuations are relevant, as those occurring in colloidal suspensions. Actually, there are some subtleties in discretizing the random stress and random heat flux that appear in the equations of fluctuating hydrodynamics of Landau and Lifschitz [10].

3.3 Lattice Boltzmann (LB)

The original idea was to replace the continuum macroscopic picture by a set of particles that move from site to site on a fixed regular lattice, colliding and changing their velocities according to certain collision rules, whose only restriction is that particle number, momentum and energy should be conserved quantities. This was the basics of lattice gas automata (LGA) [27, 28], which provided an adequate symmetry of the lattice, was a valid fluid-dynamical model yielding the correct Navier-Stokes hydrodynamics at a coarse-grained scale. However, fundamental problems were displayed by this technique, namely that isotropy and Galilean invariance were both broken by the lattice and large density fluctuations we appearing appear.

The lattice Boltzmann method was first developed empirically in 1988 [29] from LGA and it was introduced to circumvent its major shortcomings. Later it has been demonstrated that the LBE can be directly derived from the continuous Boltzmann equation [30]. For review of the method see [31, 32, 21]. The main purpose is to incorporate the physical nature of fluids from a more statistical point of view. Similar to LGA, LB consists in a set of particles that move in a space restricted to the nodes of a regular lattice. Particle distributions propagate from node to node, with a rate proportional to a discrete velocity c_k and interchange mass and momentum

with other particle distributions in the corresponding node before the next propagation step. Such mass and momentum interchange is performed according to an underlying picture of the Boltzmann transport equation, and the dynamics is implemented directly to the distribution functions instead on the individual particles. In the LB standard implementation mass and momentum are imposed to be conserved quantities what give the basics of the hydrodynamic behavior. In order to reproduce the Navier-Stokes equation the conservation of the stress tensor is required as well.

LB constitutes a much more efficient method than its precursor LGA and it is extensively employed by a large community. It has proved to be especially useful in studying flows in complex geometries, like porous media, or the dynamics of colloidal suspensions, as well as in studies of multicomponent systems. The model suffers, however, of some intrinsic problems. One is that energy conservation is not fulfilled by the present LB method, such that it is generally restricted to isothermal applications, although there are recent energy conserving generalizations of the model. Another problem is that 'thermal fluctuations' are not present in the model, while they are in LGA due to the discreteness in the number of particles. This lack of fluctuations can be in some cases an advantage since these are always a source of statistical inaccuracy. But thermal fluctuations are a required element in the correct description of a large number of physical problems, like the Brownian motion of suspensions or in the decay of spontaneous stress fluctuations. For these cases, a possible solution has been proposed based on linear fluctuating hydrodynamics according to which the stress tensor should include a noise term. This can therefore be directly added in the simulation code [33].



Fig. 6: Interaction sketches in 2 dimensions for the LB velocity vectors(left) and for the DPD particles (right).

3.4 Dissipative particle dynamics (DPD)

The DPD model was introduced in 1992 [34, 35, 36] as an attempt to free the LGA from the lattice. The first idea of DPD is to consider soft and finite interactions into a standard simulation with MD. The scales of time and space that can be reached are quite large, and phenomena related to processes on mesoscopic scales can be reproduced. The second idea of DPD is that these soft and finite interactions have dissipative and stochastic contributions, as well as a weak conservative term. The introduction of these dissipative and random interactions between DPD particles can be understood if each particle is representing not only one molecule, but rather a

group of them. Therefore, a DPD particle models the center of mass of a mesoscopic portion of the fluid, large enough to be a thermodynamic subsystem, but still subjected to thermal fluctuations [37, 38]. The state of the fluid is described by N particles with continuous positions velocities. Similar to molecular dynamics (MD), the particles time evolution is given by the integration of the Newton's equation of motion.

 \mathbf{v}_i of the particles are of Langevin type

$$d\mathbf{r}_{i} = \mathbf{v}_{i}dt,$$

$$d\mathbf{v}_{i} = \frac{1}{m}\sum_{j\neq i} \left(\mathbf{F}_{ij}^{C} + \mathbf{F}_{ij}^{D} + \mathbf{F}_{ij}^{R}\right) dt,$$
 (17)

where i ($i = 1, \dots, N$) labels the particles and m is the mass of a particle. The pair force that particle j exerts on particle i has three contributions: a conservative force \mathbf{F}_{ij}^{C} , a dissipative force \mathbf{F}_{ij}^{D} , and a random force \mathbf{F}_{ij}^{R} . These forces are interpreted as coarse grained averages over microscopic degrees of freedom.

The dissipative and random forces combined act as a thermostat. The dissipative force is proportional to a friction constant and cools the system, whereas the random force heats it up. To qualify as a fluid, DPD should be Galilean invariant and isotropic. The Galilean invariance requires that the forces depend only on relative variables $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$ and $\mathbf{v}_{ij} = \mathbf{v}_i - \mathbf{v}_j$. Isotropy requires that the forces transform under rotations as vectors. Moreover, the drift term of the Fokker-Planck equation must be linear in the velocity variable, and the diffusion term independent of it [39]. These requirements are satisfied if the dissipative force F^D is linear in the velocities and the random force F^R is independent of the velocity. A simple form of the forces satisfying these criteria is

$$\frac{1}{m} \mathbf{F}_{ij}^{C} dt = a w_{C}(r_{ij}) \hat{\mathbf{r}}_{ij} dt$$

$$\frac{1}{m} \mathbf{F}_{ij}^{D} dt = -\gamma w_{D}(r_{ij}) (\hat{\mathbf{r}}_{ij} \cdot \mathbf{v}_{ij}) \hat{\mathbf{r}}_{ij} dt,$$

$$\frac{1}{m} \mathbf{F}_{ij}^{R} dt = \sigma w_{R}(r_{ij}) \hat{\mathbf{r}}_{ij} dW_{ij},$$
(18)

where r_{ij} and $\hat{\mathbf{r}}_{ij}$ are respectively the modulus and the unit vector parallel to \mathbf{r}_{ij} . The coefficients a, γ and σ are positive constants that control the strength of the conservative repulsion, the friction and noise amplitudes. The range interaction functions w are bounded positive functions of the relative distance r_{ij} , and vanish for $r > r_c$.

The physical interpretation of the dissipative force \mathbf{F}_{ij}^D is as follows. When particles *i* and *j* are approaching/receding, the quantity $(\hat{\mathbf{r}}_{ij} \cdot \mathbf{v}_{ij})$ is negative/positive which implies that both particles feel a viscous force slowing down their relative motion in the $\hat{\mathbf{r}}_{ij}$ direction. The random force is a physical consequence of the mesoscopic description, and account for thermal fluctuations. This term describes a Gaussian white noise such that $\langle dW_{ij}(t) \rangle = 0$. The symmetry property $\mathbf{F}_{ij} = -\mathbf{F}_{ji}$ ensures that the total momentum is conserved, and enforces that $dW_{ij} = dW_{ji}$. The noise term is interpreted through the Ito calculus rule, and normalized as

$$dW_{ij}dW_{i'j'} = (\delta_{ii'}\delta_{jj'} + \delta_{ij'}\delta_{ji'})dt,$$
(19)

i.e. $dW_{ij}(t)$ is an infinitesimal of $\mathcal{O}(dt^{1/2})$ [40]. This form respects the symmetry under particle interchange. Furthermore, the relation between the dissipative and the random forces

has a precise form in order that the system in equilibrium displays the Maxwell-Boltzmann distribution. This implies the fluctuation dissipation relations [35, 37]

$$m\sigma^2 = 2\gamma k_B T$$
 , $w_D(r) = w_R^2(r) \equiv w(r)$, (20)

where a single range function has been defined.

The technique has been constructed such that both the number of particles and the total momentum are conserved quantities. Therefore, there is a transport equation for the momentum density field, coupled to the continuity equation. The macroscopic behavior of this particle model is then hydrodynamic, and not just diffusive as it occurs with Brownian dynamics.

DPD was first applied to analyze rheological properties of colloidal suspensions and polymer solution. It has also been applied to more complex situations, like microphase separation, dynamics of a drop at a liquid-solid interface, flow and rheology in the presence of polymers grafted to walls, colloidal adsorption onto polymer coated surfaces, amphiphilic mesophases, model membranes or geometrical packing of filler in composites. Later extensions of the original DPD model have been proposed in order to include various aspects. The inclusion of an internal energy variable [41, 42, 43] allows define an algorithm in which energy is conserved global and locally. Considering a spin variable, a version on which angular momentum is conserved is formulated [44]. Another generalization has been introduced for modeling viscoelastic flows [45] by including an elastic variable for each fluid.

3.5 Multiparticle collision dynamics (MPC)

The most recent mesoscale simulation technique treated in this chapter was first introduced by Malevanets and Kapral [46, 47] in 1999. The main idea was to modify a method largely employed for the simulation of gas flows the so-called Direct Simulation Monte Carlo method (DSMC), replacing binary collision by multi-particle collisions in a prescribed collision volume. This method has been called multiparticle collision dynamics (MPC) or stochastic rotation dynamics (SRD).

The MPC fluid is modeled by N point particles. Each of these particles is characterized by its position \mathbf{r}_i and velocity \mathbf{v}_i , and labeled with i = 1, ..., N. Positions and velocities are continuous variables, which evolve in discrete increments of time. The mass m_i associated with the particles is usually taken to be the same, but more generally, different masses can be assigned. The MPC algorithm consists of two steps, streaming and collision, which are illustrated in Fig. 7. In the streaming step the particles do not interact with each other (see Fig. 7a), they move ballistically according to their velocities during a time increment h, to which I will refer as *collision time*. Thereby, the evolution rule is

$$\mathbf{r}_i(t+h) = \mathbf{r}_i(t) + h\mathbf{v}_i(t). \tag{21}$$

In the collision step, the particles are sorted into collision boxes (see Fig. 7b), and interact with *all* other particles in the same collision box. This multibody interaction takes place through the collision box center of mass velocity. This is

$$\mathbf{v}_{cm,i}(t) = \frac{\sum_{j}^{(i,t)} (m_j \mathbf{v}_j)}{\sum_j m_j},\tag{22}$$

such that $\mathbf{v}_{cm,i}(t)$ is the velocity of the center of mass of all particles j, which are located in the collision box of particle i at the considered time t. The collision boxes are typically the unit



Fig. 7: *Diagram of the MPC dynamics in 2 dimensions. (a) Streaming step, (b) particles sorted into collision boxes and (c) rotation of the particle velocity relative to the center of mass.*

cells of a d-dimensional cubic lattice with lattice constant a, although other geometries would be possible. The collision is then defined as a rotation of the velocities of all particles in a box in a co-moving frame with its center of mass. Thus, the velocity of the i-th particle after the collision is

$$\mathbf{v}_{i}(t+h) = \mathbf{v}_{cm,i}(t) + \mathcal{R}(\alpha) \left[\mathbf{v}_{i}(t) - \mathbf{v}_{cm,i}(t)\right],$$
(23)

where $\mathcal{R}(\alpha)$ is a stochastic rotation matrix. The rotation by a fixed angle α , occurs around an axis which is stochastically chosen in each collision. This implies that each particle changes during the collision the magnitude and the direction of its velocity (see Fig. 7c), in such a way that the total momentum and kinetic energy are still conserved within every collision box, together with the effect of a Gaussian thermal noise. Furthermore, and in order to ensure Galilean invariance, a random displacement of the collision grid is performed before each collision step [48, 49]. This grid displacement also softens the particle interaction, and enhances collisional transport and the fluid-like behavior of the MPC fluid. The MPC transport properties will then be determined by the most relevant model parameters which are typically h, the collision time, α , the rotation angle, and ρ , the average number of particles in a collision cell. The main drawback of MPC is the the fluid is constrained by the ideal gas equation of state, which makes it difficult to apply for example to multicomponent flows, but it is not a problem in most other applications of the method.

The MPC fluid has been already successfully applied to a large range of soft matter and biology problems where hydrodynamic interactions are relevant, such as polymers or liquid crystals in external flows [50, 51], active matter [52, 53], sperm dynamics [54], or red blood flow circulation [55], and its use is currently expanding in the community.

4 Conclusions

Computer simulation methods are nowadays a very powerful, versatile, and essential tool to help to unravel the intricated physical and biological mechanisms of nature, as several example of this are provided in other chapters of this book. Microscopic simulation methods such as Molecular Dynamics and Mote Carlo are very versatile, and still very broadly employed to study the properties of a large spectrum systems and conditions. The most significant limitations of these methods is that the times and sizes accessible by simulations is limited, although it is getting progressively more realistic. With MD the system can also eventually be trapped in a particular area of the phase space, this is when typical relaxation times that would drive the system out of such state, are larger than the accessible computing times. Ergodicity would in this case not be satisfied. Another limitation is the difficulty of bridging different length and time scales that can be present in the same system, for example in system in solution where the scale of the solute and the solvent are usually order of magnitude apart. This is solved by several hydrodynamic simulation methods where the solvent degrees of freedom are very efficiently treated. Dissipative particle dynamics, Lattice Boltzmann, and multiparticle collision dynamics offer three very different approaches to treat the solvent dynamics. The particular characteristics of the problem under study determine which of the methods is more appropriate. In this way multiphase flows are extensively treated with Lattice Boltzmann methods, the presence of temperature gradients with multiparticle dynamics, and fluids with non-ideal equations of state or transport amphiphilic systems with dissipative particle dynamics. All these techniques have been improved from their original formulation, but they are still being constantly tested and generalized such that new aspects can be included. A few comparative studies between different techniques have been performed, but there is no systematic ultimate study which would apply to all possible systems. This also shows that this is a lively and promising field.

A first glimpse over the definitions and possibilities of computer simulations is offered in this chapter. The most important considerations required for an interested reader without any previous experience in simulations have been introduced, such that further understanding and implementation details are directed to more specialized literature.

References

- N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller, J. Chem. Phys. 21, 1087 (1953).
- [2] B. J. Alder and T. E. Wainwright, J. Chem. Phys. 27, 1208 (1957).
- [3] A. Rahman, Phys. Rev. 136, A 405 (1964).
- [4] M. Levitt and A. Warshel, Nature 253, 94 (1975).
- [5] R. Car and M. Parrinello, Phys. Rev. Lett. 55, 2471 (1985).
- [6] M. P. Allen and D. J. Tildesley, *Computer Simulations in Liquids* (Clarendon, Oxford, 1987).
- [7] D. Frenkel and B. Smit, *Understanding Molecular Simulation: From Algorithms to Applications* (Academic Press, San Diego, 2002), 2nd ed.
- [8] D. J. Evans and G. P. Morriss, *Statistical Mechanics of NonEquilibrium Liquids* (Academic Press, London, 1990).
- [9] D. W. Heermann, *Computer Simulations Methods in Theoretical Physics* (Springer, Berlin, 1990).
- [10] L. D. Landau and E. M. Lifshitz, Fluid Mechanics (Pergamon Press, 1959).
- [11] J. P. Hansen and I. Mc Donald, Theory of simple liquids (Academic, New York, 1986).
- [12] R. G. Larson, *The Structure and Rheology of Complex Fluids* (Oxford University Press, New York, 1999).
- [13] M. Doi and S. F. Edwards, *The Theory of Ploymer Dynamics* (Oxford University Press, Oxford, 1986).

- [14] A. W. Lees and S. F. Edwards, J. Phys. C 5, 1921 (1972).
- [15] S. Wiegand, J. Phys.: Condens. Matter 16, R357 (2004).
- [16] C. J. Wienken, P. Baaske, U. Rothbauer, D. Braun, and S. Duhr, Nature Comm. 1, 100 (2010).
- [17] B. Hafskjold, T. Ikeshoji, and S. K. Ratkje, Molecular Phys. 80, 1389 (1993).
- [18] F. Müller-Plathe, J. Chem. Phys. 106, 6082 (1997).
- [19] R. Kapral, Adv. Chem. Phys. 140, 89 (2008).
- [20] G. Gompper, T. Ihle, D. M. Kroll, and R. G. Winkler, Adv. Polym. Sci. 221, 1 (2009).
- [21] U. D. Schiller, T. Krüger, and O. Henrich, Soft Matter 14, 9 (2018).
- [22] W. W. Wood and F. R. Parker, J. Chem. Phys. 27, 720 (1957).
- [23] L. B. Lucy, Astron. J. 82, 1013 (1977).
- [24] J. J. Monaghan, Annu. Rev. Astron. Astrophys. 30, 543 (1992).
- [25] H. A. Posch, W. G. Hoover, and O. Kum, Phys. Rev. E 52, 1711 (1995).
- [26] W. G. Hoover and H. A. Posch, Phys. Rev. E 54, 5142 (1996).
- [27] J. Hardy, Y. Pomeau, and O. de Pazzis, J. Math. Phys. 14, 1746 (1973).
- [28] U. Frisch, B. Hasslacher, and Y. Pomeau, Phys. Rev. Lett. 56, 1505 (1986).
- [29] G. R. McNamara and G. Zanetti, Phys. Rev. Lett. 61, 2332 (1988).
- [30] X. He and L. S. Luo, Phys. Rev. E 56, 6811 (1997).
- [31] S. Succi, *The Lattice Boltzmann Equation: for fluid dynamics and beyond* (Clarendon, Oxford, 2001).
- [32] R. R. Nourgaliev, T. N. Dinh, T. G. Theofanous, and D. Joseph, Int. J. Multiphase Flow. 29, 117 (2003).
- [33] A. J. C. Ladd, H. Gang, J. X. Zhu, and D. A. Weitz, Phys. Rev. E 52, 6550 (1995).
- [34] P. J. Hoogerbrugge and J. M. V. A. Koelman, Europhys. Lett. 19, 155 (1992).
- [35] P. Español and P. Warren, Europhys. Lett. 30, 191 (1995).
- [36] R. D. Groot and P. B. Warren, J. Chem. Phys. 107, 4423 (1997).
- [37] C. Marsh, G. Backx, and M. Ernst, Phys. Rev. E 56, 1976 (1997).
- [38] C. Marsh, G. Backx, and M. Ernst, Europhys. Lett. 38, 411 (1997).
- [39] H. Risken, The Fokker Planck Equation (Springer, Berlin, 1989).
- [40] C. W. Gardiner, Handbook of Stochastic Methods (Springer, Berlin, 1983).
- [41] J. Bonet-Avalós and A. D. Mackie, Europhys. Lett. 40, 141 (1997).
- [42] P. Español, Europhys. Lett. 40, 631 (1997).
- [43] M. Ripoll and M. H. Ernst, Phys. Rev. E 71, 041104 (2005).
- [44] P. Español, Phys. Rev. E 57, 2930 (1998).
- [45] B. I. M. ten Bosch, J. Non-Newtonian Fluid Mech. 83, 231 (1999).
- [46] A. Malevanets and R. Kapral, J. Chem. Phys. 110, 8605 (1999).
- [47] A. Malevanets and R. Kapral, J. Chem. Phys. 112, 7260 (2000).
- [48] T. Ihle and D. M. Kroll, Phys. Rev. E 63, 020201(R) (2001).
- [49] T. Ihle and D. M. Kroll, Phys. Rev. E 67, 066705 (2003).
- [50] M. Ripoll, P. Holmqvist, R. G. Winkler, G. Gompper, J. K. G. Dhont, and M. P. Lettinga, Phys. Rev. Lett. 101, 168302 (2008).
- [51] C. Huang, G. Sutmann, G. Gompper, and R. Winkler, Europhys. Lett. 93, 54004 (2011).
- [52] A. Wysocki, R. G. Winkler, and G. Gompper, Europhys. Lett. 105, 48004 (2014).
- [53] M. Wagner and M. Ripoll, EPL 119, 66007 (2017).
- [54] J. Elgeti, U. B. Kaupp, and G. Gompper, Biophys. J. 99, 1018 (2010).
- [55] J. L. McWhirter, H. Noguchi, and G. Gompper, Proc. Natl. Acad. Sci. 106, 6039 (2009).

B1 Protein Folding and Protein Stability

J. Fitter Institute of Complex Systems: Molecular Biophysics (ICS-5) Forschungszentrum Jülich GmbH

A. Stadler Institute of Complex Systems: Neutron Scattering (ICS-1) Forschungszentrum Jülich GmbH

Contents

1	Introduction		
2	Ear	ly experiments on protein folding	2
3	Con	cepts and models in protein folding	4
	3.1	The thermodynamic hypothesis	4
	3.2	Energy landscapes, intermediate states, and barrier-less folding	5
	3.3	Protein folding kinetics	7
4	Clas	ssical experimental techniques in protein folding	8
5	The Stat	Relevance of Conformational Entropy for Protein Folding	g and 10
	5.1	Myoglobin: A Model System for Protein Folding Studies	10
	5.2	Conformational Entropy Measured with Neutron Spectroscopy	11
Ref	erence	S	14

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

It is still one of the great unsolved problems in life science that we do not understand in detail how an amino acid sequence of a polypeptide chain is transformed into a well defined threedimensional protein structure (also called "folding problem"). This fundamental process and the attempts to establish a more complete understanding of the so called protein folding, is a prominent topic in the field of biophysical chemistry and biochemistry. In addition the analysis of protein folding and stability has also become more and more relevant for medicine and biotechnology. There are at least three major reasons that an increasing theoretical and experimental interest in protein folding exists¹: (1) Due to efforts on genome sequencing projects the acquisition of DNA sequences is increasingly faster nowadays. Compared to that, the acquisition of three-dimensional proteins structures is still slow, limited by the time consuming process of searching for proper crystallization conditions. In this respect the knowledge how the linear sequence of amino acids is translated into spatial information is the "missing link". (2) There is a tremendous interest in the over-expression of recombinant proteins for industrial, biotechnological, and research applications. (3) Incorrect folding or misfolding of proteins is often related to protein aggregation and fibrillogenesis, which is connected to a number of serious diseases, such as BSE, or Huntington's and Alzheimer diseases. Essentially the folding problem can be subdivided in three basic problems². First, there is the question of the folding code. What are the forces that dictate the protein structure for a given amino acid sequence? Second, there is the challenge of a protein structure prediction. How can we predict a native three-dimensional protein structure from its amino acid sequence? And a third aspect deals with the understanding of the folding process itself. How are the routes or pathways during protein folding and why does the folding happens in some cases very fast and in others much slower (folding kinetics)? All these aspects are still intensively under investigation which is also supported by the fact that the number of research reports and review articles on various aspects of protein folding has increased dramatically over the last decades².

Herein we recapitulate first early major milestones in protein denaturation and protein folding research (Sec. 2) which started almost hundred years ago. In the following section (Sec. 3) we present some seminal concepts and models in protein folding which have been established in the last decades. In Section 4 we have a closer look on classical techniques that are used to monitor structural properties and the kinetics of folding and unfolding transitions. In the last two sections we discuss shortly the potential of using single molecule techniques in protein folding (Sec. 5) and the difference between *in vitro* and cellular (co-translational) protein folding.

2 Early experiments on protein folding

The process of protein folding is already known for more than eighty years. The view that protein denaturation is an unfolding process and not a hydrolysis of peptide bonds or a dehydration of proteins was first proposed 1929 (see Table 1)⁴. A real starting point of modern work on protein folding was the confirmation of the "thermodynamic hypothesis" through the seminal studies of Christian Anfinsen and his colleagues⁵ (see Figure 1). From his experiments he concluded that the native structure of a protein is the thermodynamically stable structure. As a consequence one can deduce that in general the native structure does not depend on whether the protein was synthesized on a ribosome and subsequently folded with the help of chaperones (cellular folding) or, if instead the protein was refolded from a



Fig. 1: Schematic representation of the disulfide-cross-linked RNase unfolding and successful refolding under reductive conditions. Here the protein spontaneously finds the only correct pairing out of 105 possible pairings of eight sulfhydryl groups that form four disulfide linkages⁵.

previously unfolded protein in aqueous buffer in a test tube (*in vitro* folding). This important finding paved the way for routine protein folding studies inside test tubes rather than studies inside cells, where the latter are by far more difficult. A subsequent achievement of early studies on protein folding was the discovery that small proteins in general show a two-state folding reaction without observable intermediate states. In order to obtain a true two-state folding it was important to work with a fully unfolded state which was reached at high concentrations of a strong denaturant, i.e. 6 M guanidine hydrochloride (GndHCl)⁶. Due to the fact that the process is in most cases fully reversible the application of equilibrium thermodynamics provided a wealth of valuable insights into the mechanisms of proteins folding and stability⁷. In the late sixties of the previous century Cyrus Levinthal pointed out that if the folding would proceed without intermediate states, then the time needed to search randomly all possible backbone conformations is far longer than the life of the universe (Levinthal's paradox). A reasonable solution to this problem is given by the assumption that there must be folding intermediates and pathways, even if they are not (easily) observable or only short-lived. Therefore, starting from the early seventies in the last century large effort was put into the search for folding intermediates and with the advent of fast kinetic methods such intermediates were finally observed⁸⁻¹⁰. In this respect the observations of Oleg Ptitsyn and co-workers on acid treated α -lactalbumin were ground-breaking. They found a special intermediate that was characterized as a "molten globule" state with a native like compact structure, a high content of secondary structure elements, but with poor tertiary structure¹¹. Such structural properties were also observed in transient intermediate states during the folding of certain proteins, especially globular proteins that undergo a hydrophobic collapse. More recent models and concepts in protein folding, in particular the so-called "new view of folding" in terms of energy landscapes, will be discussed in more detail in the next section.

Year	Principle Investigators and the related findings	
1902	E. Fischer, F. Hofmeister: Proteins are chains with covalently linked amino	
	acids	
1920-1940	H. Wu, L. Pauling, A. Mirsky: Protein denaturation is a purely conformational	
	change, it is in principle reversible. The denaturation corresponds to an	
	unfolding process and not to some chemical alterations of the protein.	
~1950	50 Neurath, R. Lumry, H. Eyring: Protein denaturation is thermodynamical	
	reversible which leads to the "Thermodynamic hypothesis".	
1956-1961 C. B. Anfinsen: The amino acid sequence of ribonuclease A contains		
	information needed to make the correct four disulphide bonds and the correct	
	three-dimensional structure of the protein (confirmation of the	
	"Thermodynamic hypothesis")	
1958-1961	J. Kendrew, M. Perutz: First reported three-dimensional protein structure	
	(myoglobin, haemoglobin)	
~1968	C. Levinthal: "It seems to be impossible that an unfolded protein can fold	
	spontaneously by a random process on a biological time scale" (Levinthal's	
	paradox). Consequence: folding is a hierarchical and stepwise process!	
~1981	<i>O. Ptitsyn:</i> He and his co-workers reported that the acid form of α -lactalbumin	
	is compact, with almost fully developed secondary structure, but with missing	
	the tertiary structure ("molten globule" state as a folding intermediate).	

Table 1: Milestones and breakthroughs of protein denaturation and folding

3 Concepts and models in protein folding

3.1 The thermodynamic hypothesis

The simplest scenario for a transition from the native to the unfolded state is characterized by a single energy barrier separating both states (see Fig. 2). Experimental evidence for pure two-state behavior is visible through mono-exponential unfolding or refolding kinetics. The application of equilibrium thermodynamics requires full reversibility of the transition reactions⁶ (F \leftrightarrow U). This requirement is fulfilled for the majority of smaller single-domain proteins, while large multi-domain proteins often exhibit irreversible unfolding transitions. The latter suffer in many cases from distinct aggregation of the unfolded state, which is in competition to a proper refolding¹². The folded protein is typically stabilized by the free energy ΔG^{0} :

$$\Delta G^0 = G^U - G^N = \Delta H - T \Delta S \tag{1}$$

As shown in equation 1, the protein stabilization is determined by entropic as well as by enthalpic contributions. While enthalpic contributions often play a role through intermolecular bonds (H-bonds, S-S-bonds, salt bridges) the entropy change upon folding is dominated by the conformational freedom of the protein structure and by the interaction of the protein with hydration water. A protein structure is controlled by a subtle balance of stabilizing and destabilizing forces, and the resulting energy which shifts the equilibrium to the native state (under physiological conditions) is often not larger than 1-4 kJ/mol (i.e. a few *RT* at 25 °C). An illustrative example of opposing "forces" which determine the protein structure is given by the entropy. On the one hand, the conformational entropy of the unfolded

state, where the structure exhibits the highest degree of conformational freedom, is much larger as compared to the native state. The latter has much more constraints and therefore conformational entropy stabilized the unfolded state $(S_{unf} > S_{fol} \rightarrow \Delta S > 0)$. On the other hand the entropic contribution of the hydration water is larger for the native state than for the unfolded state. Since the native state typically exposes a smaller number of hydrophobic residues to the solvent, less water molecules are restricted in their mobility as compared to the unfolded state. As a consequence the hydration entropy stabilizes the native state $(S_{unf} < S_{fol} \rightarrow \Delta S < 0)$.



Fig. 2: *Free energy profile for a two-state reaction where the native* N *and the unfolded* U *state are separated by a free energy barrier (transition state* T).

All the above mentioned considerations hold for the case of reversible processes. The height of the transition state barrier (T) determines how fast the system will reach equilibrium. In the case of irreversible unfolding transitions (F \rightarrow U) where equilibrium thermodynamics is not applicable, only information on the unfolding kinetics and thereby information on the kinetic stability is accessible. In this case the barrier height of the transition state (ΔG_N^{0T}) is the relevant parameter^{13, 14}.

3.2 Energy landscapes, intermediate states, and barrier-less folding

The folding of a protein involves formation and breakage of a multitude of intra-molecular contacts. Due to this fact a single-trajectory view (i.e. one more or less well defined route from the unfolded to the native state) of protein folding is an over-simplification of the folding process at a molecular level. A more realistic view is given by a perception of protein folding that is characterized by a myriad of many different folding routes. Such a picture emerges from a multidimensional rugged energy landscape with a funnel like shape (see Fig. 3) in which individual molecules find different routes from the unfolded to the native state¹⁵. Different routes might be probable, since the unfolded state is structurally rather heterogenic and individual unfolded structures represent slightly different starting conditions for the folding process. Related to this we can also expect different populations of potential intermediate states. For proteins that would fold without observable intermediates (i.e. long-

lived states) the energy surface would be rather smooth. The relevance of intermediates in the folding process is still in debate. In some cases intermediate states represent important structures with a high content of native contacts which are assumed to be productive to reach the final native state ("on-pathway intermediates"). In other cases mainly non-native interactions are trapped and therefore theses states are called "off-pathway" intermediates¹⁶. However, also "off-pathway" intermediates are of interest since they reflect properties of the energy landscape and might represent substrates for chaperones (folding helper proteins) or play a role for protein aggregation in the cellular context³.



Fig. 3: Schematic presentation of a folding energy landscape (folding funnel). Unfolded states of a protein are characterized by a higher level of free energy (shown in red at the top of the funnel). Some of the proteins are trapped transiently in intermediate states (shown in green in local energy minima) before they reach the one and only global minimum of the native state (shown in cyan). Figure taken from reference³.

Recently the idea of a folding process without barriers between the native and the unfolded state (see Fig. 2) has received a lot attention¹⁷. Elucidating the origin or quantifying the magnitude of this free energy barrier is not an easy task. This problem is caused by the fact that Evring's rate equation describes only gas phase reactions and is not valid (in particular the fundamental pre-exponential factor; for details see Section 3.3) for reactions in the condensed phase^{1,18}. However, many proteins exhibit clear experimental evidence for a twostate process in which only the fully unfolded and the native states are populated. In contrast to this, recently so-called "strange" kinetics (e.g., stretched exponentials) have been observed which is not consistent with a single exponential or a sum of a few exponentials¹⁹. Such kinetics cannot be explained by a single energy barrier but instead indicates very small or even absent barriers. In this case the so-called "downhill" folding is typically very fast and proceeds through an array of temporary conformations with a broad distribution of characteristic times^{20, 21}. For experimentalists it is still challenging to probe rapid transitions of protein subpopulations which may proceed on different folding routes. Recent experimental developments, such like ultra-rapid mixing techniques, opened the door to fast measurements and provide a closer connection with theory (see Section 4).

3.3 Protein folding kinetics

As we have already mentioned in the previous sections, the knowledge about the kinetics of the folding process is extremely valuable for the understanding of the energy landscape, for example in terms of local saddle points (transition states) or of local minima (intermediates). Assuming a simple two-state process the free energy barrier separating the native and the unfolded state is essentially rate limiting the process. The unfolding rate constant k_u is given by:

$$k_{\mu} = k_0 \cdot e^{-\Delta G_N^{0T} / RT} \tag{2}$$

In this equation k_0 represents the pre-exponential factor and ΔG_N^{0T} the free energy of the unfolding barrier (see Figure 2). Examples of unfolding kinetics, induced by chemical denaturants (e.g. GndHCl), which follow a mono-exponential behavior (see equation 2) are shown in Figure 4. In a similar manner the folding rate constant k_f is related to the free energy barrier height ΔG_U^{0T} . As demonstrated in a so-called Chevron-plot (left panel in Fig. 4) the dependence of k_u on the concentration of GndHCl is described by the following relation:

$$k_u = k_u^w \cdot e^{m_u \cdot D} \tag{3}$$

Here k_u^w is the unfolding rate in buffer without GndHCl, *D* is the molar denaturant activity (which is proportional to denaturant concentration), and the m_u -value describes the sensitivity of changes in unfolding rates with the denaturant concentrations. The m_u -values, which were obtained from Chevron-plots, are assumed to be proportional to changes in the accessible surface area (ASA) between the native and the transition state.



Fig. 4: Left: The unfolding kinetics has been measured for a protein at different concentrations of denaturant (GndHCl). The highest denaturant concentration corresponds to the fastest unfolding transition (blue solid symbols). All time-resolved unfolding transitions are fitted reasonably well with a single exponential. Right: In the Chevron-plot the dependence of unfolding (fitted by the blue line) and refolding (fitted by the red line) rates on the denaturant concentrations are shown. Unfolding (refolding) rates of the protein in the absence of denaturant can be extrapolated (see intersection of the linear fitting curve with the ordinate).

Chevron-plots are furthermore rather sensitive for deviations from pure two-state behavior. If an intermediate is populated (either in folding or in refolding) the Chevron-plot typically exhibits a curvature in the respective limb²². In addition to information on the free energy barriers kinetic data also provide a measure of the equilibrium free energy for the folded state (see equation 1 and Fig. 2):

$$\Delta G^0 = G_U^{0T} - G_N^{0T} = RT \cdot \ln(k_f^w / k_u^w) \tag{4}$$

With the knowledge of unfolding and refolding rates valuable information on the stabilization mechanisms of proteins can be deduced. Stabilization of a protein from a thermophilic organism with respect to its mesophilic homologue can be either kinetic or thermodynamic, as demonstrated in Figure 5.



Fig. 5: The effective stability of a protein can be increased either by increasing the difference between the free energy levels of the native (N) and the unfolded (U) state (thermodynamic stabilization) or by increasing barrier height of the transition state (T) (kinetic stabilization).

4 Classical experimental techniques in protein folding

The starting point for almost all *in vitro* folding studies is to accumulate the protein in the unfolded state. In order to obtain a large population of unfolded states, typically high concentrations of chemical denaturants (e.g., GndHCl, urea), low or high pH, or extreme temperatures as well as high pressure is used. An important first step in understanding protein folding is to characterize the conformational properties of the unfolded states. These states are highly dynamic with rapidly inter-converting species of similar energy (see Fig. 3). As known from low resolution data, unfolded states are less compact and have only very little native contacts as compared to the native state. Such results emerge from the application of various different methods employed in protein folding (Tab. 2). For example the hydrodynamic radius or the radius of gyration (results from dynamic light and small-angle X ray scattering) is much larger for unfolded states. A similar observation is made by tryptophan fluorescence measurements. Tryptophan residues that are buried in the interior of the native state structure, become solvent exposed and exhibit a decrease in the emission intensity and a typical red shift of the spectrum upon unfolding (see Fig. 6). Obviously protein unfolding in many cases

Technique	Timescale	Structural parameter probed
Fluorescence	ns-s ^a	
(i) Intrinsic		Environment of Trp and Tyr
(ii) FRET		Inter-residue distance
(iii) Anisotropy		Rotational mobility of the polypeptide
(iv) ANS binding		Exposure of hydrophobic surface area
Circular dichroism	ns-s ^a	
(i) Far UV		Secondary structure formation
(ii) Near UV		Environment of aromatic residues or co-factors
Infra-red	ns-s ^a	Secondary structure formation
spectroscopy		
Real-time NMR	~s	Information on the environment of individual
		residues
Pulsed Hydrogen	ms-s	Hydrogen-bond formation in specific amino-acid
exchange (NMR)		side-chains
Dynamic light	~min	Dimension in terms of hydrodynamic radius of the
scattering		polypeptide chain
Small-angle X-ray	> ms	Dimension and shape of the polypeptide chain
scattering		

is related to structural expansion of the polypeptide chain, while protein folding is characterized by a distinct structural compaction ("hydrophobic collapse").

Table 2: Classical experimental techniques used to measure protein folding. All techniques are described in more detail, for example in the Protein Folding Handbook edited by Buchner and Kiefhaber¹. ^aThe effective timescale depends on the method used to initiate folding: temperature jump (ns); ultra-fast mixing (μ s), stopped flow (ms), manual mixing (s).

CD - and infra-red spectroscopy allow to measure the content of secondary structure elements of a protein. For most proteins the disappearance of helical and beta-sheet structure happens concurrently with the three-dimensional structure expansion of the protein, which indicates a rather cooperative unfolding process (Fig. 6). Interestingly the investigation of unfolded states with various techniques revealed that often a residual native structure remains under unfolding conditions. The content of this native-like structure depends strongly on the way how the unfolding was induced (temperature, pressure, GndHCl, etc.). Besides structural parameters which are probed by the respective method, the achievable time resolution is of importance. Depending on the dead time of the methods used to initiate folding or unfolding and on the data acquisition time of the measuring technique, the effective time resolution ranges from nanoseconds to minutes (for some H-exchange methods even hours or days)²³. In addition to structural information the characterization of dynamical properties of the native, and even more important of the unfolded state, is of interest. Protein dynamics as measured for the unfolded state can exhibit valuable information about the conformational entropy and about to what extent the polypeptide chain is expanded with respect to the native state. The ever-growing arsenal of biophysical methods available to experimentalists allows nowadays extremely fast transitions to be monitored with very low populated species. In this respect single molecule techniques play a key role since they offer the potential to follow folding events of one molecule at the time. This aspect will be discussed in more detail in the following section.



Fig. 6: Steady state unfolding transitions were measured with intrinsic tryptophan fluorescence (upper left panel) and with far UV CD spectroscopy (upper right panel). From these data the temperature induced transitions were characterized in terms of a fraction of folded protein as a function of the temperature (lower panel). The unfolding process appears rather cooperative with respect to secondary and tertiary structure and represents a typical "all or nothing" behavior (same melting temperatures T_m obtained from both methods, for details see²⁴).

5 The Relevance of Conformational Entropy for Protein Folding and Stability

5.1 Myoglobin: A Model System for Protein Folding Studies

Myoglobin (Mb) serves as a model system to study basic physical properties of protein folding²⁵ and protein dynamics²⁶. The fully folded protein consists of eight α -helices that are interconnected by loops. The protein acts as a scaffold for a heme-group that can bind one O₂ molecule reversible, see Fig. 7. Mb with bound heme-group is termed holo-Mb, while the heme-free protein is called apo-Mb. The ancient Greek words 'holo' and 'apo' mean 'complete' and 'without', which refers to the presence or absence of the heme-group.

The overall structure of native apo-Mb is well defined, but the heme-binding pocket formed by the E and F helices is unstructured and flexible²⁷. During the folding process the long helices G and H are formed first via a hydrophobic collapse²⁵. In one molten globule state (MG1) that occurs during the kinetic folding process the A helix is bound to the G & H helices²⁸. In a second molten globule state (MG2) that occurs after the MG1 state on the folding pathway the helices A,G,H are further stabilized and the helix B is additionally bound²⁸. Therefore, the secondary structure content of the MG2 state is higher than in the MG1 state. Folding times of the apo-Mb molten globules are in the order of some micro- to milliseconds, while the full assembly of the native protein occurs after a few seconds²⁵.



Fig. 7: Kinetic folding process of Mb. The unfolded protein chain first collapses within some μs and forms the molten globule MG1. A second molten globule state MG2 is found during the folding process. Structural rearrangements within MG1 and MG2 take place in the order of several ten to hundred ms that lead to the folded apoMb structure. The integration of the heme-group leads to the biological functional holo-Mb structure. Adapted from Jamin 2005 Protein and Peptide Letters, 12, 229-234

The apo-Mb system has the tremendous advantage that the partially folded molten globules MG1 and MG2 can be stabilized under equilibrium conditions by choosing specific solvent conditions²⁹. Furthermore, a third partially folded molten globule can be stabilized that has less secondary structure content than the MG1 and MG2 states²⁹.

5.2 Conformational Entropy Measured with Neutron Spectroscopy

Energy-resolved neutron scattering is amongst the few methodologies that have provided experimental data in the area of biological molecular dynamics³⁰. In quasielastic incoherent neutron scattering (QENS) experiments it is the incoherent scattering of H atoms that is mainly observed and analysed. In the time-scales examined, H atom motions reflect those of the chemical groups to which they are bound. H-atoms are, therefore, indicators of internal and global dynamics of proteins in solution³¹⁻³³ on the picosecond to nanosecond time-scale. In a first set of experiments fast dynamics of apo- and holo-Mb on the picosecond time-scale have been studied by QENS²⁹. The recorded QENS spectra contain the relevant information on internal protein dynamics. The elastic incoherent structure factor (EISF) $A_0(q)$ contains information on the amplitude of confined protein motions. The EISF was interpreted using the Gaussian approximation according to

$$A_{0}(q) = A \cdot e^{-\langle x^{2} \rangle q^{2}} \cdot (1-p) + p$$
(5)

where $\langle x^2 \rangle$ is the mean square displacement (MSD) of the observed confined motions and *p* accounts for a fraction of slow moving hydrogen atoms, which appear as immobile within the energy resolution of the neutron spectrometer. The measured EISFs and the fits with the model function are shown in Fig. 8.





Fig. 8: (*left*) *EISF* of the investigated conformational states on IN6. (right) Entropy difference per residue of apo-Mb and holo-Mb in different states. The difference of conformational entropy ΔS_{conf} measured by neutron scattering (filled circles), the thermodynamic entropy difference $\Delta S = \Delta S_{conf} + \Delta S_{hydr}$ calculated with the Gibbs-Helmholtz equation with the known thermodynamic parameters (empty circles), and the difference $\Delta S_{hydr} = \Delta S - \Delta S_{conf}$ (filled squares) due to the entropy content of the hydration water and slow-moving protein dynamics.

The difference in conformational entropy in proteins can be determined from QENS measurements as suggested by Receveur et al.³⁴ and Fitter³⁵ according to

$$\Delta S_{conf} = 3R \ln \left(\sqrt{\frac{\left\langle x_u^2 \right\rangle}{\left\langle x_f^2 \right\rangle}} \right) \tag{6}$$

where $\langle x_u^2 \rangle$ and $\langle x_f^2 \rangle$ are the MSDs of the unfolded state and the folded/ partially folded conformation at the same temperature, and R = 8.3144 J/K/mol is the gas constant. In our work, the acid denatured state is considered as a reference for the unfolded state. The determined values of ΔS_{conf} at 16 °C are shown in Fig. 8 as a function of the measured α helical content of the proteins. We find approximately a linear increase of ΔS_{conf} with the helical structure content. The thermodynamic entropy difference $\Delta S = \Delta S_{conf} + \Delta S_{hydr}$ at 16 °C calculated from known thermodynamic parameters (see references in^{29, 32}) is also shown in Fig. 8. The thermodynamic entropy difference of the conformational entropy difference of the protein ΔS_{conf} and of the entropy difference of the hydration water ΔS_{hydr} .



ESG

Fig. 9: Schematic illustration water in the close vicinity of the protein surface for unfolded, partially folded and folded native conformations. Water molecules close to hydrophobic regions of unfolded and partially folded regions that are solvent exposed are ordered and less-mobile, while hydration water in close contact with hydrophilic regions of folded proteins is less-ordered and more mobile.

As expected, the conformational entropy of the protein is reduced with increasing degree of protein folding: the folded structure is less flexible than the partially folded MGs and the unfolded acid denatured state is the most flexible, see Fig 8. Different behavior is observed for water molecules that are in the hydration shell of the protein. Water molecules in the hydration layer of the folded proteins are found to be more disordered and mobile than in the unfolded state, which leads to negative ΔS_{hydr} values, see Fig. 8. Fig. 9 illustrates disordered hydration water around fully folded native proteins, whereas hydration water around unfolded and partially disordered protein regions is less mobile and more ordered. The physical reason is that water molecules form ordered structures around solvent exposed hydrophobic amino-acids³⁶. With decreasing α -helical content hydrophobic residues from the protein core are getting more solvent exposed and thus induce ordering of the hydration shell.

Protein dynamics is of great importance for protein folding, as the polypeptide chain needs to explore the accessible conformational space during the folding transition. See chapter B2 for an overview concerning the field of protein dynamics. The most relevant thermodynamic quantity concerning protein folding that can be determined from the neutron scattering experiments is the difference in conformational entropy ΔS_{conf} between two specific structural states. Only few methods are able to determine ΔS_{conf} experimentally. NMR, for example, allows determining residue resolved information on the picosecond to nanosecond time scale in terms of the order parameter, which is related to an entropy contribution. See chapter A7 for an introduction to NMR. Neutron scattering (chapter A4), on the other hand, is a direct method that probes average protein dynamics on the time scale accessible by the neutron spectrometer. As all dynamic processes contribute to ΔS_{conf} the average information provided by neutron scattering is actually not a drawback.

References

- Buchner, J. & Kiefhaber, T. Protein Folding Handbook Vol. 1-5. (Wiley-VCH Verlag, Weinheim; 2005).
- Dill, K.A., Ozkan, S.B., Shell, M.S. & Weikl, T.R. The protein folding problem. Annu.Rev.Biophys. 37, 289-316 (2008).
- Radford, S.E. Protein folding: progress made and promises ahead. *Trends Biochem Sci.* 25, 611-618 (2000).
- 4. Dill, K.A. Dominant forces in protein folding. Biochemistry 29, 7133-7155 (1990).
- Anfinsen, C.B. Principles that govern the folding of protein chains. *Science* 181, 223-230 (1973).
- 6. Tanford, C. Protein denaturation. Adv. Protein Chem 23, 121-282 (1968).
- Privalov, P.L., Khechinashvili, N.N. & Atanasov, B.P. Thermodynamic analysis of thermal transitions in globular proteins. I. Calorimetric study of chymotrypsinogen, ribonuclease and myoglobin. *Biopolymers* 10, 1865-1890 (1971).
- Tanford, C., Aune, K.C. & Ikai, A. Kinetics of unfolding and refolding of proteins. 3. Results for lysozyme. *J Mol.Biol* 73, 185-197 (1973).
- Tsong, T.Y., Baldwin, R.L. & Elson, E.L. The sequential unfolding of ribonuclease A: detection of a fast initial phase in the kinetics of unfolding. *Proc.Natl.Acad.Sci.U.S.A* 68, 2712-2715 (1971).
- Ikai, A., Fish, W.W. & Tanford, C. Kinetics of unfolding and refolding of proteins. II. Results for cytochrome c. *J Mol.Biol* 73, 165-184 (1973).
- Dolgikh, D.A. et al. Alpha-Lactalbumin: compact state with fluctuating tertiary structure? *FEBS Lett.* 136, 311-315 (1981).
- 12. Strucksberg, K.H., Rosenkranz, T. & Fitter, J. Reversible and irreversible unfolding of multi-domain proteins. *Biochim.Biophys.Acta* **1774**, 1591-1603 (2007).
- 13. Duy, C. & Fitter, J. Thermostability of irreversible unfolding alpha -amylases analyzed by unfolding kinetics. *J Biol.Chem* **280**, 37360-37365 (2005).
- 14. Sanchez-Ruiz, J.M. Protein kinetic stability. Biophys. Chem. 148, 1-15 (2010).
- 15. Onuchic, J.N., LutheySchulten, Z. & Wolynes, P.G. Theory of protein folding: The energy landscape perspective. *Annu Rev Phys Chem* **48**, 545-600 (1997).
- 16. Brockwell, D.J., Smith, D.A. & Radford, S.E. Protein folding mechanisms: new methods and emerging ideas. *Curr.Opin.Struct.Biol* **10**, 16-25 (2000).
- 17. Akmal, A. & Munoz, V. The nature of the free energy barriers to two-state folding. *Proteins* **57**, 142-152 (2004).
- Jackson, S.E. & Fersht, A.R. Folding of chymotrypsin inhibitor 2. 2. Influence of proline isomerization on the folding kinetics and thermodynamic characterization of the transition state of folding. *Biochemistry* **30**, 10436-10443 (1991).
- 19. Sabelko, J., Ervin, J. & Gruebele, M. Observation of strange kinetics in protein folding. *Proc.Natl.Acad.Sci.U.S.A* 96, 6031-6036 (1999).
- 20. Munoz, V. Conformational dynamics and ensembles in protein folding. *Annu.Rev.Biophys.Biomol.Struct.* **36**, 395-412 (2007).
- Ivarsson, Y., Travaglini-Allocatelli, C., Brunori, M. & Gianni, S. Mechanisms of protein folding. *Eur.Biophys.J* 37, 721-728 (2008).
- 22. Parker, M.J., Spencer, J. & Clarke, A.R. An integrated kinetic analysis of intermediates and transition states in protein folding reactions. *J Mol.Biol* **253**, 771-786 (1995).
- Bartlett, A.I. & Radford, S.E. An expanding arsenal of experimental methods yields an explosion of insights into protein folding mechanisms. *Nat.Struct.Mol.Biol* 16, 582-588 (2009).
- Fitter, J. & Haber-Pohlmeier, S. Structural stability and unfolding properties of thermostable bacterial alpha-amylases: a comparative study of homologous enzymes. *Biochemistry* 43, 9589-9599 (2004).
- 25. Dyson, H.J. & Wright, P.E. How Does Your Protein Fold? Elucidating the Apomyoglobin Folding Pathway. *Accounts Chem Res* **50**, 105-111 (2017).
- Frauenfelder, H., McMahon, B.H. & Fenimore, P.W. Myoglobin: The hydrogen atom of biology and a paradigm of complexity. *Proc Natl Acad Sci U S A* 100, 8615-8617 (2003).
- 27. Eliezer, D. & Wright, P.E. Is apomyoglobin a molten globule? Structural characterization by NMR. *J Mol Biol* **263**, 531-538 (1996).
- 28. Jamin, M. & Baldwin, R.L. Two forms of the pH 4 folding intermediate of apomyoglobin. *J Mol Biol* **276**, 491-504 (1998).
- Stadler, A.M., Koza, M.M. & Fitter, J. Determination of Conformational Entropy of Fully and Partially Folded Conformations of Holo- and Apomyoglobin. *J Phys Chem B* 119, 72-82 (2015).
- Fitter, J., Gutberlet, T. & Katsaras, J. Neutron Scattering in Biology Techniques and Applications. (Sringer-Verlag, Berlin; 2006).
- Grimaldo, M. et al. Hierarchical molecular dynamics of bovine serum albumin in concentrated aqueous solution below and above thermal denaturation. *Phys Chem Chem Phys* 17, 4645-4655 (2015).
- Stadler, A.M., Demmel, F., Ollivier, J. & Seydel, T. Picosecond to nanosecond dynamics provide a source of conformational entropy for protein folding. *Phys Chem Chem Phys* 18, 21527-21538 (2016).
- Monkenbusch, M. et al. Fast internal dynamics in alcohol dehydrogenase. J Chem Phys 143 (2015).
- Receveur, V. et al. Picosecond dynamical changes on denaturation of yeast phosphoglycerate kinase revealed by quasielastic neutron scattering. *Proteins* 28, 380-387 (1997).
- 35. Fitter, J. A measure of conformational entropy change during thermal protein unfolding using neutron spectroscopy. *Biophys.J.* **84**, 3924-3930 (2003).
- 36. Ball, P. Water as an active constituent in cell biology. Chem Rev 108, 74-108 (2008).

B2 Protein Dynamics

R. Biehl Institute of Complex Systems, JCNS-1 & ICS-1 Forschungszentrum Jülich GmbH

Contents

1	Intr	oduction2		
2	Local movements			
	2.1	Atomic vibration		
	2.2	Sidechain movements with Atomic Resolution: Time Resolved X-ray Crystallography		
	2.3	Microsecond dynamics of Sidechains and Picosecond Dynamics of the Backbone observed by Nuclear Magnetic Resonance Spectroscopy		
	2.4	Configurational Transition in ClpP on Second Timescale7		
	2.5	Amide ¹⁵ N Backbone Dynamics in Adenylate Kinase		
3	Domain motions			
	3.1	Phosphoglycerate kinase as a classical hinge		
	3.2	Fast antibody fragment motion: flexible linkers act as entropic spring		
	3.3	Intrinsically Disordered Proteins (IDP)		
	3.4	Diffusion in crowded environment with fast process on short scales		
	3.5	Cooperative Rotation of the F1-ATPase Motor		
4	Sum	15 umary		
Refe	erence	s15		

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Proteins are biological macromolecules present in all cells and in body liquids. They work as nanomachines of live to produce material, move objects to the place where they are needed, degrade toxic chemicals, regulate the velocity of processes or protect the cell e.g. from viruses as part of the immune system. Proteins are synthesized as linear polypeptides of 21 different amino acids that are connected by peptide bonds. The primary structure of the protein as the sequence of amino acids is coded in the sequence of DNA. The side chain defines the amino acid properties as acidic or base, hydrophilic or hydrophobic, polar or nonpolar. During the protein synthesis the protein strand folds to a unique 3 dimensional structure (e.g. α -helices, β -sheets, disordered regions, hairpin). The structures is mainly stabilized by hydrogen bonds, hydrophobic amino acids in the core of a domain and hydrophilic amino acids at the surface.

Protein domains are preserved regions of the protein structure with hydrophobic core and hydrophilic shell and can evolve function even if separated from the protein. The domain size reaches from 40 amino acids to several hundred amino acids with an average of approximately 100^{1,2}. Domains of similar structure are found in different species and similar functions show often similar domain structure, which is conserved during evolution and one speaks of domain families conserved during molecular evolution³.

Activity of proteins is often related to an active center where actual function takes place. This can be catalysis of a chemical reaction, binding of a substrate or the change of a chemical potential. An early model about protein activity and specificity was the "lock and key" model, which assumes an exact fit of the protein active site to the substrate due to complementary geometrical shapes but with a rigid conformation as found in crystal structures⁴. To explain also the stabilization of the transition state with bound substrate in a different configuration compared to the unbound state the later "induced fit" model ⁵ allows a reshaping of the binding site to the substrate including local configurational changes of amino acids or large structural changes as for allosteric transitions. Still the protein is viewed as a rigid structure in liganded and unliganded case. Today it is realized that proteins are quite flexible objects, which show configurational changes on all length- and timescales. To reach a buried active site it is often necessary to open a cleft that the substrate can enter and to release the product. Binding of substrates in specific pockets allows to bring them close together in a specific configuration for e.g. phosphate or hydrogen transfer in a protected environment. Necessary conformational changes can be the rate-limiting step in catalysis. In other cases as for kinesin walk on myesin the configurational change is the aim of a chemical process to allow transport of cargo⁶. Therefore protein dynamics on all length scales is a key to understand how conformational changes are related to function and which mechanisms are involved to allow the rich functionality of proteins.

There are two different but linked types of dynamic motions. In terms of an energy landscape view thermal motions are motions that cover the configurational space at thermal energy kT in equilibrium⁷. The accessible configurational space can spread over a single deep minimum or a broader rugged valley where in both the depth defines the occupancy of a configuration in the valley within Boltzmann distribution. Kinetic motions try to find the equilibrium from a higher energy level and are directed towards equilibrium. The higher energy level can be due to an excitation (e.g. photolysis) or binding of a ligand that changes the local energy landscape. Nevertheless thermal motions occur also in the excited state and may help to overcome energy barriers on the kinetic pathway.

The fastest motions in proteins are bond vibrations, side chain rotation at the protein surface or torsion of buried methyl groups with sub-angstrom amplitudes on picosecond timescales. Rearrangements of amino acids to adjust the orientation of functional groups may require local flexibility of neighboring amino acids in a cooperative manner that slows down the process. Movements of secondary structure elements or rearrangement of groups of amino acids on nanosecond timescale allow the adaption of the protein structure to bind specific ligands. Slower motions with larger angstrom amplitudes are relative motions of complete domains as hinge bending movements or swapping of domains depends strongly on the local environment and can be fast as several nanoseconds or slow as up to seconds dependent on the needed rearrangements and the involved interactions. Allosteric transitions, functional conformational changes, folding and unfolding will happen on microsecond timescale and nanometer length scales. In general all these motions are dependent on each other and are coupled. The local atomic fluctuations lubricate the domain motions on larger length scales and domain motions change the shape of the protein.

In the following an overview over some configurational movements is given together with a basic explanation of experimental techniques allowing detecting timescale and amplitude.

2 Local movements

Local movements comprise movements of single atoms or small atom groups. Due to thermal excitation, bonded atoms show vibrating movements as eigenmodes of their specific configuration. Each sidechain of an amino acid can have different configurations with respect to the backbone dependent on the space required for a change e.g. in orientation. Amino acids at the surface have more configurational freedom compared to completely buried amino acids. Configurational changes can be due to thermal movements or due to specific processes as binding of a substrate.

2.1 Atomic vibration

The fastest movements with highest energy are atomic vibrations as they are common for any molecule. Figure 1a shows as an example the geometry of torsion around the C-C bond of a methyl group, a symmetric stretching of the C-H bonds and a symmetric bending of the c-H bonds of the sidechain of alanine. These are only a few possible vibrations of alanine sidechain and each amino acid has different vibrational frequencies dependent on the atomic structure. For an overview see Barth et al⁸. The exact frequency depends not only on the geometry as for free molecules of the same architecture, but also on the direct environment and neighboring amino acids. Hydrogen bonds or polar interactions alter the vibration frequency. Figure 1c shows an example spectrum of $(PPG)_n$ an synthetic polypeptide with a similar structure to natural collagen 9. The frequency range reaches from 100 cm⁻¹ to about 2000 cm^{-1} (12.4 meV – 250 meV) describing motions on a timescale of 0.01ps to 0.3 ps. The difference in the spectra is related to a partial deuteration that changes the frequency and the amplitude of specific vibrations due to the change of the hydrogen mass. In this way the specific exchange allows to separate some of the vibrational modes present in the sample. Typical instruments with the necessary large energy transfer are neutron time of flight instruments.

For proteins the Amide I absorption band between 1600 cm⁻¹ and 1700 cm⁻¹ is of importance because it allows the determination of the relative content of secondary structure in the protein. It can be measured by infrared spectroscopy (e.g. FTIR). In this absorption band the C=O stretch vibration of the peptide backbone is dominant¹⁰. As shown in **Figure 1**b the hydrogen bond between the C=O and the N-H of a different amino acid stabilize the secondary structure elements. A specific hydrogen bond modifies the stretching vibration in a way that is characteristic for the local geometry defined by the secondary structure.

Measuring the different contributions allows extracting the fractional content of secondary structure elements.



Figure 1: Example geometries of torsion, symmetric stretching or bending vibrations of a methyl group as the side group of alanine. b) Secondary structure elements β -sheets and α -helices (in yellow, red) with side chain elements. The secondary structure is stabilized by the hydrogen bonds (thin blue lines). c) Vibrational spectrum of synthetic polypeptide (PPG)n associated in a right-handed supercoiled triple helical arrangement as found in natural collagen measured at the time-focused crystal analyzer spectrometer (TFXA) at ISIS, UK by Middendorf et al.⁹

2.2 Sidechain movements with Atomic Resolution: Time Resolved X-ray Crystallography

Conventional x-ray crystallography measures the diffraction patterns in different orientations of a single crystal with hundreds to thousands of Bragg- reflexes from which the threedimensional electron density is calculated. The crystal structure of atom positions is generated as found e. g in the Brookhaven Protein Databank (PDB). Configurational changes due to activity can only be examined if it is possible to grow a crystal in both states as e.g. liganded with the substrate or an inactive replacement of the substrate and without the substrate. Diffusion trapping can be used to find the binding site of a substrate that is able to diffuse into a substrate free crystal. A larger configurational change due to substrate binding is difficult to access in this way and the timescale of a process cannot be accessed.

Time resolved crystallography tries to measure diffraction patterns with a delay time to observe the time evolution of a process. To access a dynamical process the protein must be active in the crystalline state and the process needs to be triggered inside the crystal with some reasonable amount of concentration. One way is to use a pump-probe method with measurements after a trigger event. Measurements with defined time delay relative to the trigger event e.g. photolysis of a process by a laser pulse allows to follow the process with different delays. On the other hand a complete series with defined delays can be acquired if the repetition rate for a measurement is high enough. To overcome the need of different orientations for sub-second resolution a polychromatic Laue X-ray diffraction technique can be used¹¹. Larger motions as domain movements cannot be observed as the crystal structure limits the configurational freedom to move over larger distances.

Schotte et al. have reported 100 picosecond time resolution x-ray crystallography of a myoglobin mutant after photolysis by an orange laser flash¹². Figure 2 A shows the electron density map of the unphotolyzed myoglobin and 100ps after the flash showing 3 larger configurational changes (arrows) and distributed smaller changes in the protein.



Figure 2: Electron density map of myoglobin before (magenta) and after photolysis (green). Overlapping densities are shown in white. The stick model included shows the unphotolyzed state. Large arrows indicate 3 large changes while small arrows indicate small rearrangements in the whole protein. Sequence of an enlarged view of A. B-E) Times are 100 ps, 316 ps, 1 ns, 3,16 ns, after the laser flash. F and G are 31.6ns and 3.16 μ s after a longer ns laser flash. Circles indicate location of CO molecules. Figure from ¹².

Figure 2 B-E show a sequence of delay times after the laser flash. The CO molecule dissociates and is trapped 2Å apart from the original binding site (positions 2 + 3). Phe29 is displaced but relaxes back to its original position within 316 ps. After several nanoseconds His64 relax to their deoxy position and the CO molecule has migrated to position 4 and 5 where it is trapped for several microseconds.

Here the conformational change due to photolysis is demonstrated with a correlated motion of sidechains rapidly sweeping away the CO molecule from its preferred binding site.

The experiment demonstrates that fast sub nanosecond reorientations are possible even in a restricted Crystal structure.

2.3 Microsecond dynamics of Sidechains and Picosecond Dynamics of the Backbone observed by Nuclear Magnetic Resonance Spectroscopy

Nuclear magnetic resonance measures the resonance frequency characteristic for a spin flip in a magnetic field, which depends on the respective atom type (¹H, ¹³C, ¹⁵N) and e.g. the local electronic environment leading to the chemical shift (dependent on binding partners, bond length and angles) or J-coupling (interaction of different spins). Chemical exchange phenomena as conformational changes, exchange with the solvent or ligand binding modifies this local environment. **Figure 3** shows as an example for a two site chemical exchange the chemical shift dependence on the rate of exchange relative to the measurement time. Microsecond-millisecond changes can be detected due to modification of the chemical environment and lead to a change of the specific chemical shift. Molecular motions on picoseconds to nanoseconds as sidechain vibrations can be detected by measuring spin relaxation.

Figure 3: An overview over timescale in NMR for the case of a two site chemical exchange with a population of 3:1 between A and B. The difference in chemical shifts $\Delta \omega \approx \omega_A - \omega_B = 100 \text{Hz}$ determines a "shutter time" as the NMR time scale by $\sqrt{2}/\pi\Delta\omega$. a) The exchange time $(k_{AB}+k_{BA})^{-1}$ (indicated as times in the figure) is slower allowing to detect 2 separate peaks with intensities related to population of A and B. b) The exchange is on same timescale showing coalescence to a single broadened peak at population averaged position, because within the measurement time single transitions are likely. c) The exchange is much faster and a single sharp peak is observed. Within the measurement time the transition occurs several times. The position depends on population of A and B. For known chemical shifts of A and B the peak position can be used to determine the populations. Figure used from ¹³.

2D-NMR measures the coupling between spins of different atoms in a molecule. Methods like correlation spectroscopy (COSY), ECOSY, HSQC, EXSY allow the detection of correlations between atoms that are connected over several bonds. E.g. Nuclear Overhauser effect spectroscopy (NOESY) allows detecting correlations that are not connected by a bond, but which are not to far separate in space. **Figure 4** shows as an example magnetization exchange spectroscopy (EXSY) for the observation of a slow process on the second timescale. For different configurations two peaks are observed with a cross peak as a direct indication of transition between configurations within the measurement time.

Measuring with a delay time allows the determination of rate constants related to the transition. Faster motions on the piconanosecond timescale are accessible by spin relaxation

spectroscopy. The relaxation rate depends on molecular motions of atoms in the local environment that interact via dipole-dipole interactions or by chemical shift anisotropy and the global rotational motion¹³. Bond vector motions (e.g. amide N-H bond of the backbone in a ¹H-¹⁵N experiment) are separated into contribution from global rotation on the timescale of nanoseconds, which is equal for all residues, and internal dynamics on the picosecond timescale to extract local backbone dynamics. Limitations arise when the dynamics of internal motion becomes slow and reaches the timescale of rotational motion.

Figure 4: Magnetization exchange spectroscopy (EXSY) can observe the slow exchange between conformational states A В. In a standard correlation and experiment observing the backbone amide, the magnetization is transferred from ${}^{1}H$ to ¹⁵N, the ¹⁵N chemical shift is measured and afterwards the magnetization is transferred back to measure the ¹H magnetization by complex RF-pulse sequences. Dependent on the configuration correlation peaks A or B are observed (equivalent to b) T=0). In





the EXSY experiment the ¹H experiment is done after a delay time T>0 (see a) and a cross peak is observed in cases were during the delay time a transition between A and B occurred (b T>0). Varying the delay time allows measurement of the relaxation between the two configurations A and B. c) The time dependence of a relaxation with decreasing intensities for A and B. The cross peak intensity increases at short times. All intensities decrease at long times because of spin relaxation. From the relaxation the forward and reverse rate constants can be extracted. The technique is applicable for rate constants in the range from 0.5 s⁻¹ to 50 s⁻¹. Figure used from ¹³.



Figure 5: Spin relaxation: a) At equilibrium the spins are aligned to the magnetic field. After a specific perturbation (RF pulse) the spin is aligned antiparallel for longitudinal magnetization R1 or perpendicular for transverse magnetization R2. After a delay time allowing relaxation an inversion pulse is applied and after a further delay the non-relaxed magnetization can be measured. b+c) The magnetization is interpreted in a model free approach by $S^2+(1-S^2)exp(-t/\tau)^{-39}$. c) R1 relaxation depends on global tumbling e.g. due to rotational diffusion and local fluctuations. R2 depends only on the local relaxations. Local relaxations can be extracted by combined measurements. Figure used from ¹³.

2.4 Configurational Transition in ClpP on Second Timescale

A slow transition in the oligometric protease ClpP, consisting of 14 subunits arranged in a heptameric ring with a total mass of 300kDa, was observed by magnetization exchange spectroscopy (EXSY)[15]. ClpP is a cylindrical self-subdividing protease with a chamber containing the proteolytic active site were through axial pores the already unfolded substrate is inserted and proteolysis takes place (see Figure 6a). The assignment of the δ 1 methyl peaks for I149 and I151 shown in Figure 6b was done by site directed mutagenesis and presents two peaks for each of the residues, a broad and a narrow peak indicated as F and S for fast and slow relaxing. By EXSY it was established that the two peaks correspond to two different configurations and that a transition occurs with a rate constant of 60 s⁻¹ at 0.5 °C. Both residues show the same rate constants suggesting that the same process is observed. I149 and 1151 are present in all monomers building a ring at the connecting interface between the heptameric rings. The observed configurational change could be related to opening of pores, which allow the release of the product after proteolysis. This was tested by mutagenesis in introducing at position 153 a cysteine that can build a disulfide bridge between monomers if oxidized. A series of kinetic measurements showed fast product release from the reduced form, but product release was not observed in the oxidized form. The disulfide bridges quench the release process. The relaxation process observed by NMR is likely related to the release of the proteolysis products through pores in the equatorial plane.



Figure 6: *a)* ClpP Protease with two monomers shown in yellow and blue. Isoleucines are indicated as red and green circles. Blue arrows indicate substrate-entering pores. b) TROSY correlation spectrum with the two δ 1 methyl peaks (indicated F and S for fast and slow in the magnification on left of the small rectangle right) of 1149 and 1151 residue. The occurrence of two peaks is a result of the slow exchange between two conformations. c) Relaxation of the peak intensities with the relaxation of the cross peak. All relaxation rates are approximately equal. Figures from ^{14,15}

2.5 Amide ¹⁵N Backbone Dynamics in Adenylate Kinase

Atomic fluctuations on picosecond-nanosecond timescale can facilitate larger motions on slower timescale as was shown for adenylate kinase by spin relaxation spectroscopy of the ¹⁵N bonds by Henzler-Wildman et al. ¹⁶. A mesophilic and thermophilic homologue were examined at different temperatures. The measured spin relaxation spectra were analyzed by models derived from model free analysis, determining the order parameter S^2 and a relaxation time from the relaxation curve of each residue. In the basic model relaxation is described by $C(t)=S^2+(1-S^2)\exp(-t/\tau)$. The S² describes the change from completely rigid structures (S²=1) over flexible structures (0.85 for secondary structure or 0.5 for unstructured regions) to completely uncorrelated motions ($S^2=0$) and is most sensitive to local packing¹⁷. Different models and additional measurements were used to remove contributions from lid movements on us scale, anisotropic diffusion or coupled motions. Rotational correlation time was found as 14 ns⁻¹ and 18 ns⁻¹ at 20°C for thermophilic adenylate kinase in open and closed configuration in accordance to calculation of HYDRONMR¹⁸ based on the PDB structure. **Figure** 7a shows the mesophilic adenvlate kinase colored according to the value of S^2 with values between 0.92 as rigid and 0.65 for flexible regions. It was found that hotspots of flexibility are found where backbone conformation must change for lid motion as indicated by arrows in the figure. E. g. hinge 8 corresponds to a kink in the long α -helix if the lid is closed during substrate binding. Figure 7b shows the S^2 values for all residues with marked hinges. With increased temperature the overall flexibility increases. At 20°C the thermophilic adenylate kinase shows smaller order parameter values compared to the mesophilic homologue. Similar order parameters for both homologues are found if both are measure in the same distance from their living temperature (30° below 54° C and 109° C), supporting the hypothesis that fluctuations are related to stability.



Figure 7: a) Flexibility of the mesophilic adenylate kinase at 20°C. Colors indicate the value of the order parameter S^2 with grey for residues for which S^2 cannot be measured (due to fast hydrogen exchange or spectral overlap). b) Temperature dependence of S^2 for the thermophilic adenylate kinase at 20°C (blue), 50°C (black) and 80°C (red). Numbers with arrows indicate in both plots identified hinges. Figures from ¹⁶

3 Domain motions

Domain motions are correlated motions between different domains or inside of domains like bending and torsion of the domain. The mechanisms are described as shear motions along an interface or hinge motions where such an interface is missing, but unclassifiable motions are allowed. The connection between domains can be a broad soft hinge as in the case of phosphoglycerate kinase, a single α -helix as in the case of lactoferrin or a disordered amino acid sequence of the protein chain as in the case of immunoglobulin G1 or mercury ion reductase.

Configurational changes changing the shape of a protein by rearranging complete domains can be observed by methods, which allow the detection of correlated motions over larger distances. Förster resonance energy transfer (FRET) measures the distance distribution of chromophores attached to cysteines at specific sites on the protein surface or introduced by site directed mutagenesis¹⁹. The excitation of the donor is followed by a transfer of the energy to the acceptor - overlap of absorption and emission spectra is essential - and detection of the later emission. The efficiency of energy transfer is dependent on the distance and allows measuring a distance distribution between the chromophores respectively the points where the chromophores are bound to the protein. Single molecule fluorescence polarization anisotropy (smFPA) measures dynamic changes of the orientation via time resolved polarization measurements and yields information about size, shape and rotational dynamics²⁰. Electron spin resonance can be used to measure the interaction between two unpaired spins similar to NMR, but because of the stronger interaction the distances can be larger²¹. Because of the absence of unpaired electrons in proteins, spin labels have to be inserted by site directed mutagenesis. Another technique is small angle scattering by neutrons or x-rays (SANS, SAXS). Both techniques allow examination of the low-resolution structure and the ability to compare structural changes in solution due to changes of pH, salt concentrations, temperature or substrate addition. Time resolved SAXS can reach sub-millisecond resolution and can be combined with stopped flow experiments or trigger events as in time resolved x-ray crystallography²². Neutron spin echo spectroscopy (NSE) is a technique that is able to access

the timescales from 0.1 up to several hundred nanoseconds and simultaneously covers the length scale relevant for protein domain movements as in SAXS or SANS of several nanometers distance between domains²³. NSE measures the temporal correlation of configurational changes under utilization of the neutron spin to detect tiny velocity changes during the scattering process. The measured intermediate scattering function can be interpreted as a time correlation between small angle scattering patterns with nanosecond resolution. Main contributions to the correlation come from translational and rotational diffusion due to the spatial correlation of different proteins diffusing in the solution and internal domain dynamics on nanosecond timescale. In the following we present exemplary results to demonstrate a small variety of possible motional patterns.

3.1 Phosphoglycerate kinase as a classical hinge



Figure 8: Form factor measured by SANS from PGK and PGKsub (Kratky plot: Q versus $Q^2I(Q)$). Lines show the calculated form factor from the crystal structure (blue) and from structures deformed along the softest normal modes to fit the experimental data. The inset shows the protein with the hinge in yellow, the main domains in blue and red and the substrates as spheres. Figure from ²⁷.

Phosphoglycerate kinase (PGK) is an enzyme that is involved in glycolysis. It relocates a phosphate group from 1,3biphosphoglycerate, an intermediate product in glycolysis, to ADP to synthesize ATP^{24} . PGK is composed of two separated domains connected by a hinge region as shown in Figure 8 ²⁵. 1,3-biphosphoglycerate and ADP are bound at opposite positions in the cleft at the two domains. The active site is located at the hinge between the C- terminal and N-terminal domains. Evidence for a hinge bending motion induced by substrate binding was found by Bernstein at al. by comparing crystallographic structures of different species with and without bound substrates bringing the substrates closer together²⁶. The crystal structure without substrate has an open cleft configuration with key residues Arg-39 and Gly-376 in the active center separated by about 1.18

nm. The proposed mechanism of induced fit due to substrate binding closes the cleft by a 32° hinge closure to the active configuration with bound substrate. This cleft closing motion mainly brings key residues Arg-39 and Gly-376 of the active center together with the substrates to a distance of 0.35 nm as found in the closed cleft crystal structure.

Figure 8 shows SANS measurements of PGK in solution with and without substrate demonstrating that the solution structures are different from the substrate bound crystal structure²⁷. Modeling the structure by deformations along softest elastic normal modes (torsion and 2 perpendicular bends of the hinge) allowed modeling of the deformation due to substrate binding in solution. The distance between the active residues was reduced from 1.14 nm without substrate to 0.82 nm with bound substrate, but still to far to allow activity within a static structure.

NSE measurements show the relaxation of the intermediate scattering function dependent on diffusion and internal dynamics as shown in **Figure 9** left. At low scattering vectors Q (observing large length scales) the protein looks point like and a single exponential relaxation is observed. t larger Q (observing length scale of protein size) additional contribution arise from internal dynamics. At long times diffusion is observed, that can be described within a

single model for all Q values by rotational and translational diffusion constants. The additional component at short times is described by a Q dependent amplitude and a relaxation time τ , which both depend on the construction of the specific hinge and influence function.



Figure 9: Left: Semi logarithmic plot of I(Q,t)/I(Q,t=0) for selected Q values (PGK, black and red; PGKsub, green; data are shifted for clarity and are equal 1 for t=0). Red and green dashed lines represent the initial slope extrapolated to long times. The blue and black dashed lines represent the long-time limit extrapolated to t=0. The long-time limit corresponds to rigid-body diffusion including inter-particle effects. The difference between extrapolated long time diffusion at t=0 to the initial slope amplitude is the internal dynamics contribution $A(Q)exp(-t/\tau)$. Right: Q dependence of A(Q) compared to model calculations. Figure from ²⁷.

The observed relaxation times are 60 ns for the substrate free and 45 ns for the substrate bound PGK. **Figure 9** shows at the right the amplitude A(Q) with model calculations based on the softest normal mode deformations. From the deformation of the hinge the distance of the active residues can be calculated. It was found that the thermal driven deformation of the protein at the hinge is strong enough to reach configurations that allow activity as the active residue come close enough together.

In summary, the NSE investigation has demonstrated that the approach to a functional configuration of PGK needs to be attributed to the dynamic fluctuations of the main domains instead of an earlier proposed induced fit by substrate binding. Thus, in the case of PGK, hinge dynamics enables function.

3.2 Fast antibody fragment motion: flexible linkers act as entropic spring

The fragment motion of IgG was examined by Stingaciu et al by NSE^{28} . In the experiment, aside from strong diffusion contributions, a clear signature of the internal dynamics on a timescale of 6 - 7 ns was observed. The combination of translational and rotational diffusion of the rigid protein describes the observed long-time dynamics on an absolute scale, despite the fact that we have a mixed population of monomers and dimers, as found in serum.

The fragment motion shows itself in terms of a strong decay of S(Q,t) at short times, well separated from the overall diffusional relaxation. The data were analysed in terms of an Ornstein-Uhlenbeck like relaxation within a harmonic potential. Such an approach neglects the details of the complex linker region interaction on the residue level. Nevertheless, the dynamics is described in an excellent way. Even for a short flexible linker the ensemble average seems to be sufficient to produce the characteristics of a harmonic spring. The obtained spring constant ($\approx 10 \ pN/nm$) appears to be realistic compared to an entropic spring of similar length. The resulting forces stay below the limits that would be necessary to



Intermediate scattering functions I(Q,t)/I(Q,0) of IgG. Dashed lines show translational and rotational diffusion and are close to a single exponential for times t > 15 ns. Solid lines include additional internal dynamics of each domain in its own harmonic potential with a 7ns relaxation time.

unfold the secondary structure of the attached fragments. The observed effective friction appears to be close to the friction of a free unbound fragment in the solvent. None of the directions seems to be suppressed.

The pre-existing equilibrium hypothesis²⁹, with multiple local minima in the configurational energy landscape, will have long transition times between the minima but is compatible with fast 6-7 ns motions in the local minima. Consequently, NSE observes the fast dynamics pre-existing in

equilibrium states. The conformational flexibility in a pre-existing equilibrium configuration might be needed to adapt faster to a specific antigen when the antigen approaches the binding site of a fragment.

3.3 Intrinsically Disordered Proteins (IDP)

About 40 % of all proteins in the human body are intrinsically disordered. These IDP's do not exhibit any well-defined folded structure that could be crystallized. Their structural and dynamic properties reach from very soft structures over folded elements connected by extended and flexible loops to fully disordered polypeptide chains. The biological role of the IDP is founded in their high conformational adaptivity, enabling them to respond rapidly to environmental changes or other macromolecules allowing them to fold into different states or



Figure 10: SAXS data from MBP at 4.5 mg/ml. a) The red solid lines is a fit with the Debye equation for a Gaussian chain. (b) Kratky plot. The line is a result of the scattering from the most probable conformational ensembles shown in c. c) Representative coarse-grained conformations of MBP as determined by inverse Monte Carlo. The structures are rotated by 90 °C in the lower part of the figure. The color code relates to six different realizations of the ensemble, which represent the data best. Figure from ³².

even a rigid 3D structure. For these properties dynamics is essential. Specifically the sampling of the energy landscape and the exploration of the large conformational space are driven by conformational motions of the unfolded peptide chain. On the other hand IDPs show the same

dynamics as proteins during the early stage of folding. Since the intrinsic disorder prevents crystallographic structure determination, only low resolution SANS and SAXS information about the average structure in the disordered state exists.

Myelin basic protein (MBP) is a major component of the Myelin sheaths in the central nervous system³⁰. In the human body MBP is of significant importance as there are many neurological disorders, as e.g. multiple sclerosis, that are related to MBP mal function. Lipid free MBP is not completely unfolded but retains some elements of the alpha helix and beta sheet (about 60 % of the protein is unfolded)³¹.

Figure 10 a+b display X-ray form factors of MBP that display strong similarities to polymer form factors as the high Q data show a power law close to Q^{-2} that is characteristic for Gaussian chain polymers in theta solvents. The small increase visible in the Kratky plot at high Q indicates a length scale where the random oriented character vanishes and the linear character of short chain segments becomes visible. A Monte Carlo simulation was used to generate a coarse-grained ensemble representing the structural characteristics of and a



Figure 11: Left) NSE Data for MBP: All spectra start at unity but are shifted consecutively by a factor of 0.8 for clarity shown up to 50 ns. Solid lines are fits to the NSE data with the structural model. The dashed lines are exponential fits for t>20 ns to extrapolate the long time rigid body dynamics. A clear separation between the internal and the global dynamics is obvious. Middle) Displacement pattern of the normal modes 7 (upper part, bending) and 8 (lower part, stretching) from the structural model according to normal mode analysis. The lengths of the vectors are increased for better visibility. Right) Amplitude of the internal protein dynamics as obtained from the fit. The solid and dashed lines are the calculated mode amplitude according to equation 4. Figures from ³².

ensemble was selected that represent the SAXS data³³. The resulting characteristic molecular shapes are displayed in **Figure 10**c. The model conformations indicate an elongated structure with a relatively compact core and flexible ends on both sites.

Figure 11 at left displays NSE spectra from 54 mg/ml solutions (times up to 140 ns could be accessed). Inspecting this figure, the two-component structure of the NSE spectra at Q-values above 0.9 nm^{-1} is visible. Thus, we deal with long time rigid body motion augmented by internal dynamics with relaxation times below 10 ns.

The structural models based on SAXS analysis were used to describe the long time translational and rotational diffusion combined with a Q-dependent motional amplitude A(Q) and an internal mode relaxation time. The characteristic internal relaxation time τ_{int} =8.4±2.0 ns is found for the whole structural ensemble and the corresponding amplitude A(Q) is displayed in **Figure 11** at the right. Normal mode analysis was used to describe the deformation of the structural models. **Figure 11** middle shows the first two normal modes as a bending of the structure, which already give a satisfactory description of the observed amplitudes. Comparing the normal modes with the structural models in **Figure 10**c it can be concluded that the normal modes describe approximately the motion from one structural

model to the next. Here it is shown that even for very flexible structure as the amino acid chain still low frequency collective stretching and bending motions of the outer part of the structure describe the essential features of the large-scale dynamics.

3.4 Diffusion in crowded environment with fast process on short scales

The natural environment of proteins is a crowded environment as in cells, extracellular fluids or during processing. Semidilute polymer solutions have been a source of rich structural and dynamical properties and mimic a crowded environment comparable to protein's natural environment in the cell or during processing. Dispersing model globular proteins α -Lactalbumin (La) and Hemoglobin (Hb), in aqueous solution of poly-(ethylene oxide) (PEO) Gupta et al.³⁴ mimic a crowded environment and use neutron spin echo (NSE) to observe the corresponding protein dynamics in semidilute polymer solution (see **Figure 12**). NSE can access the fast diffusion process (D_{fast}) prior to the slow diffusion process on long times and length scales (D_{γ}) that can be consistently described based on particle diffusion in a periodic potential. The fast diffusion process describes the diffusion inside of a mesh or trap until it jumps out of the trap and takes part in the slower diffusion process (while also being trapped again). No coupling of the particle dynamics to the polymer dynamics was observed. Instead of the solvent viscosity the effective viscosity due to the presence of unentangled polymer segments dominates the fast diffusion while diffusing inside of the traps.

3.5 Cooperative Rotation of the F1-ATPase Motor

Adachi et al. demonstrated the stepping rotation of the F-ATPase motor through angleresolved smFPA³⁵. F-ATPase is a membrane protein complex that synthesizes ATP. A proton gradient between membrane sides is used to drive the reaction as the protons cross the membrane in a transport reaction in the membrane bound F_0 -ATPase, which acts as a rotor. F_0 -ATPase reaches into the stator F_1 by a shaft (dotted line, see **Figure 13**). A fluorescent probe was attached on top of F_0 reaching through F_1 to monitor the polarization of emitted or absorbed light during activity. **Figure 13** shows on the right the intensity course during



Figure 12: Normalized intermediate scattering functions I(Q,t)/I(Q,0) as a function of Fourier time for (a) $\phi = 5\%$, (b) 10% and (c) 15% d-PEO in Hb- D_2O solution (vertically shifted). Blue lines represent the extrapolated slow diffusion process to short times. Red lines incorporate the fast in trap diffusion process at short times.

with activity the calculated orientation corresponding and revolution angle. The time evolution of a continuous motion is shown for comparison as a red line. The stepwise character of the motion is demonstrated as a true property of the ATPase and further analysis results in 120° steps with a dwell time dependent on ATP concentration in the solution as e.g. 5.5 s at 20 nM ATP.



Figure 13: Left) Schematic sketch oft the F-ATPase complex in the membrane with F_1 in yellow-red and F_0 in light-blue, green in the membrane. The rotor (blue) goes through F1 up to the top. (http://uni-marburg.de/uT9Bj). Middle) F_1 -ATPase motor with a single fluorescent probe attached to the rotor reaching through F_1 . Right) Time trajectories of the fluorescence intensity (black, top) and calculated fluorophore angle between 0° and 180° (green, bottom). The accumulated rotation angle (blue, bottom) was obtained by assuming that all steps were counterclockwise. Reprinted from ^{20,35}.

4 Summary

In the past, the function of biological assemblies was discussed in terms of structure, which was in most cases derived by X-ray crystallography. In recent years the biological community became more and more aware of the importance of motions and dynamics in proteins that can play an important role in understanding function.

The importance of protein dynamics may be highlighted in the frame of drug design. While in the past in general the development of drugs was done using static crystallographic structures implying the lock-and-key model, during the last 10 years the state-of-the-art involves ensemble docking. Ensemble docking considers conformational changes and searches for conformations where a drug can bind to the active site of the protein, meaning that metastable protein states identified in molecular dynamics (MD) simulation are individually targeted^{36–38}.

References

- (1) Islam, S. A.; Luo, J.; Sternberg, M. J. E. "Protein Eng. Des. Sel. 1995, 8 (6), 513.
- (2) Wheelan, S. J.; Marchler-Bauer, A.; Bryant, S. H. Bioinformatics 2000, 16 (7), 613.
- (3) Jacob, F. Science (80-.). 1977, 196 (4295), 1161.
- (4) Fischer, E. Berichte der Dtsch. Chem. Gesellschaft 1894, 27 (3), 2985.
- (5) Koshland, D. E. Proc. Natl. Acad. Sci. U. S. A. 1958, 44 (2), 98.
- (6) Vale, R. D. Cell 2003, 112 (4), 467.
- (7) Onuchic, J. N.; Luthey-Schulten, Z.; Wolynes, P. G. Annu. Rev. Phys. Chem. 1997, 48, 545.
- (8) Barth, A. Prog. Biophys. Mol. Biol. 2000, 74 (3-5), 141.
- (9) Middendorf, H. D.; Hayward, R. L.; Parker, S. F.; Bradshaw, J.; Miller, a. *Biophys. J.* 1995, 69 (2), 660.
- (10) Susi, H.; Byler, D. M. Enzyme Structure Part K; Methods in Enzymology; Elsevier,

1986°	Vol	130
1,000,	v 01.	1.50.

- (11) Ren, Z.; Bourgeois, D.; Helliwell, J. R.; Moffat, K.; Šrajer, V.; Stoddard, B. L. J. Synchrotron Radiat. 1999, 6 (4), 891.
- Schotte, F.; Lim, M.; Jackson, T. A.; Smirnov, A. V; Soman, J.; Olson, J. S.; Phillips, G. N.; Wulff, M.; Anfinrud, P. A. *Science (80-.).* 2003, 300 (5627), 1944.
- (13) Mittermaier, A. K.; Kay, L. E. Trends Biochem. Sci. 2009, 34 (12), 601.
- (14) Mittermaier, A.; Kay, L. E. Science (80-.). 2006, 312 (5771), 224.
- (15) Sprangers, R.; Gribun, A.; Hwang, P. M.; Houry, W. A.; Kay, L. E. Proc. Natl. Acad. Sci. U. S. A. 2005, 102 (46), 16678.
- (16) Henzler-Wildman, K. a; Lei, M.; Thai, V.; Kerns, S. J.; Karplus, M.; Kern, D. *Nature* 2007, 450 (7171), 913.
- (17) Zhang, F.; Brüschweiler, R. J. Am. Chem. Soc. 2002, 124 (43), 12654.
- (18) García de la Torre, J.; Huertas, M. L.; Carrasco, B. J. Magn. Reson. 2000, 147 (1), 138.
- (19) Heyduk, T. Curr. Opin. Biotechnol. 2002, 13 (4), 292.
- (20) Weiss, S. Nat. Struct. Biol. 2000, 7 (9), 724.
- (21) Sahu, I. D.; McCarrick, R. M.; Lorigan, G. A. Biochemistry 2013, 52 (35), 5967.
- (22) Graceffa, R.; Nobrega, R. P.; Barrea, R. A.; Kathuria, S. V; Chakravarthy, S.; Bilsel, O.; Irving, T. C. J. Synchrotron Radiat. 2013, 20 (Pt 6), 820.
- (23) Biehl, R.; Richter, D. J. Phys. Condens. Matter 2014, 26 (50), 503103.
- (24) Scopes, R. K. Enzym. 1973, 8, 335.
- (25) Bryant, T. N.; Watson, H. C.; Wendell, P. L. Nature 1974, 247 (5435), 14.
- (26) Bernstein, B. E.; Michels, P. A.; Hol, W. G.; Micheis, P. A. M.; Hol, W. G. Lett. to Nat. 1997, 385 (6613), 275.
- (27) Inoue, R.; Biehl, R.; Rosenkranz, T.; Fitter, J.; Monkenbusch, M.; Radulescu, A.; Farago, B.; Richter, D. *Biophys. J.* 2010, 99 (7), 2309.
- (28) Stingaciu, L. R.; Ivanova, O.; Ohl, M.; Biehl, R.; Richter, D. Sci. Rep. 2016, 6, 22148.
- (29) Pauling, L. J. Am. Chem. Soc. 1940, 62 (10), 2643.
- (30) Harauz, G.; Ishiyama, N.; Hill, C. M. .; Bates, I. R.; Libich, D. S.; Farès, C. Micron 2004, 35 (7), 503.
- (31) Polverini, E.; Fasano, A.; Zito, F.; Riccio, P.; Cavatorta, P. *Eur. Biophys. J.* **1999**, *28*, 351.
- (32) Stadler, A. M.; Stingaciu, L.; Radulescu, A.; Holderer, O.; Monkenbusch, M.; Biehl, R.; Richter, D. J. Am. Chem. Soc. 2014, 136 (19), 6987.
- (33) Bernado, P.; Mylonas, E.; Petoukhov, M. V; Blackledge, M.; Svergun, D. I. J Am Chem Soc 2007, 129, 5656.
- (34) Gupta, S.; Biehl, R.; Sill, C.; Allgaier, J.; Sharp, M.; Ohl, M.; Richter, D. Macromolecules 2016, 49 (5), 1941.
- (35) Adachi, K.; Yasuda, R.; Noji, H.; Itoh, H.; Harada, Y.; Yoshida, M.; Kinosita, K. Proc. Natl. Acad. Sci. U. S. A. 2000, 97 (13), 7243.
- (36) Korb, O.; Olsson, T. S. G.; Bowden, S. J.; Hall, R. J.; Verdonk, M. L.; Liebeschuetz, J.

W.; Cole, J. C. J. Chem. Inf. Model. 2012, 52 (5), 1262.

- (37) Andrusier, N.; Mashiach, E.; Nussinov, R.; Wolfson, H. J. Proteins 2008, 73 (2), 271.
- (38) Huang, S.-Y.; Zou, X. Proteins 2007, 66 (2), 399.
- (39) Lipari, G.; Szabo, A. J. Am. Chem. Soc. 1982, 104 (17), 4546.

B 3 Crowded Protein Solutions: Dynamics, Clustering and Phase Behavior

G. Nägele Soft Matter Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

Introduction Globular Proteins		2	
		3	
2.1	Examples: BSA, lysozyme and apoferritin	3	
2.2	Direct and hydrodynamic interactions	4	
Crowding and Salt Effects in BSA Solutions			
3.1	Static scattering and equilibrium structure	7	
3.2	Short-time collective diffusion and solution viscosity	9	
3.3	Generalized Stokes-Einstein relations	12	
Short-range attractive and long-range repulsive Proteins			
4.1	Phase behavior and structural properties	14	
4.2	Dynamics in cluster-fluid and dispersed-fluid phases	16	
Patchy Colloid Model of Lysozyme			
5.1	Kern-Frenkel type model	18	
5.2	Phase diagram calculation	18	
Con	clusions and Outlook	21	
	Intro Glot 2.1 2.2 Crov 3.1 3.2 3.3 Shor 4.1 4.2 Patc 5.1 5.2 Cond	Introduction Globular Proteins 2.1 Examples: BSA, lysozyme and apoferritin 2.2 Direct and hydrodynamic interactions 2.2 Direct and hydrodynamic interactions 3.1 Static scattering and equilibrium structure 3.2 Short-time collective diffusion and solution viscosity 3.3 Generalized Stokes-Einstein relations Short-range attractive and long-range repulsive Proteins 4.1 Phase behavior and structural properties 4.2 Dynamics in cluster-fluid and dispersed-fluid phases 5.1 Kern-Frenkel type model 5.2 Phase diagram calculation Conclusions and Outlook	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Physics-based studies of concentrated aqueous solutions of bio-particles have become an important part of soft matter science [1–8]. In particular, structural and dynamical properties of concentrated (crowded) protein solutions have attracted much attention due to their biological and pharmaceutical importance, e.g., for the understanding of cellular functions and the improvement of drug delivery. Proteins constitute identical solute units, surpassing any synthetic colloidal suspension in terms of monodispersity. In this respect, they are well suited to the application of simplifying theoretical models used with good success for suspensions of synthetic colloidal particles. However, to date it is still very challenging in colloidal science to accurately describe the structure, phase behavior and dynamics of concentrated protein solutions, as many proteins have irregular shapes, complex internal conformations, and heterogeneous surface properties in form, e.g., of discrete distributions of surface charges and hydrophobic patches. The irregular protein surfaces imply an orientation-dependent protein-protein interaction energy with repulsive and attractive contributions. This complicates considerably the description and calculation of hydrodynamically influenced transport properties including diffusion coefficients, mean-squared displacement (MSD), and the solution viscosity.

Great efforts have been invested into the study of aqueous solutions of *globular* proteins which in theory and simulation can be modeled as spherical objects, with an (in general) isotropic effective pair potential describing the shorter-ranged hydrophobic and van der Waals attractions, and the longer-ranged electrostatic repulsion. By studying globular proteins, useful insights are gained into more complex protein systems such as monoclonal antibodies.

In this lecture, we describe how equilibrium structural and dynamic properties, and the phase behavior of globular protein solutions can be predicted on basis of coarse-grained colloidal models, using state-of-the-art theoretical methods and dynamic simulations. The low-molecular-weight solvent (commonly water) in these models is described as a hydrodynamic continuum, and the proteins as interacting Brownian spheres executing highly irregular trajectories owing to the thermal bombardment by the solvent molecules. Additionally to the directly transmitted interactions describable by a pair potential, the proteins are influenced *indirectly* by solvent flow perturbations created by their collective motions. These so-called hydrodynamic interactions (HIs) strongly affect the protein dynamics in crowded solutions [2, 4, 5, 9, 10]. The diffusion and structural properties of protein solutions can be measured, e.g., using synchrotron (i.e., x-ray), neutron and laser light scattering techniques, allowing for the comparison with theoretical predictions.

We will present three exemplary studies of protein solutions where theoretical, simulation and experimental works have been combined. The first study deals with crowding effects in solutions of bovine serum albumin (BSA) proteins under low and physiological salt concentration conditions [6, 11]. The second study addresses proteins having competing short-range attractive (SA) and long-range repulsive interactions (LR), such as in low-salinity lysozyme solutions where the proteins can form reversible clusters [12–14]. These so-called SALR systems have been intensely explored in the past years [8, 15–17]. They have a rich phase behavior, including equilibrium and non-equilibrium cluster phase states [7]. Finally, we explain how the equilibrium phase diagram of lysozyme solutions with higher salt content is calculated using a refined colloid model of spherical proteins with patchy interactions [3, 18]. While some necessary physical formulas are provided and explained in the lecture, our focus is on the physical explanation of various diffusion, clustering and phase transition phenomena.

For conciseness, crowding effects in single-species protein solutions are discussed. Inside a



Fig. 1: Examples of globular proteins that have been described theoretically using colloid models. Left: (Bovine) Serum Albumin (BSA). Middle: (Hen-egg) Lysozyme. Right: Apoferritin. Figures not to scale. Taken from RCSB Protein Data Bank [19].

cell, many macromolecular species of different sizes, shapes and electric charges are present. The physics of intracellular self- and tracer-diffusion is discussed in Chap. D1, and explored using Brownian Dynamics simulations. General features of self- and collective diffusion are discussed also in Chap. A3, and static crowding effects in Chap. D2

2 Globular Proteins

In this section, we first present three examples of globular protein solutions that have been described theoretically as effective colloidal spheres. This is followed by a discussion of suitable (direct) pair potentials describing the pairwise protein interactions. Additionally, solvent flow mediated hydrodynamic interactions are discussed for the experimentally relevant (colloidal) time scale where the coupled motions of proteins and intervening solvent are quasi inertia-free.

2.1 Examples: BSA, lysozyme and apoferritin

Consider Fig. 1 showing the typical (van der Waals radius) atomistic structure of a selection of experimentally and theoretically well studied globular proteins, namely BSA, lysozyme and apoferritin.

As seen, BSA is a heart-shaped globular protein with a linear extension of about 7 nm. It is readily soluble in water, and stable over a wide range of salt and protein concentrations. Its stability and low-cost derivation from bovine blood make it well-suited as an experimental model system of globular proteins. At higher salt content, however, it tends to form dimers and oligomers. The aqueous BSA solutions considered in this lecture have monovalent added salt (NaCl), and pH values in between 5.5 - 7. Under these conditions, a BSA protein is folded in its native state and has a negative charge in the range of 8 - 20 elementary charges.

Lysozyme occurs naturally in hen white egg. It is an enzyme which damages bacteria by hydrolysing cell wall attached polysacharides. The shape of lysozyme is ellipsoidal, with a linear extension of about 6 nm. Its positive charge number known from titration experiments is 8 in units of the elementary (proton) charge e, for a pH value adjusted to 7.8 using a sodium hydroxyl solution [3, 20].

While not discussed in this lecture, the figure includes apoferritin as another well-studied globular protein found in the intestinal mucosa and the liver. It has a diameter of about 13 nm, and it is composed of 24 peptide subunits forming a hollow shell in which iron atoms can be packed. A typical ferritin protein contains about 2000 iron atoms at its core, but potentially may hold up to 4500. Pores are present on the surface of the complex allowing iron atoms to enter and be released from the core. Apoferritine solutions are relatively stable, and depending on the solution ionic strength and concentration the proteins can be significantly charged [1,2].

2.2 Direct and hydrodynamic interactions

In calculating structural and dynamic properties of protein solutions, it is advisable to use a model of reduced complexity, by describing the globular proteins as charged colloidal spheres, with their interaction parameters such as the pH-dependent protein charge determined from a consistent theoretical fit of the average intensity of scattered neutrons or photons in a scattering experiment.

In the spherical colloid model of proteins, the interaction potential, $u(r; \Omega_1, \Omega_2)$, of two spheres 1 and 2 at center-to-center distance r and respective orientations quantified by the solid angles (unit vectors) Ω_1 and Ω_2 , is typically of the form

$$u(r; \mathbf{\Omega}_1, \mathbf{\Omega}_2) = u_{\rm hc}(r) + u_{\rm rep}(r) + u_{\rm att}(r; \mathbf{\Omega}_1, \mathbf{\Omega}_1).$$
(1)

It consists of a hard-core contribution

$$u_{\rm hc}(r) = \begin{cases} \infty, & r \le \sigma \\ 0, & r > \sigma \end{cases}$$
(2)

with an *effective* hard-sphere diameter σ , a longer-ranged soft repulsive contribution, $u_{rep}(r)$, due to the screened electrostatic repulsion between the charged proteins, and a shorter-ranged attractive contribution, $u_{att}(r; \Omega_1, \Omega_2)$, owing to van der Waals and anisotropic (patchy) hydrophobic interactions. Molecular dynamics simulations [21] revealed that the discrete distribution of protein surface charges matters at high salt concentrations only (typically > 1.0 M). For lower salt content, the electric repulsion is well described for $r > \sigma$ by the isotropic screened Coulomb potential (see Ref. [22] for a derivation),

$$\beta u_{\rm rep}(r) = l_{\rm B} Z^2 \left(\frac{\exp\{\kappa a\}}{1+\kappa a}\right)^2 \frac{\exp\{-\kappa r\}}{r}, \qquad (3)$$

constituting the electrostatic part of the Derjaguin-Landau-Verwey-Overbeek (DLVO) pair potential frequently used in colloid science (see, e.g., [23]). Here, $l_{\rm B} = e^2/(\epsilon_{\rm s}k_BT)$ is the Bjerrum length of the solvent with static dielectric constant $\epsilon_{\rm s}$. For water at room temperature, $l_{\rm B} = 0.7$ nm. Note that Z is an *effective* protein charge number which accounts for the reduced screening ability of the small counter-and coions (microions) surrounding each protein at larger protein concentration, causing Z to be somewhat larger than the actual (bare) surface charge number [3, 24–26]. There is, in principle, the opposite effect of counter-ion condensation [27] on the surface of strongly charge particles which reduces the value of Z. For concentrated proteins, Z is often so small that the first effect prevails. The screening of the electric repulsion between two proteins by intervening microions causes $u_{\rm rep}(r)$ to decay exponentially, with the repulsion range, $1/\kappa$, determined by the Debye electric screening parameter κ where

$$\kappa^{2} = \frac{4\pi l_{B}}{(1-\phi)} \left(\rho |Z| + 2\rho_{s}\right).$$
(4)



Fig. 2: Left: Sketch of particles motion induced stream lines illustrating the solvent-flow mediated inter-particles HIs which act quasi-instantaneously and are long-ranged in general. **Right:** Schematic particle mean-squared displacement, $\langle \Delta \mathbf{R}(t)^2 \rangle$, of directly and hydrodynamically interacting globular proteins in solution (red curve). See text for discussion.

The contribution to κ^2 proportional to the protein concentration ρ is due to monovalent counterions released from the protein surfaces, while the second contribution is due to dissociated monovalent salt ion pairs of concentration ρ_s . The factor $1/(1-\phi)$ invoking the protein volume fraction $\phi = (\pi/6)\rho\sigma^3$ corrects for the free volume accessible by the point-like microions. As shown in the following section for the example of BSA solutions, the spherical colloid model applies well for lower protein and added salt concentrations. Protein shape anisotropy and surface patchiness are relevant at large protein concentrations and larger salt content where these refined particle features influence, in particular, the protein solution dynamics and phase behavior.

A protein moving in the viscous solvent in presence of other ones creates a flow pattern that affects not only the motion of neighboring particles, but via hydrodynamic back-reflections from other particles also its own motion is changed. On the time scale $t \gtrsim \tau_d = a_{hyd}^2/d_0$ where significant diffusion-driven changes in the protein positions take place, these solvent-mediated HIs are quasi-instantaneously transmitted [28]. Here,

$$d_0 = \frac{k_B T}{6\pi\eta_0 a_{\text{hyd}}} \tag{5}$$

is the single-particle Stokes-Einstein diffusion coefficient of an isolated spherical protein in a solvent of shear viscosity η_0 at temperature T, and a_{hyd} is its hydrodynamic particle radius which can differ from the effective hard-core radius $a = \sigma/2$ [29]. For BSA in water, $\tau_d \approx 0.3$ μ s. HIs are long-ranged and non-pairwise additive at larger concentrations, meaning that the HIs between a pair of particles are changed in presence of a third one in their vicinity. These features make HIs difficult to deal with both in theoretical calculations and mesoscale dynamic simulations. The HIs between N spheres at instant configuration $X = \{R_1, \ldots, R_N\}$ of the center positions $\{R_i\}$ are described by coupled linear force-velocity relations [28],

$$\mathbf{V}_{i} = \sum_{j=1}^{N} \boldsymbol{\mu}_{ij}(X) \cdot \mathbf{F}_{j}, \qquad (6)$$

where $\mu_{ij}(X)$ is a three-dimensional mobility matrix relating the instant hydrodynamic force, \mathbf{F}_j , exerted on a particle j by the surrounding fluid to the resulting velocity change, \mathbf{V}_i , of a

particle i. Owing to the HIs, V_i and F_i are not colinear in general (see left part of Fig. 2). Eq. (6) is valid on coarse-grained time and length scales where accelerated (inertial) motions of particles and fluid are not resolved any more. Acceleration effects are operative on time scales much short than τ_d during which the proteins have moved a tiny distance only compared to their size (typically ~ 1 ps for BSA in water). The hydrodynamic mobility coefficients $\mu_{ii}(X)$ are salient input to theoretical and Brownian Dynamics simulation methods for calculating diffusion transport properties. They have been obtained, based on the linearized Navier-Stokes equation describing low-Reynolds number solvent flow, in form of multipole expansions in powers of inverse interparticle distances, combined with a lubrication correction for nearly touching spheres [30]. We refer to Ref. [28] for a pedagogical discussion of HIs in colloidal systems. Note that the underlying description of the solvent as a hydrodynamic continuum is useful not only for colloidal particles and proteins, but also for solutions of small electrolyte ions [31]. The influence of HIs is nicely illustrated by the self-diffusion of a tagged Brownian particle in the crowded environment of other ones, by examining in Fig. 2 the mean-squared displacement (MSD), $\langle \Delta \mathbf{R}(t)^2 \rangle$, of fluid-like correlated proteins as a function of time t, where $\Delta \mathbf{R}(t) = \mathbf{R}(t) - \mathbf{R}(0)$ is the positional change after the time span t, and $\langle \cdots \rangle$ denotes an equilibrium average over particle positions. For short times $t \ll \tau_d$, a particle has diffused only a tiny distance compared to its size, and the dynamically formed cage of neighboring particles (which on average is spherically symmetric) has practically not changed. The tagged particle diffuses thus in the potential minimum of neighboring ones, so that it is slowed down by the instantaneously acting HIs only, implying a linear initial temporal increase $\langle \Delta \mathbf{R}(t)^2 \rangle \approx 6 \, d_s \, t$, quantified by the short-time self-diffusion coefficient d_s which is smaller than the single-particle coefficient d_0 . At intermediate times $t \sim \tau_d$, the cage is slightly distorted from spherical symmetry so that the direct (i.e., potential-based) interactions of Eq. (1) become operative, implying a less than linear increase of the MSD [32,33]. The sub-linear increase of the MSD for intermediate times is approximately described by a fractional power law (MSD $\propto t^{\alpha}$), with an exponent α smaller than one. At long times $t \gg \tau_d$, the MSD increases again linearly with slope $6 d_l$. The long-time self-diffusion coefficient d_1 is smaller than d_s , since both direct and hydrodynamic interactions with neighboring particles contribute now to the slowing down. The linear longtime form of the MSD is very slowly approached, typically like $(6 d_1 t - \langle \Delta \mathbf{R}(t)^2 \rangle) \sim 1/\sqrt{t}$. This makes simulation calculations of d_1 elaborate even when HIs are neglected [32, 34]. The inequalities $d_1 < d_8 < d_0$ are valid in crowded protein solutions for all kinds of pair potentials. For solutions where the proteins form a crystal, glass or gel, d_1 is basically equal to zero, and the MSD reaches a long-time plateau value termed localization length.

3 Crowding and Salt Effects in BSA Solutions

We proceed by discussing the concentration dependence of collective diffusion and solution viscosity in aqueous BSA solutions, at low to moderately large concentrations of added NaCl as explored in [6, 11]. Data from dynamic light scattering (DLS) and rheometry are compared with theoretical calculations based on the spherical colloid model, for a pair potential according to Eq. (1). For the considered BSA and added salt concentrations where the electric repulsion is strong and of longer range, van der Waals and possibly existing patchy hydrophobic attractions described by $u_{\text{att}}(r; \Omega_1, \Omega_2)$ can be neglected, so that the pair potential is spherically symmetric. For the calculation of dynamic properties, easy-to-implement methods are available which account for protein-protein HIs [12, 14, 35]. The only input required by these methods is the



Fig. 3: Left: Scattering setup for probing particle correlations on the length scale $2\pi/q$. **Right:** Small-angle x-ray scattering intensity I(q) (in arbitrary units) of an aqueous BSA solution at mass concentration $\rho \times m_p = 4.5$ mg/ml, and added NaCl concentration $c_s = 1.3$ mM. Right figure taken from [6].

so-called protein static structure factor S(q), which can be determined from static small-angle x-ray scattering, on positing that the BSA proteins can freely rotate. Before discussing dynamic solution properties, we describe first the physical meaning of S(q), and how it is obtained from scattering experiments.

3.1 Static scattering and equilibrium structure

The schematic setup of a static x-ray or neutron scattering experiment is shown in Fig. 3 (left part). X-ray photons or thermal neutrons are used since their wavelengths are comparable to the protein size, allowing thus to infer information about the internal structure (form) of the proteins. Monochromatic radiation of wavelength λ impinges on a solution sample, and is scattered by the proteins at the angle ϑ into a detector measuring the average intensity, I(q), of scattered photons or neutrons. Here, \mathbf{k}_i and \mathbf{k}_s are the wavevectors of incident and scattered radiation, respectively, and $\mathbf{q} = \mathbf{k}_s - \mathbf{k}_i$ is the scattering wavevector of magnitude $q = (4\pi/\lambda) \sin(\vartheta/2)$. A measurement of I(q) at the specific wavenumber q resolves structural details of the solution on the length scale $2\pi/q$.

For a highly dilute solution where the proteins diffuse independently, $I(q) \propto \overline{f(q, \Omega)}^2$. Here, $f(q, \Omega)$ is the scattering amplitude of a protein with orientation Ω , and the overline denotes orientational averaging on assuming equal probabilities of all particle orientations. From low-concentration x-ray measurements, it was found that BSA proteins behave like homogeneously scattering prolate ellipsoids of revolution, with the long semi-axis of revolution b = 4.75 nm and the short semi-axis c = 1.75 nm. The effective protein hard-core diameter σ used in the spherical colloid model can be identified from demanding the second virial coefficient of hard ellipsoids,

$$B_2^{\text{ell}}(T) = -2\pi \int_0^\infty dr \, r^2 \left(\left\langle \exp\left[-\beta u_{\text{hc}}^{\text{ell}}(r; \mathbf{\Omega}_1, \mathbf{\Omega}_2)\right] \right\rangle_{\mathbf{\Omega}_1, \mathbf{\Omega}_2} - 1 \right) \,, \tag{7}$$

where $u_{\rm hc}^{ell}$ is the hard-core pair potential of two ellipsoids, to be equal to the second virial coefficient, $B_2^{\rm hc} = (2\pi/3)\sigma^3$, of genuine hard spheres with diameter σ . For BSA this gives



Fig. 4: In decoupling approximation, BSA proteins are treated as freely rotating ellipsoids regarding their scattering properties, and as uniformly charged effective spheres regarding their direct and hydrodynamic interactions. Figure taken from [11].

 $\sigma = 7.40$ nm. In Eq. (7), $\langle \cdots \rangle_{\Omega_1,\Omega_2}$ is the unbiased average over the solid angles of a pair of hard ellipsoids labeled 1 and 2. The second virial coefficient is the quadratic-order coefficient in the concentration Taylor expansion of the osmotic pressure of a one-component solution, and it can be determined experimentally by static light scattering and osmometry [23, 36].

We can account for the non-spherical shape of BSA proteins in the calculation of I(q) by using the translation-orientation decoupling approximation where the proteins are described as freely rotating ellipsoids regarding their scattering amplitudes, and as charge-stabilized hard spheres of diameter σ as far as their interactions are concerned. The decoupling approximation steps are pictured in Fig. 4.

The scattering intensity of a BSA solution is calculated in decoupling approximation according to

$$I(q)/\rho \approx \overline{f(q,\Omega)}^2 S(q) + \left(\overline{f^2(q,\Omega)} - \overline{f(q,\Omega)}^2\right).$$
(8)

Here, S(q) is the static structure factor of proteins modeled as effective spheres interacting by the screened Coulomb plus hard-core pair potential $u_{rep}(r) + u_{hc}(r)$. The structure factor quantifies concentration fluctuation correlations of wavelength $2\pi/q$ (see left part of Fig. 3). It is basically the spatial Fourier transform of the protein pair distribution function g(r), i.e.

$$S(q) = 1 + \rho \int d\mathbf{r} \exp\{i\mathbf{q} \cdot \mathbf{r}\} \left[g(r) - 1\right].$$
(9)

The value of g(r) is the conditional probability of finding a sphere at center-to-center distance r from another one. In simulations of $N \gg 1$ spheres in the primary simulation box of volume V, g(r) is calculated as an equilibrium average over particles configurations according to

$$\rho g(r) = \left\langle \frac{1}{N} \sum_{i \neq j}^{N} \delta(\mathbf{r} - \mathbf{R}_{i} + \mathbf{R}_{j}) \right\rangle.$$
(10)

Here, δ is the three-dimensional delta function, $r = |\mathbf{r}|$, and \mathbf{R}_i is the position vector of particle *i*. An example of g(r) for a solution of spheres in a cluster-fluid phase state is presented in the following section.

Correlations between the orientational and translational degrees of freedom of BSA proteins arising from their residual non-sphericity are neglected in the decoupling approximation. They are of importance at high salt content and for larger protein concentrations only. That the decoupling approximation of I(q) is well suited for BSA solutions is illustrated in Fig. 3 (right part) where small-angle x-ray scattering (SAXS) intensity data for the mass concentration $\rho \times m_p = 4.5$ ml/mg, with m_p denoting the protein mass, are compared with the theoretical prediction. The only adjusted parameter in the theoretical intensity is the protein effective charge number in $u_{\rm rep}(r)$, determined here as Z = 19, on using results for S(q) and g(r) calculated by an accurate analytic method (see Refs. [6,37] for details). The maximum of I(q) at $q \approx 3$ nm⁻¹ is a correlation (crowding) feature related to the principal maximum, $S(q_m)$, of the structure factor at the wavenumber q_m characterizing the linear extension, $\sim 2\pi/q_m$, of nearest neighbor cages of proteins formed around each protein sphere in crowded solutions (cf. Fig. 6, right part).

3.2 Short-time collective diffusion and solution viscosity

Different from self-diffusion which stands for the Brownian motion of a particle in the statistically homogeneous environment of other ones, collective diffusion for a given wavenumber q is the relaxation of thermally induced, low-amplitude sinusoidal concentration fluctuations of wavelength $2\pi/q$ caused by the coordinated motion of many particles. Experimentally, information about collective diffusion is contained in the dynamic extension, S(q, t), of the static structure factor (with S(q, 0) = S(q)), referred to as dynamic structure factor or intermediate scattering function. The function S(q, t) quantifies spatio-temporal correlations of concentration fluctuations of wavelength $2\pi/q$. It is the spatial Fourier transform,

$$S(q,t) = \int d\mathbf{r} \exp\{i\mathbf{q}\cdot\mathbf{r}\}G(r,t), \qquad (11)$$

for q > 0, of the *cumulative* van Hove function

$$G(r,t) = \left\langle \frac{1}{N} \sum_{i,j}^{N} \delta(\mathbf{r} - \mathbf{R}_{i}(t) + \mathbf{R}_{j}(0)) \right\rangle.$$
(12)

The latter is the conditional probability density of finding, at time t, a particle i at distance r from another one $(j \neq i)$ or itself (j = i) at the earlier time t = 0. Notice that $G(r, 0) = \rho g(r) + \delta(\mathbf{r})$, i.e. up to a concentration factor, G(r, t) is the time-dependent generalization of the static pair distribution function but with self-correlations (j = i) included.

The dynamic structure factor is the key quantity measured in neutron spin echo (NSE), x-ray photon correlation spectroscopy scattering (XPCS) and dynamic laser light scattering (DLS) experiments, for the respectively accessed wavenumber and correlation time intervals.

Short-time collective diffusion is probed experimentally by measuring S(q,t) for correlation times t small compared to the single-particle diffusion time τ_d , during which S(q,t) decays exponentially according to

$$S(q, t \ll \tau_{\rm d}) = S(q) \exp\left[-q^2 D(q) t\right] .$$
(13)

Here,

$$D(q) = d_0 \frac{H(q)}{S(q)} \tag{14}$$



Fig. 5: Left: Collective diffusion coefficient, d_c , of aqueous BSA solutions as function of protein mass concentration $\rho \times m_p$, and for zero and physiological NaCl concentrations ρ_s . **Right:** Low-shear solution viscosity, η , divided by the solvent (water) viscosity η_0 . Comparison of experimental data (symbols) and theoretical curves. Figures adjusted from [6].

is the short-time collective diffusion function characterizing the initial decay of concentration fluctuations of wavelength $2\pi/q$, given by the ratio of the so-called hydrodynamic function H(q) and S(q). The hydrodynamic function quantifies the influence of HIs on short-time diffusion, and it can be calculated as the equilibrium average

$$H(q) = \left\langle \frac{1}{N\mu_0} \sum_{l,j=1}^{N} \left(\hat{\boldsymbol{q}} \cdot \boldsymbol{\mu}_{lj}(X) \cdot \hat{\boldsymbol{q}} \right) \exp\{i \boldsymbol{q} \cdot (\boldsymbol{R}_l - \boldsymbol{R}_j)\} \right\rangle,$$
(15)

where $\mu_0 = d_0/k_B T$ is the single particle mobility coefficient, and $\hat{q} = q/q$ the unit vector in direction of q. Without HIs, H(q) is a constant equal to one, independent of q and the protein volume fraction ϕ . Deviations between D(q) and $d_0/S(q)$ are thus a hallmark of the influence of HIs. For large q where short distances are probed and cross-correlations in Eq. (15) are absent, H(q) is equal to the short-time self-diffusion coefficient d_s in units of d_0 . For small wavenumbers where only macroscopically large distances are resolved, it can be shown that $H(q \to 0)$ is equal to the mean (short-time) sedimentation velocity, in units of the singlesphere sedimentation velocity, of monodisperse Brownian spheres that are slowly settling in the direction of q, under the action of a weak constant force field [28]. Consequently,

$$d_{\rm c}^{(s)} = d_0 \, \frac{H(q \to 0)}{S(q \to 0)} \tag{16}$$

is the (short-time) *collective* diffusion coefficient, quantifying the initial relaxation of a constant density gradient. In principle, one should distinguish the short-time coefficient from its long-time analogue $d_{\rm c}^{(l)}$ for which $d_{\rm c}^{(l)} \leq d_{\rm c}^{(s)}$, and which is defined by $d_{\rm c}^{(l)} = -(d/dt) \log S(q,t)/q^2$ in the hydrodynamic limit $q \to 0$ and $t \to \infty$ with q^2t being constant. This collective long-time coefficient has its appearance in the celebrated constitutive Fick law,

$$\mathbf{j}(\mathbf{r},t) = -d_{\mathbf{c}}^{(l)} \nabla \rho(\mathbf{r},t), \qquad (17)$$

relating on a macroscopic scale the particles diffusion flux, **j**, to the driving concentration gradient. However, since the equilibrium distribution of weakly sedimenting, monodisperse particles remains practically unperturbed even in a strongly crowded environment, $d_c^{(l)}$ is practically equal to $d_c^{(s)}$. Therefore, we can simply refer to a single collective diffusion coefficient d_c , without having to distinguish between short- and long-time values. This is in sharp contrast to self-diffusion where for concentrated systems d_l is substantially smaller than d_s , e.g. $d_l/d_s \approx 0.1$ at the freezing concentration where particles start to form crystals.

We point out that $S(q \to 0) = k_B T (\partial \rho / \partial \Pi)_T$ is proportional to the isothermal osmotic solution compressibility, with $\Pi(\rho, T)$ denoting the osmotic solution pressure. For strongly repelling particles, the osmotic compressibility and hence S(0) decrease strongly with increasing concentration, giving rise to values of $d_{\rm c}$ significantly larger than the one of d_0 . An example in case is given in Fig. 5 (left part), showing experimental and theoretical results for the (mass) concentration dependence of $d_{\rm c}$, for aqueous BSA solutions with vanishing and physiological amounts of added salt. The experimental data were obtained from DLS short-time measurements for which, owing to the wavelengths of the optical photons long compared to the protein size, S(q,t) is probed for $q \ll q_m$. Consider first the result for d_c without added salt: At small concentrations, S(0) decreases faster with growing ρ than H(0). While S(0) accounts for the direct interactions influence on d_c via the osmotic compressibility, H(0) accounts for the slowing influence of the HIs. After a sharp increase at low concentrations, $d_{\rm c}$ passes through a distinct maximum and decreases then gradually with further growing ρ , as a consequence of the dominant influence of HIs for larger concentrations. The maximum of d_c is located roughly at the concentration where the number of counterions released from the protein surfaces matches that of the residual electrolyte ion pairs in the solution different from dissociated NaCl ions. With increasing amount of added NaCl, the maximum of $d_{\rm c}$ becomes smaller and shifts to larger ρ . At physiological salt conditions where electrostatic screening is strong, d_c grows only mildly and nearly linearly with increasing concentration. The depicted theoretical predictions for d_c were obtained using a versatile method of calculating D(q), termed BM-PA method as the acronym for Beenaker - Mazur plus Pairwise Additivity HIs method that combines analytic simplicity with high accuracy. Details about this method are included in Refs. [29, 35] and references therein. In view of the accuracy of the BM-PA method, the remaining differences between the experimental and theoretical data for $d_{\rm c}$ can be ascribed to the simplifying description of the BSA proteins as charged spheres. The extrapolation of the experimental d_c data to zero concentration leads to a value for $d_0 = d_c(\rho \to 0)$ larger than the one used in the theoretical model according to Eq. (5) with $\sigma = 7.40$ nm. If instead of the theoretical d_0 value the experimentally extrapolated one is used, good agreement between experimental and theoretical data for d_c is achieved.

In the right part of Fig. 5, experimental data (symbols) for the low-shear solution viscosity η are depicted. According to

$$\frac{F}{A} = \eta \frac{\partial u}{\partial y}, \tag{18}$$

the viscosity relates the force per area (shear stress), F/A, applied to the moving upper plate of the shear cell, to the resulting constant gradient, $\partial u/\partial y$, of the solution flow velocity u(y)inside the cell (see figure inset). As seen, the viscosity grows monotonically with increasing concentration, and it is smaller with than without added salt. With added salt, the nearestneighbor cages are less rigid, and they can be sheared with less free energy consumption. While in decent overall agreement, the theoretical predictions based on the spherical colloid model underestimate the viscosity reduction induced by the addition of salt. The theoretical method of calculating η is explained in Ref. [6].





Fig. 6: Left: Experimental - theoretical test of the Kholodenko-Douglas GSE relation, for BSA solutions of salt concentrations $\rho_s = 0$ and 150 mM. **Right:** Schematic form of the structure factor S(q) of purely repulsive particles. Left figure taken from [6].

3.3 Generalized Stokes-Einstein relations

The methods for calculating d_c and η require S(q) as the only input (see Fig. 6, right part). They can be further used to scrutinize, for crowded protein solutions, the validity of a generalized Stokes-Einstein (GSE) relation proposed by Kholodenko and Douglas (KD) [38],

$$\frac{d_{\rm c}(\rho)\eta(\rho)}{d_0\eta_0}\sqrt{S(q\to 0;\rho)} \approx 1, \qquad (19)$$

relating the viscosity to the collective diffusion coefficient, and to the square-root of S(0). For zero protein concentration, the GSE relation reduces to the single-particle Stokes-Einstein relation in Eq. (5). The (approximate) validity of a GSE relation such as the one by Kholodenko and Douglas is of interest not only from a theoretical viewpoint since it also has experimental applications. Solutions of laboriously isolated proteins are often available in small amounts only not sufficient for a mechanical rheometry measurement. A valid GSE relation could be then profitably used to infer η indirectly from scattering measurement of d_c and S(0), since in scattering experiments only small amounts of proteins are required.

The left part of Fig. 6 includes an experimental-theoretical test of the KD-GSE relation applied to BSA solutions. While the relation holds decently well for high-salinity solutions, the bad news is its distinct violation for low-salt conditions and non-zero concentrations. We mention that additional GSE relations between rheology and diffusion properties have been proposed in the literature, and critically scrutinized. A general finding for GSE relations is that their applicability depends crucially on the range and character of the pair potential [12, 35, 39].

4 Short-range attractive and long-range repulsive Proteins

In many protein solutions, the direct interactions include, under appropriate low salt conditions, both short-range attraction (SA) and long-range electrostatic repulsion (LR) contributions of comparable strengths. Examples in case are lysozyme [8], and some monoclonal antibody protein solutions [40, 41]. The competing interaction contributions in these so-called SALR



Fig. 7: Left: Hard-sphere plus double-Yukawa (HSDY) potential for parameters $z_1 = 10$, $z_2 = 0.5$, $K_1 = 1.63$, and various values of $\alpha = K_1 - K_2$ with $\alpha = -\beta u(\sigma^+)$. **Right:** Generalized Lennard-Jones plus repulsive Yukawa (LJY) potential, with $\nu = 50$, A = 2, $\xi = 1.8$ and attraction strength values ϵ as indicated. Dotted red lines in left part are SA and LR potential contributions for $\alpha = 1$. From [13] and [14], respectively.

systems can lead to to the thermodynamically *reversible* formation of particle clusters of preferred size. SA triggers particles aggregation which continues until a cluster has accumulated enough repulsion strength (i.e. electric charge) to prevent its further growth. Depending on different combinations of the SA and LR contributions, and the protein concentration, the solution can have different equilibrium and non-equilibrium phase states, including dispersed-fluid and cluster-fluid phases, and random and cluster-percolated states [7, 12–17].

There are two interaction potentials that have been widely used to describe theoretically the phase behavior and structure of SALR systems (see Fig. 7), The first one is the hard-sphere plus double-Yukawa (HSDY) pair potential [7, 12, 13]

$$\beta u(x) = \begin{cases} \infty, & x \le 1 \\ -K_1 \frac{e^{-z_1(x-1)}}{x} + K_2 \frac{e^{-z_2(x-1)}}{x}, & x > 1, \end{cases}$$
(20)

where $x = r/\sigma$ is the inter-particle distance, r, in units of the particle diameter σ , and $\beta = 1/k_{\rm B}T$. Moreover, z_1 and z_2 determine the range of the attractive and repulsive Yukawa potential parts in units of σ , respectively, and K_1 and K_2 are the respective short-range attractive and long-range repulsive potential strengths in units of $k_{\rm B}T$. The depth of the attractive well quantifying the net attraction is $\alpha = -\beta u(x = 1^+) = K_1 - K_2$. As seen in the left part of Fig. 7, with increasing α a shallow potential barrier develops at $x_{\rm max} \approx 1.3$, followed for $x > x_{\rm max}$ by the slowly decaying LR potential part of range $1/z_2 = 2$ which is 20 times larger than the SA range $1/z_1 = 0.1$. Structural and dynamic predictions based on this spherical particles protein model have been recently compared with NSE, small-angle neutron scattering, and rheometry data for aqueous lysozyme solutions. Under low salt conditions and at smaller protein concentrations, good agreement is observed [42].

The second SALR potential is the Lennard-Jones-Yukawa (LJY) potential. It combines the usual repulsive Yukawa potential describing the long-range, weakly screened electrostatic repulsion with a generalized Lennard-Jones potential describing the short-range attraction and



Fig. 8: Left: $\tilde{T} - \phi$ generalized phase diagram for HSDY and LJY systems with competing short-range attraction and long-range repulsion (SALR). Right: Corresponding cluster size distribution functions N(s). Figures taken from [7].

repulsion [13, 14, 16, 17]. Explicitly,

$$\beta u(x) = 4\epsilon \left[\left(\frac{1}{x}\right)^{2\nu} - \left(\frac{1}{x}\right)^{\nu} \right] + \frac{A\xi}{x} \exp\left(-x/\xi\right) \,. \tag{21}$$

This potential is well suited for mesoscale dynamic simulation methods such as the multiparticles collision dynamics (MPC) method for which systems with $\nu = 50$ have been studied [14]. Here again is $x = r/\sigma$, with σ being now the soft particle diameter, ϵ the strength of short-range attraction and repulsion, and A characterizing the strength of long-range repulsion. The repulsion strength A is proportional to the square of the effective protein charge Ze, and ξ is accordingly the electrostatic Debye screening length in units of σ (cf. Eq. (3)). Fig. 7 displays the shapes of the two SALR potentials, for selected attraction and repulsion strength parameters given in the caption. The minimum of the depicted LJY potential curves is located at $x_{\min} \approx 1.01$, at which $\beta u(x_{\min}) \approx 2 - \epsilon$. Thus, the repulsion range is about 18 times longer than the attraction range.

4.1 Phase behavior and structural properties

A generalized temperature - volume fraction $(\tilde{T} - \phi)$ phase diagram for the two considered SALR model systems is shown in Fig. 8 (left part), which is taken from Ref. [7]. Here, $\phi = (\pi/6)\rho\sigma^3$ is the particle volume fraction, and \tilde{T} is a reduced temperature representing the ratio of thermal energy and strength of attraction, i.e. $\tilde{T} = k_B T/u(x = 1^+)$ for the HSDY potential, and $\tilde{T} = 1/\epsilon$ for the LJY potential. The reduced temperature and volume fraction are divided by the respective values, \tilde{T}_c and ϕ_c , of the critical point of the gas-liquid binodal curve of the *reference* system, which includes only the purely attractive short-range portion of the full SALR potential. The gas-liquid phase coexistence line (binodal) of the reference system with pure attraction is drawn as the black line in the diagram. See Ref. [7] for details about the phase diagram calculation, and the selected reference system potential. Four different phases are distinguished in the figure. For large \tilde{T} and small ϕ , there is a dispersed-fluid phase (blue region) where most particles are not part of a cluster, i.e. they are monomers. Increasing ϕ leads to the formation of a random percolated state (green area). For the reference system with attractive interactions only, state points (\tilde{T}, ϕ) located below the binodal line would cause



Fig. 9: Left: Simulation snapshot of a LJY equilibrium-cluster phase system at $\phi = 0.05$, A = 2, $\xi = 1.8$ and $\epsilon = 7$. **Right:** Pair distribution function g(r) at $\phi = 0.05$ and $\epsilon = \{7, 8\}$. From [14].

a macroscopic liquid-gas like phase separation into a protein-rich and protein-poor phase. In SALR systems, however, the phase separation is frustrated by the LR repulsion, leading for lower ϕ to the formation of a cluster-fluid phase consisting of equilibrium particle clusters of preferred size (red area). The size distribution, and the dynamic shapes and lifetimes of the clusters depend on the selected system parameters. For larger volume concentrations, namely for $\phi \gtrsim 0.35 \times \phi_c$, the clusters percolate into a system-spanning network (yellow area).

The discussed phase states can be identified by the characteristic shape of the respective cluster size distributution (CSD) function N(s) which can be calculated, e.g., using Monte-Carlo [7] and MPC simulations [14] where configurations of N_p particles in a periodically replicated simulation box are generated.

The CSD function is the average fraction of particles (proteins) contained in a cluster formed by s particles (i.e. with cluster size s). It is defined by

$$N(s) = \left\langle \frac{s}{N_{\rm p}} n(s) \right\rangle, \tag{22}$$

where $\langle \ldots \rangle$ is the average over representative particle configurations, and n(s) is the number of clusters of size s within a given configuration. Note that $\sum_{s=1}^{N_p} N(s) = 1$. A cluster is defined by applying a distance criterion according to which the center-to-center distance, $|\mathbf{R}_{ij}|$, between two particles i and j of the same cluster obeys the condition $|\mathbf{R}_{ij}| < r_{\text{cluster}}$, by using an appropriately defined cut-off distance r_{cluster} [7, 14]. As noticed in the right part of Fig. 8, the dispersed-fluid phase is distinguished by a monotonically decreasing N(s), representing a state where monomers are the most abundant species in the system. The CSD function of the cluster-fluid phase has a local maximum at the preferred cluster size, whereas the N(s) of the random percolated state is overall slowly decaying, with a pronounced system-size peak at $s = N_p$. The random percolated state is defined according to the standard criterion that at least 50% of the sampled configurations contain a system-spanning cluster. In comparison, the CSD function of the cluster-percolated phase state has a preferred cluster size peak, in addition to the system-size peak at $s = N_p$. The formation of percolated clusters is a prerequisite of gelled or glassy non-equilibrium states.

Fig. 9 (left part) includes a configurational snapshot of a LJY system in the cluster-fluid phase, generated in a MPC simulation [14]. The various colors indicate particles belonging to clusters

of different sizes. The corresponding CSD function (not shown here) has a pronounced peak at the preferred cluster size s = 7. The right part of Fig. 9 displays the corresponding orientationally averaged particle-particle pair correlation function g(r), for the cluster-fluid LJY systems of interaction strengths $\epsilon = 7$ and 8, indicating in both cases strong local ordering. The figure displays the preferred geometric particle configurations associated with respective peaks in g(r), such as an octahedral arrangement of four particles (peak at $x \approx 1.66$) and a cubic arrangement (peak at $x \approx 1.43$). Enlarging ϵ from 7 to 8 sharpens and raises the peaks in g(r), and leads to the appearance of additional ones. As seen in the figure, owing to SA the relative likelihood, $g(x = 1^+)$, of finding two particles at contact distance x = 1 is very large.

4.2 Dynamics in cluster-fluid and dispersed-fluid phases

In comparison with the large body of work on the structure and phase behavior of SALR dispersions, comparably little is known about their dynamic behavior. The main challenge for simulation and theoretical studies is the consideration of many-particles HIs which strongly influence the dispersion and intra-cluster dynamics. We discuss here simulation results for collective diffusion in LJY cluster-fluid phase systems obtained using the multiparticles collision dynamics (MPC) method. This is a dynamic simulation technique which combines the Molecular Dynamics evolution of particles, interacting here by a SALR potential, with a coarse-grained treatment of solvent particles by a so-called stochastic rotational dynamics evolution. The MPC method correctly reproduces the low-Reynolds-number hydrodynamic solvent flow with associated many-particles HIs, and the thermal Brownian particles concentration fluctuations, for a somewhat reduced hydrodynamic particle radius. Refs. [14, 43, 44] include detailed information about this method. There are alternative simulation methods that can be used for calculating dynamic properties of directly and hydrodynamically interacting Brownian particles systems. They all have their respective pros and cons. We only mention Brownian Dynamics simulation techniques (see Refs. [2, 35, 39, 45, 46]) where the solvent-mediated HIs are accounted for by hydrodynamic mobility matrices such as those in Eq. (6).

MPC simulation results (symbols) for the inverse of the collective diffusion function D(q) of LJY dispersions are shown in Fig. 10, for various attraction strength values ϵ . The D(q)'s have been determined from the short-time form of the calculated dynamic structure factor S(q,t) using Eq. (13). The small-wavenumber peak in $d_0/D(q)$ at $q \sim 2/\sigma$ indicates intermediate range, cluster-like correlations. The height of this peak grows strongly with increasing attraction strength. A small-q peak is likewise present in S(q), and in the hydrodynamic function H(q) (not shown here). The function H(q) has pronounced undulations, signaling the influence of HIs on collective diffusion. The overall slowing down of collective diffusion in SALR systems due to HIs is illustrated by the comparison of $d_0/D(q)$ with $S(q) = d_0/D(q)|_{\text{no-HIs}}$ (dashed lines). Neglecting HIs leads to an overestimation of D(q) by a factor of about 2. The LJY systems for $\epsilon \leq 6$ are in the dispersed-fluid phase state for which there is quantitative agreement between the simulation data and the theoretical predictions for D(q) based on the BM-PA method (solid curves in the inset). This method, however, is not applicable to cluster-fluid phase systems for which elaborate simulation calculations are required.

The evolution of a cluster of seven LJY particles generated in a MPC simulation is shown in the right part of Fig. 10. The cluster undergoes substantial shape changes during the considered time span, with a particle leaving the cluster at time $t = 19 \tau_d$. The simulation sequence illustrates that clusters are dynamical objects that change their shape and particle number in the course of time. Quite importantly, the HIs between an ensemble of mobile particles are not screened,


Fig. 10: Left: MPC simulation results (symbols) for the inverse of the short-time diffusion function D(q) at $\phi = 0.05$, and different attraction strength values ϵ . For $\epsilon = \{7, 8\}$, the dispersions are in the cluster-fluid phase state (filled symbols), and for $\epsilon = \{2, 4, 6\}$ in the dispersed-fluid state. Dashed curves: S(q) for $\epsilon = \{6, 8\}$. Solid curves in the inset: BM-PA theory results for dispersed-fluid phase systems. **Right:** MPC generated time sequence of a 7-particles cluster, tracked until a particle labeled in green is dissociated from the cluster. Adjusted from [14].

15

i.e. the range of the HIs is not shortened with increasing concentration of mobile particles. Only the interaction strength is reduced. Consequently, a cluster-fluid dispersion cannot be simply treated as a polydisperse system of rigid cluster bodies. The mean lifetime of clusters is significantly affected (prolonged) by the HIs, and it grows strongly with increasing attraction strength [14].

5 Patchy Colloid Model of Lysozyme

10

 $q\sigma$

While the description of globular proteins as effective spheres with spherically symmetric interaction is useful for protein solutions of lower salt content and particle concentration, the comparison of theoretical predictions for D(q) and η with NSE scattering and rheometry measurements on lysozyme solutions clearly shows that for larger salinity and concentration the slowing influence of coupled rotational-translational motion caused by non-sphericity and patchy interaction needs to be accounted for [42] (see also Ref. [47] for platelet dispersions). The importance of this coupling was shown additionally in MPC simulations of dispersions of spherical Brownian particles with short-range patchy attraction [48].

Additionally to dynamic properties, the phase behavior of globular protein solutions is affected by surface patchiness at larger salt content. We discuss now a sphere model of aqueous lysozyme solutions which combines a charge-induced isotropic repulsion potential with a non-isotropic patchy attraction potential, $u_{\text{att}}(r; \Omega_1, \Omega_2)$, that is of Yukawa-type in its radial dependence, and of Kern-Frenkel type in its orientational modulation [49]. Experimentally known values for the temperature T_c , volume fraction ϕ_c , and second virial coefficient $B_2(T_c)$ of the gas-liquid type critical point of lysozyme solutions are used to determine the strength and range of the patchy attraction. Moreover, using second-order thermodynamic perturbation theory, the metastable gas-liquid coexistence (binodal) and spinodal curves, and the liquid-crystal coexistence curve are determined for this model [3, 18].



Fig. 11: Left: Simplification steps: From an atomistic model of lysozyme to a coarse-grained model where hydrophobic surface regions are colored in red, and finally to a sphere model with uniform surface charge distribution and spherical hydrophobic surface patches (in red). **Right:** Two patchy proteins at vector distance \mathbf{r} , and with solid angles Ω_1 and Ω_2 characterizing their respective orientation. There is attractive interaction only if \mathbf{r} intersects a hydrophobic patch on each sphere. Taken from [3, 18].

5.1 Kern-Frenkel type model

Lysozyme is approximately an ellipsoidal polypeptide of volume $v_0 = 20.58 \text{ nm}^3$. It is treated here as a spherical particle of equal volume v_0 , which determines the effective diameter in the hard-core contribution to the total potential in Eq. (1) as $\sigma = 3.6 \text{ nm}$. The isotropic electrostatic repulsion is described by $u_{\text{rep}}(r)$ according to Eq. (3), using pH = 7.8 and the associated bare charge number Z = 8 for zero protein concentration. As noted earlier, the effective charge number Z increases above this bare value with increasing protein concentration, by less than ten percent at $\phi = 0.4$ [18].

The anisotropic attraction of two lysozyme proteins due to their hydrophobic surface patches is described by the Kern-Frenkel type potential (see Fig. 11)

$$\beta u_{\text{att}}(r; \mathbf{\Omega}_1, \mathbf{\Omega}_2) = -\tilde{\epsilon}_{\text{att}} \left(1 + \psi \left[1 - \frac{T}{T_c} \right] \right) \frac{\exp\{-z_{\text{att}} \left(r/\sigma - 1 \right) \}}{r/\sigma} \times d(\mathbf{\Omega}_1, \mathbf{\Omega}_2) , \quad (23)$$

for $r > \sigma$, with

$$d(\mathbf{\Omega}_1, \mathbf{\Omega}_2) = \begin{cases} 1, & \text{if } \mathbf{r} \text{ intersects two hydrophobic patches} \\ 0, & \text{otherwise.} \end{cases}$$
(24)

In using Eq. (23), we assume that the radial dependence of the patchy attraction, which is of range $1/z_{\text{att}}$ in units of σ , can be factorized from the function $d(\Omega_1, \Omega_2)$ depending on the relative orientation of the two particles. The temperature-dependent factor in the above potential can be considered as a simple approximation for the temperature dependence of attractive hydrophobic interactions. Here, T_c is the experimentally well-assessed critical temperature of liquid-gas like coexistence. The attraction strength $\tilde{\epsilon}_{\text{att}} > 0$ and ψ are two physical parameters determined by the experimental critical point.

5.2 Phase diagram calculation

To calculate the gas-liquid and fluid-solid coexistence curves of lysozyme solutions of larger salt content, one can employ second-order perturbation theory in the so-called compressibility approximation [50,51]. In this way, time-consuming computer simulation based phase diagram calculations are avoided.

Owing to the significant electrostatic screening induced by the added salt ions, we can use hard spheres as the reference system, by expanding the free energy, F(N, V, T), of the protein system up to quadratic order in the orientationally pre-averaged perturbation potential $u_p = u - u_{hc} = u_{rep} + u_{att}$. This results in (see Refs. [3, 23] for details)

$$f(T,\phi) \approx f_{\rm hc}(T,\phi) + \frac{12\phi^2}{\sigma^3} \int_0^\infty dr r^2 g_{\rm hc}(r) \langle \beta u_{\rm p}(r;\Omega_1,\Omega_2) \rangle_{\Omega_1,\Omega_2} - \frac{6\phi^2}{\sigma^3} \left(\frac{\partial\phi}{\partial\Pi_{\rm hc}}\right)_{\rm T} \int_0^\infty dr r^2 g_{\rm hc}(r) \langle \left[\beta u_{\rm p}(r;\Omega_1,\Omega_2)\right]^2 \rangle_{\Omega_1,\Omega_2}.$$
(25)

Here, $f = (\beta F/N)\phi$ is the free energy per protein, in units of k_BT and multiplied by ϕ , and Π_{hc} is the dimensionless osmotic pressure of the reference hard-sphere system, in units of the ideal gas pressure at the same temperature and protein concentration. We refer to f as the free energy, for brevity. For simplicity, the unbiased orientational average, $\langle \cdots \rangle_{\Omega_1,\Omega_2}$, over the solid angles Ω_1 and Ω_2 is taken for the perturbation potential.

The equilibrium coexistence of two phases requires equal temperature T, osmotic pressure Π , and chemical potential μ in both phases. The latter two conditions determine the protein volume fractions of the coexisting phases. Therefore, the fluid (f) - solid (s) phase coexistence curve is determined from the conditions:

$$\Pi_{\rm f}(T,\phi_{\rm f}) = \Pi_{\rm s}(T,\phi_{\rm s}) \quad \text{with} \quad \Pi = \phi^2 \left(\frac{\partial(f/\phi)}{\partial\phi}\right)_{\rm T}$$
(26)

$$\mu_{\rm f}(T,\phi_{\rm f}) = \mu_{\rm s}(T,\phi_{\rm s}) \quad \text{with} \quad \beta\mu = \left(\frac{\partial f}{\partial \phi}\right)_{\rm T}.$$
(27)

At sufficiently low T, a liquid-like (l) and a gas-like (g) protein solution phase of high and low volume fraction ϕ_l and ϕ_g , respectively, are coexisting, with the liquid-gas coexistence curve (binodal) determined by the conditions

$$\Pi_{l}(T,\phi_{l}) = \Pi_{g}(T,\phi_{g}) \text{ and } \mu_{l}(T,\phi_{l}) = \mu_{g}(T,\phi_{g}).$$
 (28)

To compute the coexistence curves from these conditions using the free energy perturbation expression in Eq. (25), the Newton-Raphson method with line search can be used (see Ref. [3] for details). For the super-critical hard-sphere (hc) reference system, accurate analytic expressions are known for the reduced free energy $f_{hc}(T, \phi)$, and the orientationally averaged pair distribution function $g_{hc}(r)$ in both the fluid [52–54] and fcc phase regions [55, 56]. These expressions are used as input in the free energy expression. Different free energy and pair distribution function expressions are needed for the fluid and crystalline phases since the orientational symmetry is reduced in going from the fluid to the solid phase. Note that solid lysozyme solutions have a tetragonal crystal structure. In our simplifying model, however, the ellipsoidal lysozyme proteins are described as patchy spheres for which the hard-sphere reference system with its fcc crystal phase can be used.

The (upper) critical point (T_c, ϕ_c) of the liquid-gas coexistence curve in the $T - \phi$ phase diagram is a concentration saddle point of the free energy $f(T, \phi)$, where its second and third derivatives with respect to ϕ vanish simultaneously. The spinodal curve of diverging isothermal compressibility encloses the mechanically unstable fluid phase region, and it is determined by $\partial^2 f(T, \phi)/\partial \phi^2 = 0$.



Fig. 12: Left: Phase diagram of an aqueous lysozyme solution at pH = 7.8, and with 0.5 M added NaCl. Circles: Experimental metastable gas-liquid coexistence curve (binodal). Triangles: Experimental fluid-crystal coexistence curve. Solid lines: Calculated coexistence curves using $\psi = 5$. The spinodal separates the mechanically unstable region from the metastable one below the binodal. Right: Gas-liquid coexistence curves at four different salt concentrations. Symbols: Experimental binodal data. Solid curves: Theoretical predictions using constant $\psi = 5$. Taken from [3, 18].

A characteristic property of the patchy attraction potential is the surface coverage χ , which is the fraction of the sphere surface covered by the *m* attractive spherical caps,

$$\chi = m \sin^2 \left(\delta/2 \right) \,, \tag{29}$$

where 2δ is the opening angle of a red patchy cone (see Fig. 11, right part). Only the square of χ appears in the employed orientational average of u_p , with all the details of the patchy surface averaged out. In the phase diagram calculations discussed here, m = 2 is used, and δ is treated as adjustable parameter in addition to the strength $\tilde{\epsilon}_{att}$ and range z_{att} of the attraction potential part ($\psi = 5$ is fixed). The potential parameters $\{z_{att}, \tilde{\epsilon}_{att}, \delta\}$ are determined from matching the critical point values (T_c, ϕ_c) of the gas-liquid phase coexistence curve to the experimental ones. As a third constraint determining the three unknown parameters, it is required for the calculated second virial coefficient that

$$B_2(T_c) = -2.7 \times B_2^{hc}, \qquad (30)$$

independent of the salt concentration for values $\rho_s > 0.25$ M, with $B_2^{hc} = (2\pi/3)\sigma^3$ denoting the second virial coefficient of the hard-sphere reference system. This requirement is supported by the experimental finding, being in accord with extended principle of corresponding states [36,57], that there is a narrow band of values $b_2(T_c) \equiv B_2(T_c)/B_2^{hc} = -2.7 \pm 0.2$ for which the solution separates into gas- and liquid-like phases [58,59]. For the considered lysozyme solutions of higher salt content, the liquid-gas coexistence is *metastable* with respect to fluid-solid phase coexistence. The metastability follows also from the present phase diagram calculations. The constraint in Eq. (30) is reasonable, since $f(T, \phi) = f_{id}(T, \phi) + 4b_2(T)\phi^2 + O(\phi^3)$ for smaller ϕ values, with $f_{id}(T, \phi)$ denoting the ideal gas free energy. Any viable free energy expression such as the one in Eq. (25) should reproduce the exact low-concentration form. The protein osmotic pressure near the critical point is rather low, so that it can be approximated using a small- ϕ expansion. This makes it plausible why the value of $b_2(T_c)$ provides a criterion for the gas-liquid critical point. The parameter ψ determines essentially the width of the gas-liquid binodal. It is fixed here to $\psi = 5$, independent of the salt concentration.

In the left part of Fig. 12, the calculated phase diagram is shown for $\rho_s = 0.5$ M, and compared with the corresponding experimental finding for lysozyme solutions [3]. The potential parameter values determined by fitting the experimental critical point on the binodal are $z_{att} = 3.0$, $\tilde{\epsilon}_{att} = 3.1$, and $\delta = 73^{\circ}$ corresponding to a surface coverage of 71%. The surface coverage of the patches, the range of attraction of about 30% of the protein diameter, and the temperature dependence of u_{att} are consistent with experimental results. The range of hydrophobic attraction, in particular, is in excellent agreement with force measurements between two hydrophobic plates [60]. The calculated binodal and spinodal curves agree well with the experimental data for $\phi < 0.25$, while for larger concentrations some deviations are seen. The calculated fluidcrystal coexistence curve is in decent accord with the experimental data (triangles). In region (f), a stable fluid phase is observed. There is a fluid-crystal coexistence region in (f + s), and a metastable gas-liquid coexistence region in (g + 1). A purely crystalline phase region is found to the right of the nearly vertical red line at $\phi \approx 0.44$ in the phase diagram.

In the right part of Fig. 12, the calculated gas-liquid binodals are depicted for various added salt concentrations, illustrating the good agreement with the experimental data points. As expected, the binodal broadens, and is shifted towards larger temperature values with increasing salt concentration.

The consistent results obtained from the discussed patchy spheres model point to the importance of patchy hydrophobic attraction between the lysozyme proteins for the phase behavior.

6 Conclusions and Outlook

We have discussed the dynamics, structure, and phase behavior of concentrated solutions of globular proteins. The purpose of this lecture was to show that theoretical and simulation methods developed originally for colloidal suspensions can be applied successfully to crowded protein solutions, with the proteins described in a minimalist way as (patchy) Brownian spheres. A thorough quantitative analysis of non-spherical shape and patchiness effects in crowded protein solutions can be an important task in future simulation studies where these effects are individually explored.

In contrast to crowded living cells, we emphasize that the presented material deals with Brownian particles systems at global thermodynamic equilibrium. A living cell operates out of equilibrium, having *active* exchange of energy and matter with its surrounding, e.g. through ion channels operating in response to external and internal stimuli. Additionally to passive diffusional transport in the cytosol, there is active transport of macromolecules by motor proteins operating along microtubular tracks. Metabolic activities, e.g., can fluidize large macromolecular components in the cytoplasm that otherwise would be in a glass-like state [33]. Studies of the dynamics and phase behavior of concentrated protein solutions such as those discussed in this lecture contribute to reveal and quantify means to distinguish living from dead matter.

For conciseness, we dealt here with the dynamics and structure of single-component protein solutions. A cell with its huge number of different biological macromolecules is a highly multicomponent and polydisperse system in which additional diffusion mechanisms are operative. Regarding self-diffusion, each protein species α is characterized by its own MSD, and by respective short-time and long-time self-diffusion coefficients $d_{s\alpha}$ and $d_{l\alpha}$. The latter can vary strongly among the various species, e.g. when particles of some species form a glass or gel, while the particles of more weakly interacting species are still able to diffuse through the glass or gel interstices [61, 62]. A limiting case of multi-component self-diffusion is the *tracerdiffusion* of a species being so dilute that its particles are correlated with particles of other species only. The generalization of collective diffusion to multi-component systems is termed *cooperative diffusion*, and is characterized by various cooperative (partial) diffusion coefficients $D_{c\alpha\beta}^{(s,l)}$ relating a concentration gradient in species β to the diffusion current of species α . Different from collective diffusion in concentrated one-component systems, long-time cooperative diffusion coefficients are distinctly different than their short-time counterparts. A special cooperative diffusion mechanism is interdiffusion, meaning the relaxation of thermal fluctuations in the relative concentration of two tagged species in a many-species mixture. Additionally to the translational diffusion mechanisms discussed above, the proteins perform rotational Brownian motion in the crowded solution environment [4,63]. The coupling of translational and rotational diffusion can be quite strong for anisometric particles, and spherical particles with anisotropic soft interactions, thereby affecting the transport of matter. For an elementary treatment of interdiffusion and rotational diffusion, see Ref. [54].

There are also interesting rheological effects in multispecies systems. For instance, the mixing of small and large particles in a dispersion of higher salt content significantly lowers the viscosity, owing to a HIs mechanism where small particles are dragged along hydrodynamically by big ones in their vicinity [64]. This hydrodynamic effect of reducing the viscous dissipation can play a role in the passive and active material transport inside a cell.

Acknowledgment

I am grateful to my former PhD students M. Heinen (Univ. of Guanajuato, Mexico), J. Riest (Fa. Zeiss, Oberkochen) and Ch. Gögelein (Fa. Arlanxeo, Dormagen), and to R.G. Winkler (ICS-2, FZ Jülich), for fruitful collaborations on crowded protein solutions.

References

- [1] W. Häußler, A. Wilk, J. Gapinski, and A. Patkowski, J. Chem. Phys. 117, 413 (2002)
- [2] J. Gapinski, A. Wilk, A. Patkowski, W. Häußler, A.J. Banchio, R. Pecora and G. Nägele, J. Chem. Phys. 123, 054708 (2005)
- [3] Ch. Gögelein, G. Nägele, R. Tuinier, T. Gibaud, A. Stradner and P. Schurtenberger, J. Chem. Phys. 129, 085102 (2008)
- [4] S.R. McGuffee and A.H. Elcock, PLoS Computational Biology 6, e1000694 (2010)
- [5] P. Szymczak and M. Cieplak, J. Phys. Condens. Matter 23, 033102 (2011)
- [6] M. Heinen, F. Zanini, F. Roose-Runge, F. Zhang, M. Hennig, T. Seydel, R. Schweins, M. Sztucki, M. Antalik, F. Schreiber and G. Nägele, Soft Matter 8, 1404 (2012)
- [7] P.D. Godfrin, N.E. Valadez-Perez, R. Castaneda-Priego, N.J. Wagner and Y. Liu, Soft Matter 10, 5061 (2014)

- [8] V.L. Dharmaraj, P.D. Godfrin, Y. Liu, and S.D. Hudson, Biomicrofluidics 10, 043509 (2016)
- [9] A.H. Elcock, Current Opinion in Structural Biology 20, 196 (2010)
- [10] S. Kondrat, O. Zimmermann, W. Wiechert, and E. v. Lieres, Phys. Biol. 12, 046003 (2015)
- [11] M. Heinen, Charged Colloids and Proteins: Structure, Diffusion and Rheology, PhD Thesis, Forschungszentrum Jülich GmbH, Zentralbibliothek, Key Technologies Vol. 32, Jülich (2011)
- [12] J. Riest and G. Nägele, Soft Matter 11, 9273 (2015)
- [13] J. Riest, Dynamics in Colloid and Protein Systems: Hydrodynamically structured Particles, and Dispersions with competing attractive and repulsive Interactions, PhD Thesis, Forschungszentrum Jülich GmbH, Zentralbibliothek, Key Technologies Vol. 127, Jülich (2016)
- [14] S. Das, J. Riest, R.G. Winkler, G. Gompper, J.K.G. Dhont and G. Nägele, Soft Matter 14, 92 (2018)
- [15] A. Stradner, H. Sedgwick, F. Cardinaux, W.C.K. Poon, S. U. Egelhaaf and P. Schurtenberger, Nature 432, 492 (2004)
- [16] F. Sciortino, S. Mossa, E. Zaccarelli, and P. Tartaglia, Phys. Rev. Lett. 93, 055701 (2004)
- [17] J.C.F. Toledano, F. Sciortino and E. Zaccarelli, Soft Matter 5, 2390 (2009)
- [18] C. Gögelein, Phase Behavior of Proteins and Colloid-Polymer Mixtures, PhD Thesis, Forschungszentrum Jülich GmbH, Jülich (2008)
- [19] RCSB Protein Data Bank: https://www.rcsb.org/pdb/home/home.do
- [20] C. Tanford and R. Roxby, Biochemistry **11**, 2192 (1972)
- [21] E. Allahyarov, H. Löwen, A.A. Louis and J.P. Hansen, Europhys. Lett. 57, 731 (2002)
- [22] G. Nägele, *Theories of Fluid Microstructures*, in: Soft Matter: From Synthetic to Biologicl Materials, 39th IFF Spring School 2008. Series: Key Technologies Vol. 31, Forschungszentrum Jülich Publishing, Jülich, 2006
- [23] Ch. Gögelein, D. Wagner, F. Cardinaux, G. Nägele and S.U. Egelhaaf, J. Chem. Phys. 136, 015102 (2012)
- [24] L. Belloni, J. Chem. Phys. 85, 519 (1986)
- [25] H. Ruiz-Estrada, M. Medina-Noyola and G. Nägele, Physica A 168, 919 (1990)
- [26] A-P. Hynninen and M. Dijkstra, J. Chem. Phys. **123**, 244902 (2005)
- [27] E. Trizac, L. Bocquet, M. Aubouy and H. H. von Grünberg, Langmuir 19, 4027 (2003)

- [28] G. Nägele, *Colloidal Hydrodynamics*, in Physics of Complex Colloids, edited by C. Bechinger, F. Sciortino, and P. Ziherl, Vol. 184, Proceedings of the International School of Physics 'Enrico Fermi', page 451, IOS Press, Amsterdam; SIF, Bologna, 2012
- [29] J. Riest, T. Eckert, W. Richtering and G. Nägele, Soft Matter 11, 2821 (2015)
- [30] D.J. Jeffrey and Y. Onishi, J. Fluid Mech. 139, 261 (1984)
- [31] C. Contreras Aburto and G. Nägele, J. Chem. Phys. 139, 134110 (2013)
- [32] G. Nägele, *Brownian Dynamics Simulations*, in: Computational Condensed Matter Physics, S. Blügel, G. Gompper, E. Koch, H. Müller-Krumbhaar, R. Spatschek, and R. G. Winkler, (Eds.). Forschungszentrum Jülich Publishing, Jülich, 2006
- [33] B.R. Parry, I.V. Surovtsev, M.T. Cabeen, C.S. O'Hern, E.R. Dufresne and C. Jacobs-Wagner, Cell 156, 183 (2014)
- [34] B. Cichocki and K. Hinsen, Physica A 187, 133 (1992)
- [35] M. Heinen, A.J. Banchio and G. Nägele, J. Chem. Phys. 135, 154504 (2011)
- [36] F. Platten, N.E. Valadez-Perez, R. Castaneda-Priego and S.U. Egelhaaf, J. Chem. Phys. 142, 174905 (2015)
- [37] M. Heinen, P. Holmqvist, A.J. Banchio and G. Nägele, J. Chem Phys. 134, 044532 & 129901 (2011)
- [38] A.L. Kholodenko and J.F. Douglas, Phys. Rev. E 51, 1081 (1995)
- [39] A. J. Banchio and G. Nägele, J. Chem. Phys. 128, 104903 (2008)
- [40] P. D. Godfrin, I. E. Zarraga, J. Zarzar, L. Porcar, P. Falus, N. J. Wagner and Y. Liu, J. Phys. Chem. B 120, 278 (2016)
- [41] E. J. Yearley, I. E. Zarraga, S. J. Shire, T. M. Scherer, Y. Gokarn, N. J. Wagner and Y. Liu, Biophys J 105, 720 (2013)
- [42] J. Riest, G. Nägele, Y. Liu, N.J. Wagner and P.D. Godfrin, submitted (2017)
- [43] S. Poblete, A. Wysocki, G. Gompper and R. G. Winkler, Phys. Rev. E 90, 033314 (2014)
- [44] J.T. Padding and A.A. Louis, Phys. Rev. E 74, 031402 (2006)
- [45] A.J. Banchio and J.F. Brady, J. Chem. Phys. 118, 10323 (2003)
- [46] P. Mereghetti and R.C. Wade, J. Phys. Chem. B 116, 8523 (2012)
- [47] D. Kleshchanok, M. Heinen, G. Nägele and P. Holmqvist, Soft Matter 8, 1584 (2012)
- [48] S. Bucciarelli, J. S. Myung, B. Farago, S. Das, G. A. Vliegenthart, O. Holderer, R.G. Winkler, P. Schurtenberger, G. Gompper and A. Stradner, Science Advances 2, e1601432 (2016)
- [49] N. Kern and D. Frenkel, J. Chem. Phys. **118**, 9882 (2003)

- [50] J.A. Barker and D. Henderson, J. Chem. Phys. 47, 2856 (1967)
- [51] Ch. Gögelein, F. Ramano, F. Sciortino and A. Giacometti, J. Chem. Phys. 136, 094512 (2012)
- [52] N.F. Carnahan and K.E. Starling. J. Chem. Phys. 51, 635 (1969)
- [53] L. Verlet and J.-J. Weis. Phys. Rev. A 5, 939 (1972)
- [54] G. Nägele, *The Physics of Colloidal Soft Matter*, Lecture Notes 14, Institute of Fandamental Technological Research, Polish Academy of Sciences Publishing, Warsaw, ISSN 1642-0578, 2004
- [55] W.W. Wood, J. Chem. Phys. 20, 1334 (1952)
- [56] J.M. Kincaid and J.J. Weis. Mol. Phys. 34, 931 (1977)
- [57] M.G. Noro and D. Frenkel, J. Chem. Phys. 113, 2941 (2000)
- [58] P.B. Warren, J. Phys.: Condens. Matter 14, 7617 (2002)
- [59] C.K. Poon, S.U. Egelhaaf, P.A. Beales, A. Salonen and L. Sawyer, J. Phys.: Condens. Matter 12, L569 (2000)
- [60] J. Israelachvili and R. Pashley, Nature **300**, 341 (1982)
- [61] G. Nägele, J. Bergenholtz and J.K.G. Dhont, J. Chem. Phys. 110, 7037 (1999)
- [62] T. Voigtmann and J. Horbach, Phys. Rev. Lett. 103, 205901 (2009)
- [63] K. Makuch, M. Heinen, G.C. Abade and G. Nägele, Soft Matter 11, 5313 (2015)
- [64] R.A. Lionberger, Phys. Rev. E 65, 061408 (2002)

B4 Membrane Channels & Pumps

J.-P. Machtens Cellular Biophysics Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	In	troduction	2		
2	Tr	ansport across biological membranes			
	2.1	Organization and function of biological membranes	2		
	2.2	Energetic basis of channels, transporters, and pumps	2		
	2.3	Molecular basis of transport	4		
3	El	ectrical signaling by ion channels			
	3.1	The resting potential	8		
	3.2	Gating of ion channels	9		
	3.3	The action potential	9		
4	Μ	olecular simulations of membrane transport proteins			
	4.1	Computational Electrophysiology simulations	10		
	4.2	Identification of the Cl ⁻ channel mechanism of glutamate transporters	12		
R	References				

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved

1 Introduction

The lipid bilayer is a perfect electrical insulator and provides the basis for molecular information processing at the cell membrane. Electrical signaling requires both ion-selective and dynamically regulated membrane permeabilities as well as specific concentration gradients across the membrane. Ion channels are membrane proteins responsible for creating selective ion permeabilities, whereas pumps or active transporters establish ionic gradients.

Ion channels contain molecular pores that selectively permit certain ions—for example Na⁺, K⁺, Ca²⁺, or Cl⁻—to passively cross the cell membrane along their electrochemical potential gradient. In contrast, pumps or transporters actively produce ion concentration gradients. An important example is the Na⁺-K⁺-ATPase, which hydrolyzes ATP to energize concentrative accumulation of Na⁺ and K⁺ in the outside and inside of the cell, respectively. In summary, transporters establish concentration gradients to drive ionic currents through ion channels to induce voltage signals at the cell membrane.

2 Transport across biological membranes

2.1 Organization and function of biological membranes

Every living cell is surrounded by a cell membrane that defines the cell as a unit, receives and releases signal molecules, and thus mediates communication between the intracellular and the extracellular space. Biological membranes are based on a phospholipid bilayer. This design is of exceptional simplicity and furthermore provides unique mechanical and electrical properties. The fluidity of the lipid membrane permits cell division and movement. The lipid bilayer is permeable to gases and hydrophobic molecules, but restricts the permeation of ions and polar or charged molecules—such as sugars, amino acids, and sometimes even water; the bilayer thereby effectively defines the intra- and extracellular milieu. Thus, to meet the cell's requirements, polar substances and ions have to be selectively transported across the membrane by specialized membrane proteins: ion channels, transporters, and ATP-driven pumps.

The impermeability of the lipid bilayer for charged molecules makes it a perfect electrical insulator and represents the basis for electrical signaling. Moreover, the electrical capacitance of a lipid bilayer is very low. For a typical cell with a diameter of 30 μ m, the cell capacitance is only 30 pF. This value implies that only 2.4 pC, that is about $1.5 \cdot 10^7$ elementary charges, need to be moved across the membrane to generate transmembrane voltages of about 80 mV. This number of ions can be moved across the membrane within one second even by a single ion channel (see below). Changing the transmembrane potential is fast and does not require major energy consumption. In summary, the physicochemical properties of the lipid membrane enable ion channels to play a key role in electrical signaling in virtually all cells of the body.

2.2 Energetic basis of channels, transporters, and pumps

Membrane transport proteins mediate either passive or active transport of solutes across the membrane, and they are usually categorized into channels, pumps, and transporters. Moreover, transporters can be subdivided into uniporters, symporters, or antiporters (Fig. 1). From an energetic perspective, ion channels and uniporters mediate passive (energetically favorable) transport of solutes down their concentration or electrical potential gradients. In contrast, pumps, symporters, and antiporters perform active transport of ions or small molecules "uphill" against their electrochemical gradient. From a structural perspective,

pumps and transporters undergo significant conformational changes per transported substrate, whereas ion channels—once opened—permit the transmembrane flux of several thousands of ions without further protein conformational changes per individual transported ion.



Fig. 1: Overview of membrane transport proteins. Electrochemical gradients are indicated by triangles. (A) Pumps, transporters, and ion channels are the three major types of transport proteins. Pumps harness the energy released by ATP hydrolysis to drive transport of specific ions or molecules (red circles) against their electrochemical gradient. Transporters catalyze the movement of specific ions or molecules and can be classified as uniporters, symporters, or antiporters. Ion channels enable rapid and passive movement of specific ions (or water) down their electrochemical gradient. (B) Uniporters facilitate the passive movement of a single substrate down its electrochemical gradient. Symporters and antiporters catalyze secondary active transport of one substrate against its gradient, using the energy released by the movement of another ion down its gradient using a fixed transport stoichiometry.

Ion channels permit the fastest transport of ions across the membrane, reaching transport rates of up to more than 10^8 ions per second [1]. They are usually involved in fast ion transport in the generation of electrical signals by changing the cell's membrane potential. Channels differ from transporters and pumps in the simplicity of their transport function (Fig. 1). They are characterized by an aqueous conduction pathway that permits the diffusion of ions from one membrane side to the other. Only passive ion diffusion is possible via ion channels. There are two driving forces for ion diffusion: the concentration gradient and the voltage. These two driving forces can be combined to the electrochemical gradient:

$$\Delta G = RT \ln \frac{c_i}{c_o} + zFV$$

Here, R is the gas constant, T is the absolute temperature, c_i and c_o are intra- and extracellular ion concentrations, respectively, z is the charge of the ion, F the Faraday constant, and V the transmembrane voltage. Ion channels conduct ions in one direction at negative electrochemical potentials and in the other direction at positive values. At an electrochemical potential of 0 V, ion efflux and influx are identical, resulting in zero net ion transport.

Transporters and pumps differ from ion channels in the mechanisms underlying transport. In these two groups of membrane transport proteins, substrates are moved across the membrane via a conformational change of the protein. The linchpin of their transport mechanism is that the substrate binding site is made accessible to either the inner or to the outer side of the membrane. Structural transitions are then thought to occur between inward- and outward-facing states. Thus, transporters and pumps use an alternating-access mechanism to move substrates across the membrane [2]. The conformational changes underlying alternating accessibility can be discrete, for example by providing alternating access of substrates to or from a central binding site to the internal and external space by changing the direction of only few amino acid side chains. Other transporters use large-scale elevator-like conformational changes of proteins occur at much longer timescales than diffusion processes in ion channels, transport rates of transporters and pumps are much lower than those of ion channels.

Uniporters are transporters, which passively catalyze substrate transport following the electrochemical gradient across the membrane. This process is also commonly referred to as facilitated diffusion. Symporters and antiporters move more than one substrate: the transport of one substrate occurs passively, and the transport along its electrochemical gradient provides the energy to move another substrate against its electrochemical gradients. They are called secondary active transporters because they use the preexisting gradient (generated by a primary active transport process, see below) of a driving substrate as a source of energy to actively transport another solute against its gradient. According to the relative transport directions of the substrates, symporters (cotransporters) and antiporters (exchangers) are distinguished. Symporters move all substrates in the same direction, whereas antiporters mediate the simultaneous import of one substrate and export of another substrate. Secondary active transporters normally operate with a fixed transport stoichiometry—a certain number of one substrate is only transported in co- or countertransport with a defined number of another substrate. Mechanistically, secondary active transporters are very similar to passive transporters: They are thermodynamically reversible and the transport direction depends on the combined electrochemical potentials of the driven substrate and the driving substrate.

Pumps are ATPases, which are specialized transporters that move substrates against their electrochemical gradient using the energy provided by the hydrolysis of ATP to ADP and P_i. This process is called primary active transport. The most important pump, the Na⁺-K⁺-ATPase keeps intracellular [Na⁺] low and [K⁺] high, with an exchange stoichiometry of 3 Na⁺ and 2 K⁺ per hydrolyzed ATP molecule. Other physiologically important pumps are the Ca²⁺-ATPase and the H⁺-K⁺-ATPase.

2.3 Molecular basis of transport

For many decades ion channels and transporters could only be studied on a functional level. Detailed structural insights that are the prerequisite for understanding transport processes at near atomic resolution have not existed until about 20 years ago. The first high-resolution structure of a membrane transport protein was solved from the bacterial K^+ -selective ion channel KcsA from *Streptomyces lividans* [3]. This structure explained the basis for selective ion conduction through biological membranes (Fig. 2).

Electrical signaling in mammalian cells is primarily based on regulated membrane permeabilities for the two main cations in biological media, Na^+ and K^+ . The plasma membranes of resting cells are usually Na^+ impermeable, but highly conductive for K^+ . Cell excitation is based on the transient opening of Na^+ -permeable channels. These processes thus require ion channels that reliably select between these two cations. This important task is complicated by the fact that Na^+ and K^+ carry identical charges and only differ slightly in their diameter (Na^+ , 1.9 Å; K^+ , 2.6 Å).

High-resolution structures of K^+ -selective channels demonstrate how these requirements can be fulfilled (Fig. 2). KcsA exhibits a short narrow constriction—the selectivity filter—in which ions can only enter after complete dehydration. K^+ channels are tetrameric and each subunit carries a short glycine-tyrosine-glycine (GYG) sequence in the narrow part of the selectivity filter. The carbonyl oxygens of the tyrosine side chain and of the inner glycine perfectly substitute for interactions of the K^+ ions with water molecules. The smaller Na⁺ ions interact less favorably with the selectivity filter and are thus excluded from entering the channel.



Fig. 2: Crystal structure of the KcsA potassium channel shown from within the membrane (left panel) or from outside (right panel). Blue circles depict permeating K^+ ions that are either entering, leaving, or bound to the GYG motif within the selectivity filter. In the side view, two of the four subunits of the tetrameric assembly are removed for better illustration.

This design permits very tight interactions of K^+ ions with the channel protein. Since effective electrical signaling requires high transport rates, selective ion channels do not only have to bind ions with high affinity, but also to permit rapid binding and unbinding. In K^+ -selective channels, tight binding is overcome by the occupation of the selectivity filter with multiple K^+ ions. K^+ ions that bind within the selectivity filter coming from outside or inside the cell destabilize binding of the other ions and permit the release of the outermost ion to the opposite side of the membrane. The enormous complexity of generating a K^+ -selective ion channel is illustrated by the fact that only one mechanism for K^+ selectivity has been developed during evolution. In contrast, there are multiple protein families and designs for channels selective for Na⁺, Ca²⁺, and Cl⁻ ions.

Transporters and pumps function by a completely different mechanism as ion channels. Whereas ions permeate through an aqueous conduction pathways of ion channels, ion transporters move ions across the membrane via conformational changes of the protein to ensure alternating accessibility of the substrate binding sites to each site of the membrane [2,

4]. In recent years, our mechanistic understanding of the function of transporters has been significantly advanced by experimentally determined structures of such proteins in multiple conformations (e.g. by X-ray crystallography, cryo-electron microscopy, or nuclear magnetic resonance spectroscopy) and by molecular dynamics simulations. In the following I will illustrate secondary active transport on one particular example, Na⁺-coupled glutamate transport. There are multiple other transporter families that function using similar mechanisms.

Glutamate is the major excitatory neurotransmitter in the mammalian central nervous system. After release from presynaptic nerve terminals, glutamate diffuses across the synaptic cleft and binds to specialized postsynaptic receptors. Activation of these receptors results in the generation of electrical or chemical signals in the postsynaptic neuron. Glutamatergic synaptic transmission is terminated by the uptake of glutamate into surrounding glial and neuronal cells. This transport is the main physiological task of Na⁺-coupled glutamate transporters (excitatory amino acid transporters, EAATs) of the SLC1 family [5]. Coupling glutamate and Na⁺ transport ensures the establishment of low extracellular glutamate concentrations, resulting in high sensitivity of glutamatergic synaptic transmission.

High-resolution crystal structures of the prokaryotic glutamate transporter Glt_{Ph} from *Pyrococcus horikoshii* in multiple conformations [6-12] illustrated how these transport processes occur and how movement of diverse substrates are coupled. Glt_{Ph} is assembled as trimer from three identical subunits [13] that operate independently of each other (Fig. 3). Each subunit exhibits eight transmembrane helices with two hairpin loops (HP1 and HP2) that open and close during transport.



Fig. 3: Backbone fold of Glt_{Ph} shown from the outside. Glt_{Ph} is a trimer assembled from three identical subunits. In each subunit, the static trimerization domain (green) is distinguished from the mobile transport domain (orange) with HP1 (yellow) and HP2 (red).

Transport is initiated by binding of two Na⁺ ions to the *apo* state (Fig. 4). These binding steps cause the extracellular gate, HP2 to open, thereby enabling association of glutamate or aspartate, one H⁺, and the remaining Na⁺ ion. After closure of HP2 an elevator-like conformational change occurs, the translocation of the transport domain (orange in Fig. 3, yellow in Fig. 4) across the membrane. This translocation involves an ~18 Å rotational-translational rigid body movement of the substrate-bound transport domain relative to the membrane. After release of the substrates in the inward-facing conformation the transporter returns via re-translocation to its original conformation.



Fig. 4: The EAAT glutamate/aspartate transport cycle illustrated by structures of the prokaryotic homolog Glt_{Ph} . EAATs transport glutamate together with three Na^+ and one H^+ in counter-transport with one K^+ . T_o represents the outward-facing conformation of a transporter monomer. Binding of two Na^+ (red; T_o-Na_2) is followed by glutamate (shown in purple) and H^+ binding, and association of a third Na^+ (T_o-Na_3 -GluH). Subsequently, the transport domain (yellow) moves via an intermediate conformation ($T_{int}-Na_3$ -GluH) with an elevator-like movement towards the inward-facing conformation (T_i-Na_3 -GluH). The substrate (T_i-Na_3) and sodium is then released (T_i) and retranslocation involves antiport of one K^+ . (Coordinates are from Protein Data Bank IDs: T_o , 40YE; T_o-Na_2 , 40YF; T_o-Na_3 -GluH, 2NWX; $T_{int}-Na_3$ -GluH, ref. [14]; T_i-Na_3 -GluH, 3KBC; T_i-Na_3 , 4P6H; T_i , 4P3J).

Although transport occurs at much lower rates in transporters than in channels, transporters have to solve the same conflict like a selective channel: to selectively bind the right substrate, while performing effective transport. Glt_{Ph} illustrates one possible solution to this problem: an induced-fit mechanism is used to select between different amino acids [15]. Transport substrates (glutamate or aspartate) bind after opening of HP2 first loosely to the transporter. The subsequent closure of HP2 prevents dissociation of the substrate and is necessary for the transport of the substrate across the membrane (Fig. 5). At the other membrane side, opening of the other hairpin loop, HP1, could then enable release of the substrate association and HP2 closure—are key to the amino acid selectivity of glutamate transporters because they permit selection of the main substrate from other substrates that bind initially tighter to the transport domain and can therefore not be transported at high rates.



Fig. 5: Induced-fit mechanism of substrate binding and selection in the EAAT homolog Glt_{Ph}.

3 Electrical signaling by ion channels

3.1 The resting potential

Every cell exhibits a voltage across the plasma membrane and most intracellular membranes. When measured in cells that are not generating electrical signals, this voltage is usually referred to as the resting potential. Measuring these voltages requires electrical access to the inside of the cell without major impairment of the electrical isolation by the membrane, for example by impaling a fine glass electrode into the cell. For many cells, for example skeletal muscle or glial cells, one can observe a voltage of about -80 mV after impalement and short regeneration of the membrane. In these cells, the membrane potential critically depends on the extracellular potassium concentration $[K^+]$. Plotting the resting membrane potential of a skeletal muscle fiber as a function of extracellular $[K^+]$, one observes a dependence that can be almost perfectly described by the Nernst equation [16]. This equation describes diffusion potentials at the interface between solutions with different ion concentrations. From this mathematical description, it was concluded that the cellular resting potential is—in many cells—a K⁺ diffusion potential and does not directly depend on energy consuming processes.

One can study the generation of a K⁺ diffusion potential in a model cell that only exhibits K⁺-selective ion channels and physiological internal and external [K⁺]. In this model cell, K⁺ will initially diffuse out of the membrane driven by the concentration gradient. Since only K⁺ can cross the membrane, K⁺ efflux will result in charge separation and in the development of a transmembrane voltage difference. Net K⁺ efflux will stop when the chemical and the electrical driving force compensate each other so that the electrochemical gradient ΔG for K⁺ ion is zero: $\Delta G = RT \ln \frac{c_i}{c} + zFV = 0$

At this moment net, K^+ flux across the membrane stops, and ion concentrations as well as the membrane potential are stable resulting in the Nernst equation:

$$V = \frac{RT}{zF} \ln \frac{c_o}{c_i}$$

In the human body, intracellular K^+ concentrations are about 150 mM, due to the high density of the primary active Na⁺-K⁺-ATPase. The extracellular $[K^+]$ is tightly regulated to low levels, with concentrations between 3.6–5.6 mM. Under these conditions, cells with a plasma membrane, which is selectively permeable for K^+ ions, will exhibit a K^+ diffusion potential close to the experimentally determined values. In summary, the resting potential in many cells is dominated by the K^+ diffusion potential due to the existing $[K^+]$ gradient and the abundant presence of potassium channels.

3.2 Gating of ion channels

Electrical signaling requires ion channel function to be tightly regulated. In principle, the permeability of a cell membrane and the resulting ionic currents can be tuned by either modulating the conductance of an individual ion channel, its selectivity, or its open probability. Most channels stochastically switch between a non-conductive and a conducting state with, in most cases, a single ion conductance level. In essence, it has been established that ion channels can undergo conformational changes to open or close their ion conduction pathway in an all-or-nothing manner, a process commonly referred to as channel gating. Being a stochastic process, a channel's gating state can be quantified by its open probability. Thus, the open probability is the key determinant of ion channel activity, and various gating mechanism have been developed during evolution to modulate ion channel function to match cellular needs.

Ion channels can be broadly classified into three types according to their activation mechanism and the stimuli that increase the open probability [1].

The first class includes ion channels that are regulated by physical stimuli such as voltage, temperature, or mechanical forces. These ion channels are often equipped with defined protein domains that serve as a sensor for the respective stimulus. For example, voltage-gated ion channels have a voltage sensing domain, which undergoes conformational changes upon transmembrane voltage changes, thereby opening or closing the ion channel. As we will see in the next section, voltage gating plays a key role in action potential generation.

The second class of ion channels responds to pH changes in the intra- or extracellular space. In these proton-gated channels, the open probability depends on the protonation state of titratable amino acid side chains.

The third class comprises primarily ligand-gated ion channels. These channels comprise defined structural elements, ligand binding sites, that recognize and bind specific signaling molecules such as neurotransmitters; the ligand binding energy is then used to open or close the channel pore.

3.3 The action potential

Excitable cells are capable of generating action potentials, transient changes of the cellular membrane potential in response to a depolarizing stimulus. A prerequisite for action potential generation is the abundant existence of voltage-gated sodium channels in certain regions of the cell [17]. These sodium channels exhibit unique gating properties: they are closed at negative potentials and open upon cell depolarization with steep voltage dependence. After reaching the so-called threshold potential, few voltage-gated sodium channels open. Mammalian cells exhibit small intracellular [Na⁺] so that opening of sodium channel results in passive influx of Na⁺ into the cell and in further depolarization of the cell. These particular gating properties of the voltage-dependent sodium channels result in reinforcing sodium influx (Fig. 6). Changes of the membrane potential to more positive values cause opening of additional sodium channels making the action potential a self-energizing stereotyped process that will result—once initiated—always in similar electrical signals.

Sodium currents are terminated by inactivation, an additional conformational change of these sodium channels, by which sodium channels are closed and then unable to re-open. Only the return of the membrane potential to negative values permit transitions from the inactivated to the resting closed state from which the channels can open again upon depolarization of the threshold potential. Sodium channel inactivation initiates the return of the action potential to the resting potential. Membrane repolarization is enhanced by the delayed opening of voltage-gated potassium channels that permit K^+ efflux and return to the K^+ diffusion potential.



Fig. 6: Idealized action potential. The action potential is initiated by an extracellular signal that depolarizes the cell (initiation phase). After reaching the threshold potential, Na^+ channels open, and Na^+ influx into the cell results in further depolarization of the membrane potential (action potential upstroke). Na^+ channel inactivation and opening of voltage-gated K^+ channels result in membrane repolarization.

The steep voltage dependence of sodium channels results in the generation of similar action potentials by a given type of cell independently of the stimulus. This so-called *all-or-nothing rule* of action potential generation has been known for more than 100 years. Inactivation of sodium channels prevents the generation of an action potential immediately after the end of the proceeding signal. Only after the refractory period, a new action potential can be elicited again. One can distinguish an absolute refractory period in which no action potential can be elicited and a relative refractory period with only reduced action potential amplitudes (Fig. 7).



Fig. 7: Refractory periods after the initiation of an action potential. Immediately after an action potential, no further action potential or only action potentials with strongly reduced amplitudes can be elicited by external stimuli.

4 Molecular simulations of membrane transport proteins 4.1 Computational Electrophysiology simulations

Electrophysiological experiments have been key to our understanding of many ion channels, transporters, and pumps. By directly measuring electrical currents induced by the transport of ions or substrates through individual transport proteins, patch-clamp experiments can directly probe the physical events underlying channel or transporter function and their regulation [18]. Combined with mutagenesis and biochemical approaches, these functional experiments also defined many of the structural underpinnings of membrane transport [1, 19].

To fully understand membrane transport at a structural level, a major breakthrough was achieved by the resolution of atomic structures of membrane proteins, pioneered by seminal work on potassium channels (e.g. on KcsA, Fig. 2), more recently followed by chloride channels, sodium channels and calcium channels, as well as transporters [3, 4, 20-23], by X-ray crystallography and cryo-electron microscopy. These structures have highlighted not only the variety of architectures employed to facilitate or drive ion transport, but also provided detailed insight into the interactions of ions within the pore as well as snapshots of ion channels and transporters in various functional states, thereby yielding a direct structural link to electrophysiological observations.



Fig. 8: Computational Electrophysiology setup for all-atom molecular dynamics simulations. Left, schematic of the simulation system with two bilayers, two membrane proteins (green), water (blue), anions, and cations. The two bilayers, in presence of periodic boundary conditions, define two separate aqueous compartments. In this case, a parallel alignment is used such that one membrane protein experiences an inward, and the other an outward directed electric field. Middle, atomistic simulation setup with two Glt_{Ph} glutamate transporters. Right, electrostatic potential profile along the membrane resulting from an applied ionic charge imbalance between compartments A and B, calculated using Poisson's equation. The potential difference between A and B is the transmembrane voltage.

However, the core of channel, transporter, or pump function, the actual transport process across the membrane, is an inherently dynamic process and thus difficult to track with static structural studies such as X-ray crystallography. Molecular dynamics (MD) computer simulations have therefore been successfully utilized to study mechanisms of ion channels, transporters, and pumps [24, 25]. Only the combination of structural biology, functional experiments, and molecular simulations can provide a detailed understanding of the structure–dynamics–function relationship of transporter function, state-of-the-art molecular simulation in channels or partial reactions of transporter function, state-of-the-art molecular simulation techniques can provide quantitative predictions of experimentally accessible properties (e.g. single-channel conductance or selectivity of ion channels). Thereby, the combination of experiment and simulation can provide validated insights into the detailed motions of individual atoms and ions underlying membrane transport with high spatial and temporal resolutions, which would have been impossible with experimental approaches alone [26-28].

Computational Electrophysiology is a recently established simulation technique to investigate membrane proteins in presence of a transmembrane proteins by all-atom molecular dynamics

(MD) simulations [29-31]. This method permits to control both ionic concentration gradients as well as voltage across the membrane. Thus, by applying a sustained electrochemical potential gradient, MD simulations can be used to directly simulate ion or substrate transport under near-experimental conditions.

MD simulations are commonly performed under periodic boundary conditions. In Computational Electrophysiology simulations, however, a separation of two compartments in a periodic simulation system is obtained by constructing a two-bilayer simulation system (Fig. 8). Such an arrangement of two membranes in a periodic system results in an inner and outer compartment, of which the latter is connected across the periodic boundaries. An ion/water exchange protocol is then used to apply small sustained charge imbalances Δq between these compartments; consequently, a voltage across both membranes is generated, similarly to the creation of transmembrane voltages in biological membranes. This transmembrane voltage is due to the capacitance of the lipid bilayer, and the transmembrane voltage is related to the charge imbalance through the capacitor equation. Since the membrane capacitance is relatively small, an imbalance of a few ions is sufficient to evoke a physiological voltage across the membrane in an atomistic MD simulation system.

4.2 Identification of the Cl⁻ channel mechanism of glutamate transporters

As explained above, EAAT glutamate transporters are secondary active transporters, that is glutamate uptake is coupled to the co-transport of three Na⁺ and one H⁺, in exchange for one K⁺ ion. Interestingly, these transporters also operate as anion-selective channels [32]. However, while providing important insights into the elevator transport mechanism underlying secondary active glutamate transport, none of the available crystal structures provided any information on the location of the anion channel and the anion-conducting conformation of the transporter.

In electrophysiological experiments on glutamate transporters, two distinct current components can be distinguished: thermodynamically uncoupled anion (e.g. Cl⁻) fluxes and electrogenic glutamate transport [33]. These Cl⁻ channels are activated upon application of transport substrates such as Na⁺ and glutamate. The opening and closing of these Cl⁻ channels is assumed to be mechanistically linked to transitions within the glutamate transport cycle [34]. In classical electrophysiological studies, several side chain mutations were identified that affect anion permeation properties such as unitary conductance or relative anion selectivities, and it was hypothesized that the anion pore was dynamically formed during the glutamate transport cycle [35]. Since mutations affect both channel and glutamate transport activities at the same time, experimental structure–function investigations alone turned out to be insufficient to identify how glutamate transporters operate as Cl⁻ channels.

Using Computational Electrophysiology MD simulations (Fig. 8), we have been able to resolve the mystery of the Cl⁻ channel mechanism in glutamate transporters [14]. At voltages from \pm 500 mV to \pm 1.6 V the available outward- and inward-facing Glt_{Ph} X-ray structures were non-conductive to ions on timescales of several microseconds (Fig. 9A). However, for certain intermediate conformations along the transition path from the outward- to the inward-facing Glt_{Ph} conformation, we observed a fully reversible transition to an anion-conducting Cl⁻ channel conformation. Within hundreds of nanoseconds, lateral movement of the so-called transport domain created an anion-selective and water-filled pore (Fig. 9B). Subsequently, the onset of anion permeation caused by the transmembrane voltage defined the anion permeation pathway.

These simulations reproduced the experimentally determined anion selectivity of these channels [32] and yielded Cl⁻ conductances and ion selectivities consistent with experimental

data (Fig. 9C). The model for the anion channel shows unique features distinct from other ion channels: the anion pore has a large diameter of \sim 5 Å and anions permeate in a partially hydrated state (Fig. 9D). Furthermore, the positive charge of a single arginine side chain in the pore center confers anion selectivity to these channels. To experimentally confirm the simulated Cl⁻ permeation pathway, *in silico* screening of mutations of pore-lining residues was performed using Computational Electrophysiology simulations. Several substitutions were identified that either increase or decrease anion currents. Interestingly, some mutants even converted the anion channel into a non-selective anion/cation channel. Patch-clamp experiments on corresponding mutations inserted into glutamate transporters revealed similar effects, thereby validating the simulation results [14].



Fig. 9: Chloride channel mechanism of the secondary active glutamate transporter Glt_{Ph} . (A) The outward–inward elevator transition of the glutamate/aspartate transport cycle (cf. Fig. 4). The static trimerization domain is shown in blue cartoon representation, the mobile transport domain in yellow. (B) Illustration of the Glt_{Ph} chloride channel-forming conformation. An anion pore along the interdomain interface is created via lateral movement of the transport domain in intermediate transporter conformations. Red spheres represent a single permeating Cl. (C) Cumulative permeation count from Computational Electrophysiology MD simulations of Glt_{Ph} at +800 mV or -900 mV in the presence of 1 M NaCl or NaI. (D) Pore profile of anion hydration and pore diameter. Hydration numbers are integrals of Cl'/hydrogen radial distribution functions to the first minimum.

The identification of the Cl⁻ conduction mechanism of secondary active glutamate transporters would not have been possible without the use of molecular simulations. With Computational Electrophysiology simulations, Cl⁻ permeation through these channels was directly simulated and key ion conduction parameters, which were accessible to experiments (such as unitary conductance or anion/cation selectivity) could be readily determined from the simulations. This gave rise to further experimental validation and informed a number of additional experiments, which together were able to draw a comprehensive picture of the molecular interplay of secondary active and ion channel-mediated transport by a single transport protein [14].

References

- [1] B. Hille, *Ion channels of excitable membranes*, 3rd ed. (Sinauer Associates, Sunderland, MA, 2001)
- [2] O. Jardetzky, Nature 211, 969 (1966)
- [3] D.A. Doyle, et al., Science 280, 69 (1998)
- [4] D. Drew and O. Boudker, Annu. Rev. Biochem. 85, 543 (2016)
- [5] N.C. Danbolt, Prog. Neurobiol. 65, 1 (2001)
- [6] D. Yernool, et al., Nature 431, 811 (2004)
- [7] O. Boudker, et al., Nature 445, 387 (2007)
- [8] N. Reyes, C. Ginter, and O. Boudker, Nature 462, 880 (2009)
- [9] G. Verdon, et al., Elife 3, e02283 (2014)
- [10] S. Jensen, et al., Nat. Struct. Mol. Biol. 20, 1224 (2013)
- [11] A. Guskov, et al., Nat. Commun. 7, 13420 (2016)
- [12] J.C. Canul-Tec, et al., Nature 544, 446 (2017)
- [13] S. Gendreau, et al., J. Biol. Chem. 279, 39505 (2004)
- [14] J.P. Machtens, et al., Cell 160, 542 (2015)
- [15] D. Ewers, et al., Proc. Natl. Acad. Sci. U. S. A. 110, 12486 (2013)
- [16] A.L. Hodgkin and P. Horowicz, J. Physiol. 148, 127 (1959)
- [17] A.L. Hodgkin and A.F. Huxley, J. Physiol. 117, 500 (1952)
- [18] B. Sakmann and E. Neher, *Single-Channel Recording*, (Springer, Boston, MA, 2009)
- [19] P. Läuger, *Electrogenic ion pumps*. (Sinauer Associates, Sunderland, MA, 1991)
- [20] Y. Zhou, et al., Nature 414, 43 (2001)
- [21] R. Dutzler, et al., Nature 415, 287 (2002)
- [22] Z. Yan, et al., Cell 170, 470 (2017)
- [23] E. Park, E.B. Campbell, and R. MacKinnon, Nature 541, 500 (2017)
- [24] C. Maffeo, et al., Chem. Rev. 112, 6250 (2012)
- [25] G. Enkavi, et al., Methods Mol. Biol. 924, 361 (2013)
- [26] M.O. Jensen, et al., Science 336, 229 (2012)
- [27] D.A. Köpfer, et al., Science 346, 352 (2014)
- [28] S.Y. Noskov, S. Bernèche, and B. Roux, Nature 431, 830 (2004)
- [29] C. Kutzner, et al., Biochim. Biophys. Acta 1858, 1741 (2016)
- [30] C. Kutzner, et al., Biophys. J. 101, 809 (2011)
- [31] J.P. Machtens, et al., Biophys. J. 112, 1396 (2017)
- [32] Ch. Fahlke, D. Kortzak, and J.P. Machtens, Pflügers Arch. 468, 491 (2016)
- [33] J.I. Wadiche and M.P. Kavanaugh, J. Neurosci. 18, 7650 (1998)
- [34] J.P. Machtens, P. Kovermann, and Ch. Fahlke, J. Biol. Chem. 286, 23780 (2011)
- [35] R.M. Ryan, A.D. Mitrovic, and R.J. Vandenberg, J. Biol. Chem. 279, 20742 (2004)

B 5 DNA and Chromatin

Helmut Schiessel Institute Lorentz for Theoretical Physics, Leiden University Niels Bohrweg 2, 2333 CA Leiden, The Netherlands

Contents

1	Introduction				
	1.1	The hierarchical structure of chromatin	2		
	1.2	"A genomic code for nucleosome positioning"	3		
	1.3	The space of all nucleosomal sequences	4		
	1.4	A coarse-grained nucleosome model	5		
	1.5	The Mutation Monte Carlo (MMC) method	7		
2	The mechanical genome				
	2.1	DNA mechanics dictates nucleosome positioning rules	8		
	2.2	Genetic and mechanical information can be multiplexed	9		
	2.3	Evidence for a mechanical evolution of DNA molecules	10		
3	Con	clusions	13		

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

DNA molecules are the carriers of genetic information for all life forms. DNA is typically found as a right-handed double helix with its two strands running antiparallel and its bases A, T, G, C forming pairs, A with T and G with C. Each stretch of DNA that encodes for a protein is called a gene. A protein is build from 20 different building blocks, called amino acids. With its much smaller alphabet DNA encodes for amino acids by grouping sets of three consecutive bases into information units, the codons. The precise rules how codons encode for amino acids is called the genetic code. As there are $4^3 = 64$ codons but only 20 amino acids the genetic code is degenerate. In other words, there are multiple ways to encode for one and the same amino acid (in 18 of the 20 cases). This degeneracy will be crucial in the following.

The point that these Lecture Notes want to make is that in addition to the genetic information (the genes encoding for the proteins) there is a second layer of information which is mechanical in nature. This is possible because the mechanical and geometrical properties of the DNA double helix depend on the underlying sequence of base-pairs. By choosing the "right" sequence of base-pairs a stretch of DNA can be made softer than average or stiffer than average. It is also possible to choose sequences that make the DNA molecule bent in a certain direction. The claim I want to make is that organisms have evolved their genomes to put mechanical cues along DNA molecules. Especially exciting is the fact – demonstrated below – that the classical genetic and the mechanical information can be multiplexed freely, allowing to put mechanical cues on top of genes at will, and not just on top of stretches of "junk" DNA. (We know multiplexing from daily life technologies, e.g. having two phone conversations on the same wire.)

What could be the meaning of such mechanical cues? I will argue that the cues guide the packaging of DNA molecules inside cells and by this indirectly the access to its genes. What I need to describe next is what we know about the packaging of DNA inside cells.

1.1 The hierarchical structure of chromatin

We focus here on eukaryotes (which include animals, plant and fungi). Cells of eukaryotes keep their DNA in a separate compartment, the nucleus. Eukaryotic DNA is packaged with the help of proteins into a DNA-protein complex called chromatin. Each individual DNA molecule together with the complexed proteins is called a chromosome (human somatic cells have 46 chromosomes). The structure of chromatin is hierarchical, see Fig. 1 adapted from Ref. [1], but many details of the different levels are not well understood. The first level of compaction is the wrapping of DNA molecules around protein cylinders, leading to DNA spools called nucleosomes. These are the main players of these Lecture Notes.

Just for completeness let me provide a short discussion of the higher levels. Traditionally the next level is believed to be the chromatin fiber, a rather compact structure into which the string of nucleosomes is folded. Fibers are easily seen in the test tube and lots of energy has been spent in figuring out their precise microscopic structure. But just as the debate between various competing detailed models of chromatin fibers raged at its fullest, some new experiments put serious doubts on the generally accepted believe that chromatin fibers exist in living cells [2]. Even though they are readily seen *in vitro* [3], also very recent work does not find them *in vivo* [4]. That is why I put a big question mark on top of my picture of the chromatin fiber in Fig. 1. Also the structures beyond that level are not well understood. But it is worthwhile to mention that there is currently tremendous progress in the understanding of the larger scales thanks to a new experimental method called chromosome conformation capture [5]. This leads



Fig. 1: The hierarchical structure of a chromosome: the DNA double helix is wrapped around protein cylinders to form nucleosomes, the string of nucleosomes packs into a 30 nm wide chromatin fiber that folds into the chromosome. Here the well-known X-shaped mitotic chromosome is shown with its two identical copies of the DNA molecule which forms before cell division. Details (and not only details!) of the structures beyond the nucleosome are still a matter of debate.

to new exciting developments, e.g. the idea of the loop extrusion mechanism [6]. This would have been an interesting and timely subject to speak about in a lecture. But the goal here is to discuss mechanical cues in DNA molecules. And such mechanical cues are most important at those places where DNA is bent most. This happens on the smallest compaction scale, the nucleosome, to which we now turn.

1.2 "A genomic code for nucleosome positioning"

The nucleosome is an ideal reader of mechanical information. The reason for this is twofold and can be best seen by inspecting the crystal structure of the nucleosome core particle [7]. The nucleosome core particle is the complex formed by DNA of exactly the nucleosomal wrapping length, 147 base-pairs, and the core of histone proteins, see Fig. 2(a). In the cell millions of such complexes are connected by a given DNA molecule into a string of nucleosomes connected by non-complexed stretches of so-called linker DNA, about 0 to 80 base-pairs in length. In each nucleosome 147 base-pairs, corresponding to about one DNA persistence length, are wrapped in one and three quarter turns around an octamer of histone proteins. This means that the energy of bending DNA into a nucleosome is large, about $60 k_{\rm B} T$ [1]. Even small difference in the sequence will lead to large differences in the bending energy between those sequences. This is one reason why nucleosomes are ideal readers of mechanical information. The second reason is related to the way the DNA is bound to the histone octamer. It is bound at 14 locations where the two backbones touch the surface of the histone octamer (at the so-called minor groove of the double helix). As the sugar-phosphate backbones are independent of the underlying sequence, the pure binding energy is only weakly sequence dependent. Taken together, these two features make the nucleosome the master of the so-called indirect readout. Whereas most DNA binding proteins find their target by reading the sequence directly, the affinity of a 147 base-pair long stretch to be complexed in a nucleosome reflects the ease with which it is wrapped into it.



Fig. 2: (a) Crystal structure of the nucleosome core particle (top and side view): 147 basepairs are wrapped around an octamer of histone proteins [7]. (b) Widom's "genomic code for nucleosome positioning" [9]: key base-pair steps that increase the affinity of a sequence to the nucleosome are displayed relative to the structure of one-half of the nucleosome.

It is indeed known since a long time that the affinity to nucleosomes varies with sequence. To learn which sequences have a strong affinity to nucleosomes, Travers and coworkers extracted in 1986 chromatin from chicken [8]. They added an enzyme, DNAase, that digested all the freely available DNA leaving just DNA intact that was very stably wrapped into nucleosomes. 177 of those intact sequences, about 150 base-pairs long, were sequenced. When looking at individual sequences it was not obvious why they had a higher affinity than average. But by looking at statistical properties of these sequences, i.e. averages over those sequences, certain pattern arose. It turned out that certain base-pairs steps are more likely to be found at certain positions along the nucleosome and less likely at other positions (one looks at base-pair steps because these are the objects that are deformed when wrapping the DNA around the octamer, as explained in detail further below). Looking along one of the strands in its so called 5' to 3' direction, the base-pair steps of importance are GC (a G followed by a C), AA, TT and TA. These four steps show characteristic oscillations in their occurrence frequency along the more stable nucleosomes, with a period of 10 base pairs, the DNA helical repeat, see Fig. 2(b). Specifically GC steps are found where the DNA's minor groove faces outward, whereas AA, TT and TA are all in phase with each other and peak at the positions in between, namely where the minor groove faces the histone octamer.

This suggests the possibility that genomic sequences have evolved over evolutionary time scales to position nucleosomes at certain positions and to encode for their stability or other physical properties. The late Jonathan Widom suggested that there is a "genomic code for nucleosome positioning" [9]. He said that "genomes care where nucleosomes are on average and so genomes code explicit information to bias [their positions]." After Jon's untimely death in 2011 I decided to follow up on his ideas using methods from statistical physics.

1.3 The space of all nucleosomal sequences

As a starting point to think about this complex problem, it is useful to introduce the concept of the sequence space of all DNA stretches that can be wrapped into a nucleosome. How many distinct sequences exist? There are four different bases and the length of the wrapped DNA is 147 base-pairs. This leads to 4^{147} different sequences, a huge number on the order 10^{88} (strictly speaking, due to a rotation symmetry of the nucleosome around the so-called dyad, see Fig. 2(b), this number needs to be divided by two). If you wanted to synthesise all this DNA,

you would need a big lab as all this DNA would fill five Milky Ways densely.

So how do scientists study this gigantic space? Well, strictly speaking they don't. This is because no matter if they are biologists, bioinformaticians or physicists, no matter whether they are experimentalists or theorists, what they tend to look at are typically genomes of certain organisms. Very popular is baker's yeast whose genome is about 12 million base-pairs long. This means that when one studies nucleosome positioning on that genome, one accesses only the 10^{-80} 's fraction of the nucleosomal sequence space. Even highly complex organisms like humans with their 3.2 billion base-pairs long genome can only scratch the surface of sequence space.

If we wanted, for instance, to answer the question: "which is the sequence with the highest affinity to reside in a nucleosome?" we would have no chance to find it by scanning the whole human genome (assuming that our genome has not evolved for highest nucleosome affinity which is a known fact). A slightly better approach is to start from a huge pool of random DNA molecules and then fish out of this pool the sequences with the highest affinity. This has been done in 1998 in the Widom lab. Starting from a huge pool of 5 trillion random DNA molecules (slightly longer than the nucleosomal wrapping length) the molecules were mixed with a much smaller number of histone proteins (namely one octamer per 10 DNA molecules) [10]. After the complexes had formed, the non-complexed DNA molecules were discarded and the winners were multiplied. This process was repeated 15 times. At the end of this so-called SELEX experiment there were only a few dozen types of DNA molecules left, all with a much higher affinity than average. The best of those sequences was called 601 (I do not know why, I wish I did) and it is nowadays the most common sequence used in the lab when working with nucleosomes.

You might wonder why one is left with just a handful sequences after starting with 5 trillion sequences. Why so few? The reason is that this experiment started from a random pool of sequences. If we order our sequence space such that the highest affinity sequences are in the center of the "Milky Way" and the lower affinity sequences toward the outskirts, then the starting sequences will fill randomly the whole space. There will then be only a small fraction of sequences close to the center and those few sequences are the only ones who have a chance to win the competition for the histone octamers. In short, one has a lot of waste DNA that needs to be discarded.

In these Lecture Notes I will introduce a different type of approach where practically all sequences to be "produced" will be automatically high affinity sequences. It is a computational approach that we call Mutation Monte Carlo (MMC) method [11]. It can be used in principle for any computer model of the nucleosome as long as it accounts for DNA sequence effects. Before I explain how MMC works I will introduce the model that we have tested and used [11].

1.4 A coarse-grained nucleosome model

Our nucleosome model is depicted in Fig. 3(a). It consists of a coarse-grained representation of the DNA molecule, 147 base-pairs long. The DNA molecules is forced into the configuration in which it is found in the nucleosome crystal structure, using 28 constraints that mimic the 14 binding sites in a real nucleosome. The protein core is not modelled explicitly, its presence is only accounted for by the constraints on the DNA.

The DNA double helix is represented by the so-called rigid base-pair model [12]. This model accounts only for the base-pairs of the DNA molecule that are modelled as rigid blocks. This leaves six degrees of freedom between neighbouring base-pairs, called shift, slide, rise, tilt,



Fig. 3: (a) The coarse-grained nucleosome model from Ref. [11]. Same colour scheme as in Fig. 2(b). (b) Base-pair step probability distributions (average over 10 million sequences) produced by MMC at one third of room temperature using the model from (a) [11]. The model reproduces the standard nucleosome sequence preferences, Fig. 2(b).

roll, twist, see middle of Fig. 4. If you want to create an ordinary straight piece of DNA double helix, all you need are two degrees of freedom, rise and tilt. A rise of 0.34 nm combined with a twist of about 36 degrees leads to a twisted stack of base-pairs that looks similar to the real DNA double helix (in its common B-form), see Fig. 4 left. Note that because the two sugarphosphate backbones are attached to one long site of the base-pairs, one has two grooves, a major and a minor one, going around the DNA double helix. When one looks at space-filling figures of the DNA double helix, one can clearly identify these two types of grooves. Also note that one can see clearly in such figures the parallel twisted stack of base-pairs through the gap created by the major groove. It is through this gap that DNA binding proteins typically "read" the DNA sequence by reaching inside that groove. This is how direct readout works. As mentioned above, this is *not* how the sequence preferences of nucleosomes come about.

In order to bend the DNA around the nucleosome, other degrees of freedom have to be invoked. We will mention here only the most important one, roll. Roll is the rotation around the long axis of the base-pair step. It is convention to call the roll positive if the base-pair stack is compressed towards the major groove. By periodically changing the roll from positive to negative and back with the DNA helical repeat (about 10 base-pairs) the twisted stack of base-pairs bends in one direction, see Fig. 4 right. This allows to rephrase the nucleosome positioning rules: high affinity sequences feature GC steps at positive roll positions, and AA, TT and TA steps at negative roll positions. Where do these sequence preferences come from?

In order to make any prediction one needs to go beyond a purely geometrical model by introducing also energy into the system. In fact, the rigid base-pair model has been fully parametrized in the literature. One assumes only nearest-neighbor interactions with a quadratic deformation energy between successive base-pairs [12]:

$$E = \frac{1}{2}(q - q_0) \cdot K \cdot (q - q_0).$$
(1)

Here q is a six-component vector that describes the relative degrees of freedom between two base-pairs. The intrinsic, preferred values of these degrees of freedom are given by q_0 . The properties of the (six-dimensional) springs connecting the base-pairs are given by K, a six-by-six stiffness matrix. The sequence-dependence of the model comes into play because the stiffness (K) and intrinsic shape (q_0) of a given base-pair step depend on its chemical identity. In other words K and q_0 are different for different types of base-pair steps.

These parameters have been determined in the literature, either by looking at the conformations and fluctuations of DNA-protein cocrystal structures [12] or by performing all-atom molecular



Fig. 4: The rigid base-pair model is a coarse-grained DNA representation that leaves six degrees of freedom per base-pair step (middle). A base-pair step with 0.34 nm rise and about 36 degrees twist produces a straight standard DNA double helix (left). When in addition the roll is changed periodically with the DNA's helical repeat one obtains a bent stack (right).

dynamics simulations of short DNA molecules with various sequences [13]. We are learning currently which of the parameter sets works best, often we use a hybrid version that uses both sources [14].

In summary, our nucleosome model consists of 147 base-pairs of DNA represented by the rigid base-pair model wrapped into a superhelix that mimics its configuration in the nucleosome crystal structure. This is achieved via 28 rigid constraints, two per binding site. One last additional detail: each rigid constraint consists of a fixed mid-plane for a consecutive base-pair step (corresponding to the bound phosphate of the involved backbone).

1.5 The Mutation Monte Carlo (MMC) method

Having now a nucleosome model with sequence dependent energetics we can introduce the MMC method. This method allows to scan regions in the nucleosomal sequence space that are special with respect to their elastic properties. But first let us discuss how a nucleosome with a given fixed sequence can be studied using a standard Monte Carlo scheme. What we would like to achieve is to sample the configurational space of the nucleosome according to the Boltzmann distribution. This is achieved as follows. Pick a random base pair. Perform a small rotation around a random axis together with a small translational shift in a random direction. According to Eq. 1 this changes the mechanical energies of the two involved base-pair steps by a (small) amount ΔE . If $\Delta E < 0$ accept the move. If $\Delta E > 0$ accept it only with a probability $e^{-\beta \Delta E/k_{\rm B}T}$. Continuing this process one obtains an equilibrium distribution of the nucleosomal DNA configurations from which one can determine e.g. the average elastic energy. (One technicality: whenever the chosen base-pair happens to be next to a rigid constraint, move the base-pair across the fixed mid-plane symmetrically to keep it fixed).

So far we are stuck at one point in sequence space. How can we explore that space? The trick is

extremely simple and very effective. The MMC method developed by Eslami-Mossallam [11] uses in addition to the conformational moves also mutation moves. A mutation move consists of randomly picking a base-pair and attempting to change its chemical identity. This affects again, as for the conformational moves mentioned above, the mechanical energy, Eq. 1, of the two involved base-pair steps. But instead of changing q, this changes K and q_0 of the corresponding steps. The move is accepted or rejected according to the energy change using the same rules as above.

By randomly mixing conformational and mutational moves, the system moves through sequence space and quickly arrives at nucleosome sequences (and corresponding conformations) that are much cheaper than average. One can then easily create 10 million independent high affinity sequences, rather than just a few as it is the case in experiments [10].

What is the role of temperature in such a simulation? MMC produces a set of conformations and sequences distributed according to the Boltzmann distribution at the chosen temperature. The lower the temperature the smaller is the section in sequence space that is explored focusing on sequences with higher and higher affinities. The temperature can thus be seen as a tool that allows to adjust the volume in sequence space that will be probed. By cooling the system close to zero temperature, it is even possible to identify the ground state sequence of our model nucleosome.

2 The mechanical genome

We can now ask the question: Is there – in addition to the classical genome (the genes that encode for the proteins) – a "mechanical genome", i.e. a set of mechanical cues written along DNA molecules that have formed over evolutionary time scales in parallel and independent of the classical genome? To answer this question we need at least to show three things: that the nucleosome positioning rules are mechanical in nature, that the mechanical cues can be multiplexed with the classical genetic information and that such mechanical cues do actually exist on real genomes. The next three sections demonstrate these three fundamental aspects of the mechanical genomic code using the tools introduced above.

2.1 DNA mechanics dictates nucleosome positioning rules

We have described earlier the nucleosome positioning rules, see also Fig.2(b). Assuming that the rules are caused by DNA mechanics and that the nucleosome model we introduced above, Fig. 3(a), is realistic enough to make reasonable predictions (as various tests suggest [11, 15, 16]), we need to show that sequences which follow these rules more than average sequences have indeed a higher affinity than average. The MMC approach allows to answer this question in a straightforward way. All what needs to be done is to run such a simulation on the model nucleosome and then to check whether the produced sequences follow on average the rules. We produced 10 million independent high-affinity sequences by performing a MMC simulation at 1/3 of room temperature. We then looked at base-pair step distributions obtained by averaging over all those sequences. Figure 3(b) displays the distribution for GC steps and the combined distribution for AA, TT and TA steps. This procedure indeed recovers the standard nucleosome positioning rules, Fig. 2(b). This suggests that these well-known rules are caused to a large extent by the sequence dependent elasticity and geometry of the DNA double helix.

the model predicts that GC steps peak at positive roll position. When one inspects the underlying geometrical preference of GC steps and compares it to all other steps, one learns that it is the step with the lowest value of roll. This together with the fact that it is one of the stiffest steps shows that GC is the step most "unhappy" to occupy large roll positions. So why does it peak at these positions against its own preferences? The reason is that each base-pair step is part of a larger sequence. When a GC step occupies a given position, the previous step has to end on a G and the following step starts with a C. As it happens, these neighbouring steps feature on average a low elastic energy if GC sits at large positive roll position. So it is the neighbours of GC but not GC itself that cause the peak of GC at high roll positions.

2.2 Genetic and mechanical information can be multiplexed

The next question to be considered is whether mechanical cues can be written freely on top of genes. To demonstrate this we start by looking at some randomly picked gene from a standard model organism, baker's yeast. Figure 5(a) shows a 500 base-pairs long stretch of gene YAL002W. It depicts the energy landscape that the nucleosome experiences as it is moved along that stretch of DNA. This has been calculated by a simple Monte Carlo simulation (without mutations). Note the strong undulations with a period of about 10 base-pairs. These are caused by the fact that - as one moves the DNA molecule through the nucleosome - it has to perform a corkscrew motion such that the DNA minor groove is always in contact with the binding sites. Typically a given DNA molecule has locally a preferred bending direction, just by accident. So about every 10 base-pairs along the sequence there is typically a minimum in the energy landscape, five base-pairs further a maximum and so on. This leads to the so-called rotational positioning of nucleosomes. This positioning might be important as it guides the higher order arrangement of nucleosomes. Another type of positioning, translational positioning, will be discussed in the next section. Note the vertical lines in this plot; these correspond to nucleosomes that have been mapped via a chemical method in yeast *in vivo* [17]. Most of the mapped nucleosomes fall in minima of our energy landscape, all along the yeast genome. This shows again that this model predicts properly the nucleosome positioning rules.

In the following we demonstrate that it is possible to change this rotational positioning at will without affecting the protein that the gene encodes for. How is this possible? As mentioned in the introduction 64 codons encode for only 20 amino acids. This degeneracy of the genetic code can be employed to change the mechanical properties of the DNA molecule keeping the encoded protein unchanged. Figure 5(b) displays a short stretch of the YAL002W gene (top row). The sequence is already broken into codons. Below that sequence of codons you find in red the sequence of amino acids that the gene encodes for. Below each amino acid there is a list of all the synonymous codons that represent this specific amino acid. In 18 of 20 cases there is in fact more than one codon available.

We make use of this degeneracy of the genetic code in a modified version of the MMC method. We now allow only synonymous mutations, i.e. we swap between codons that form a synonymous set. This way we can change the mechanical properties of the DNA molecule without affecting the sequence of amino acids that the base-pair sequence encodes for. More specifically, we focus on the well-positioned nucleosome on base-pair position 826 in Fig. 5(a). This nucleosome has also been mapped *in vivo* at precisely that position [17]. We would like to demonstrate that one can shift this local energy minimum to any position one likes, e.g. base-pair position 827 and so on, by synonymous mutations.

To achieve this we perform synonymous MMC simulations for the nucleosome placed on the



Fig. 5: (a) Energy landscape for a nucleosome on a stretch of gene YAL002W from baker's yeast calculated using the model from Fig. 3(a). The vertical lines correspond to nucleosomes mapped in vivo [17]. (b) Top row: a stretch of the same gene as in (a) broken into a sequence of codons. Middle row (red): the sequence of encoded amino acids. Bottom (purple): list of synonymous codons for each given amino acid. (c) Synonymous energy landscapes below the original landscape (inside magenta box). This plot shows that a local minimum can be placed anywhere on that stretch of gene from base-pair 826 to 831 [11].

positions where we want to create new local minima. Figure 5(c) demonstrates that this method works by displaying energy landscapes obtained by this method for the positions 827 to 831 [11]. Using a more sophisticated method we have in the meantime extended this calculation genome wide and were able to show that nucleosomes can be rotationally positioned anywhere on the yeast genome in at least 99.95% of the cases. This is ongoing work and we are not yet sure whether the remaining 0.05% correspond really to locations where one cannot position a nucleosome at all or whether we have to improve our method further.

2.3 Evidence for a mechanical evolution of DNA molecules

So far we have shown that in principle mechanical cues could have been written into DNA molecules and that even on top of genes. But the question remains whether this has really happened on actual genomes. And if the answer is yes: what are then the biological functions of such cues? These questions are not straightforward to answer. For instance, when you look again at Fig. 5(a) you can see a wildly oscillating energy landscape that a nucleosome would experience as one pushes it along the DNA molecule. But does this landscape constitute some kind of meaningful signal? As a matter of fact, the landscape of a completely random base-pair sequence looks pretty much the same.

To isolate meaningful signals out of genomes it turns out to be crucial to look at genome wide averages. Only then one beats the (possibly random) oscillations and starts to find in fact some very strong mechanical cues along DNA molecules. However, our nucleosome model is too slow to calculate genome-wide energy landscapes as one would need to perform a Monte Carlo simulation at each position along the genome, in order to allow the wrapped DNA stretch to
sample equilibrium configurations. Surprisingly the MMC method comes at our rescue also for this seemingly unrelated problem. The idea is to perform one long MMC simulation to learn about the sequence preferences of the nucleosome and then to use these sequence preferences as input in a simplified probabilistic model [18].

We have already determined these sequence preferences in the form of base-pair step probabilities along nucleosomes, see Fig. 3(b) for some examples. Specifically, we can use the MMC approach to learn about the probability to have base S_i at position *i* along the nucleosome with i = 1, ..., 147 and $S_i = A, T, G, C$ and the joint probability $P(S_i \cap S_{i-1})$ to have base S_i following base S_{i-1} . From these probabilities we can calculate conditional probabilities $P(S_i|S_{i-1}) = P(S_i \cap S_{i-1})/P(S_{i-1})$. We then estimate that the probability of the 147 base-pair long sequence S to be occupied by a nucleosome is given by

$$P(S) = P(S_1)P(S_2|S_1)\prod_{i=3}^{147} P(S_i|S_{i-1}).$$
(2)

This equation assumes that there are no longer-ranged effects along the DNA molecule, i.e. the probability of a given base to appear in a nucleosome only depends on the previous base but not on the precise nature of bases further away. That this is a reasonable approximation can be rigorously tested by performing an improved analysis starting from the probability distributions for triplets of bases which only slightly improve the predictions [18].

Which predictions do we actually speak about? One can estimate the energy of a sequence from the probability by taking the logarithm: $E(S) = -k_{\rm B}T \ln P(S)$. This way one can calculate the energy landscape of e.g. the YAL002W gene from above and compare it to the actual energy landscape as calculated from the full model. The deviations between the "real" energy landscape and the probabilistic one (based on duplets or triplets of bases) is on the order of one $k_{\rm B}T$, much smaller than the typical energy undulations in the energy landscape, see Fig. 5(a). So we make only a small error but what do we gain from it? It turns out that the speed-up using the probabilistic model is of the order of 10^5 . This means that we can now perform genome wide calculations.

What needs to be done to obtain clean signals is somehow to average over the genome. A wellknown way to do this is to align the same type of functional sites from all over the genome. The most promising candidate to look at is the beginning of genes as these are the places where a cell decides whether a gene is read out or not. Fig. 6(a) shows gene start sites of baker's yeast averaged over all its genes (about 6000). More specifically we show a 2000 base-pair long interval with the genes starting in the middle and going toward the right. The quantity depicted is the so-called nucleosome occupancy which is the probability that a given basepair is covered by a nucleosome. This quantity is chosen as it is experimentally accessible. We assume that there is one nucleosome and calculate its occupancy for this 2000 base-pair wide window. From our calculation (based on Eq. 2) we produce the blue curve [16] which fits astonishingly well with the experimentally determined occupancy (green curve). In the experiment [19] nucleosomes are reconstituted on yeast DNA and their positions are determined by digesting the DNA with DNAase. The excellent agreement between data and model is certainly only a fortunate coincidence; what is important here is that both approaches give qualitatively the same overall signal.

It can be clearly seen that there is a strong depletion of nucleosomes just in front of the genes. This depletion signal in yeast has been speculated to be "partially encoded in the genome's intrinsic nucleosome organisation, and that this intrinsic organisation may facilitate transcription



Fig. 6: Nucleosome occupancies around the beginning of genes in various organisms. All plots are averages over all genes, either aligned at the gene start sites, (a), or the transcription start sites, (b) and (c). (a) Baker's yeast (in vitro data [19] and prediction), (b) humans (in vitro data [20], retained nucleosomes in sperm cells [23] and prediction) and (c) prediction for various unicellular and multicellular organisms. All predictions are based on Eq. 2 [16].

initiation and assist in directing transcription factors to their appropriate sites in the genome" [19]. In short, DNA is stiffer than average before genes to keep its DNA free of nucleosomes so that the transcription machinery can always access that region if it wants to produce the corresponding protein.

This is what has been found for yeast. What about other organisms? Figure 6(b) shows the nucleosome occupancy averaged over all genes and aligned at the transcription start site (which is close to the gene start site) for the human genome. Surprisingly our model (blue curve in Figure 6(b)) shows a completely different signal, featuring a large and wide peak in the nucleosome occupancy [16]. How does this compare to experiments? At first not very favourably. The green curve in Fig. 6(b) are *in vitro* data showing a much smaller peak [20]. However, there is a profound difference between the experiment and the calculation. In the calculation we consider the probability distribution of only *one* nucleosome density. Since nucleosomes cannot sterically overlap, there is a saturation in the density around the peak. In fact, accounting in our calculations for a similar density as in the experiment, preventing steric overlap, we find a curve (dotted in Fig. 6(b)) similar to the experimental curve. But even if the nucleosome density cannot increase much around the transcription start sites, the signal is still contained in the affinity and thus stability of the corresponding nucleosomes.

What could be the biological function of these mechanical cues in the human genome? The following speculation [21] is based on the fact that humans – unlike yeast – are multicellular organisms: "[...] high nucleosome preference is directly encoded at regulatory sequences in the human genome to restrict access to regulatory information that will ultimately be utilised in only a subset of differentiated cells." So the idea is that many genes are only meant for specialised cells and that those genes should be closed off in all other types of cells. And this is achieved by encoding for stable nucleosomes around the start sites of those genes.

An exciting question to ask is whether this distinction between yeast and human is an example of a general rule in biology. This is hard to answer on the basis of experiments as there are not so many nucleosome maps available and, even if they were, it is not so easy to detect a signal because of the density saturation due to the excluded volume between nucleosomes. Using our model we looked at 50 different organisms and calculated the occupancy signal (for a single nucleosome) around all transcription start sites [16]. As you can see in Fig. 6(c), it is

indeed generally true that the DNA elasticity around transcription start sites is entirely different between unicellular organisms like baker's yeast or the green alga *Chlamydomonas reinhardtii* and multicellular lifeforms like zebrafish, mouse, human, chimpanzee and rice.

Is this the end of the story? Not quite. At least for humans the above given biological speculation turns out to be wrong. Dividing genes between house keeping genes and tissue specific genes and looking at the mechanical cues separately, one discovers the opposite of what one would have expected: the strong mechanical cues stem from the house keeping genes that all cell types need [22]. Even worse, when looking at actual *in vivo* nucleosome occupancies it was found that they do not reflect at all underlying DNA mechanics. Instead the transcription start sites of house keeping genes are typically depleted of nucleosomes [22]. This is, of course, expected, but it goes against the mechanical cues.

What is happening here? Apparently other processes, the binding of transcription factors to their specific target sites, the transcription by RNA polymerase and/or the action of chromatin remodellers (motor proteins that push and pull nucleosomes using ATP) overrule the mechanical cues around transcription start sites. The mechanical cues must therefore have a different function. The most logical explanation would be that they are of importance in a cell type that is transcriptionally not active.

Are there such cells in multicellular organisms? In fact, each animal no matter how big it is (think of an elephant!) needs eventually make itself very small when passing through the germ line into the next generation. Especially in sperm cells elephants shrink substantially (even smaller than the sperm cells of fruit flies or mice!). Small sperm cells are good swimmers and can be produced in larger numbers, a fact especially important for species where there is a competition between different males. That might be the reason why in sperm cells DNA is tightly packed with the help of protamines and all nucleosomes are evicted. But not quite: a recent finding shows that about 4% of the nucleosomes are retained in human sperm cells [23]. How does a sperm cell know which nucleosomes to keep? As we found out, it is the mechanical cues in the DNA molecules that determine which nucleosomes are retained: Regions where our model predicts the most stable nucleosomes correspond to regions where sperm cells retain nucleosomes, see Fig. 6(b) (brown curve) [16].

What is the evolutionary driving force for retaining a fraction of nucleosomes in sperm cells instead of getting rid of all of them? We can only speculate. But a likely reason is to allow for the transmission of epigenetic information via the father to the offspring (and not only by the mother where the nucleosomes are kept in the egg cell). Epigenetics is information in addition to and shorter-lived than genetic information. It is scribbled along the margins of the book of life. One way this can be achieved is by chemically modifying the histone proteins that form the octamer of the nucleosomes. This changes e.g. their stickiness affecting the accessibility to the associated DNA. And it is the genes that are important for the early embryonic development that are singled out for receiving this extra information; these carry the mechanical cues and thus retain the nucleosomes. A concrete (though controversial) experiment trained male mice using mild foot shocks to fear cherry blossom smell. Their offsprings had an aversion to this specific odour [24].

3 Conclusions

We have come a long way from the mechanics of base-pairs to the smell of cherry blossoms. The point that these Lecture Notes wanted to make is that if there are some degrees of freedom (here the DNA elasticity) that evolution can play with, it very likely makes use of it. It is, however, far from obvious what comes out of such an evolution. For example, so far the main interest in the field is to learn which nucleosomes are positioned by mechanical cues. This is done by assigning one number, the affinity of the sequence, to a given 147 base-pair long sequence. This does, however, overlook the fact that this is a much richer problem. In principle, the mechanical properties of 147 base-pairs wrapped into a nucleosome could give some nucleosomes distinct sets of physical properties, setting them far apart from standard nucleosomes. For instance, it has been shown that nucleosomes which are strongly asymmetric with respect to their two DNA halves act as polar barriers for transcribing RNA polymerases [25]. Using our model together with the MMC approach we have started to build designer nucleosomes and might be used as "force sensors" [26]. An exciting question to ask is whether and where such special nucleosomes have evolved on real genomes and to what purpose.

References

- H. Schiessel, *Biophysics for Beginners: a Journey through the Cell Nucleus* (Pan Stanford, Singapore, 2014)
- [2] K. Maeshima, S. Hihara, M. Eltsov, Curr. Op. Cell Biol. 22, 291 (2010)
- [3] P.J.J. Robinson, L. Fairall, V.A.T. Huynh, D. Rhodes, Proc. Natl. Acad. Sci. U. S. A. 103, 6506 (2006)
- [4] H.D. Ou et al., Science 357, 370 (2017)
- [5] E. Lieberman-Aiden et al., Science 326, 289 (2009)
- [6] G. Fudenberg, M. Imakaev, C. Lu, A. Goloborodko, N. Abdennur, Cell Reports 15, 2038 (2016)
- [7] K. Luger, A.W. M\u00e4der, R.K. Richmond, D.F. Sargent, T.J. Richmond, Nature 389, 251 (2016)
- [8] S.C. Satchwell, H.R. Drew, A.A. Travers, J. Mol. Biol. 191, 659 (1986)
- [9] E. Segal et al., Nature 442, 772 (2006)
- [10] P.T. Lowary, J. Widom, J. Mol. Biol. 276, 19 (1998)
- [11] B. Eslami-Mossallam, R.D. Schram, M. Tompitak, J. van Noort, H. Schiessel, PLoS ONE 11, e0156905 (2016)
- [12] W. K. Olson, A.A. Gorin, X.-J. Lu, L. M. Hock, V.B. Zhurkin, Proc. Natl. Acad. Sci. U. S. A. 95, 11163 (1998)
- [13] F. Lankaš, J. Šponer, J. Langowski, T.E. Cheatham III, Biophys. J. 85, 2872 (2003)
- [14] N.B. Becker, L. Wolff, R. Everaers, Nucl. Acids. Res. 34, 5638 (2006)
- [15] L. de Bruin, M. Tompitak, B. Eslami-Mossallam, H. Schiessel, J. Phys. Chem. B 120, 5855 (2016)
- [16] M. Tompitak, C. Vaillant, H. Schiessel, Biophys. J. 112, 505 (2017)
- [17] K. Brogaard, L. Xi, J.-P. Wang, J. Widom, Nature 486, 496 (2012)
- [18] M. Tompitak, G.T. Barkema, H. Schiessel, BMC Bioinformatics 18, 157 (2017)
- [19] N. Kaplan *et al.*, Nature **458**, 363 (2009)
- [20] A. Valouev et al., Nature 474, 516 (2011)
- [21] D. Tillo *et al.*, PLoS ONE **5**, e9129 (2010)
- [22] T. Vavouri, B. Lehner, PLoS Genetics 7, e1002036 (2011)
- [23] S.S. Hammoud et al., Nature 460, 473 (2009)

- [24] B.G. Dias, K.J. Ressler, Nature Neuroscience 17, 89 (2014)
- [25] V.A. Bondarenko et al., Mol. Cell 24, 469 (2006)
- [26] M. Tompitak, L. de Bruin, B. Eslami-Mossallam, H. Schiessel, Phys. Rev. E 95, 052402 (2017)

B6 G-protein-coupled receptors

A. Baumann Institute of Complex Systems, ICS-4 Forschungszentrum Jülich GmbH

Contents

1	Introduction	2
2	G-protein-coupled receptors	2
	2.1 Structural properties of G-protein-coupled receptors	2
	2.2 Signalling properties of G-protein-coupled receptors	4
	2.3 Classification of G-protein-coupled receptor families	5
3	Biogenic amine receptors	7
	3.1 Biosynthesis of biogenic amines	7
	3.1.1 Biogenic amines derived from tyrosine	8
	3.1.2 Serotonin is derived from tryptophan	9
	3.1.3 Histamine is derived from histidine	10
	3.2 Biogenic amine receptors	11
	3.2.1 Histamine receptors	11
	3.2.2 Serotonin receptors	12
	3.2.3 Catecholamine receptors	15
	3.2.3.1 Dopamine receptors	16
	3.2.3.2 Norepinephrine and Epinephrine receptors	17
	3.2.4 Phenolamine receptors	
	3.2.4.1 Tyramine receptors	19
	3.2.4.2 Octopamine receptors	19
4	Summary	21
	References	
	Appendices	
	Bibliography	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

An important issue for proper development and survival of living organisms is to constantly survey, identify and respond to changing external stimuli. During evolution, cellular signalling systems have been established that equip organisms with a molecular framework to register, evaluate, and adequately react to a broad range of physical and chemical cues. Communication within and between cells is achieved by regulatory processes involving a variety of functionally distinct proteins. These proteins can be assigned to superordinate families, e.g., ion channels, transport proteins, membrane receptors, enzymes etc. In the human body, the largest gene family codes for membrane receptors that are ideally suited to register external signals. These proteins are commonly referred to as G-protein-coupled receptors (GPCRs). They are located in the cell's surface membrane. GPCRs can register signals as diverse as photons, odorants, peptides, lipids, small organic compounds, or nucleotides which act as "ligands" on these proteins. Typically, the interaction of a receptor with a ligand activates an intracellular signalling cascade that induces either biochemical or electrical cellular responses. In this tutorial I will introduce the main branches of currently known GPCR families and provide general information on the structural properties of the proteins. The focus will be on the subgroup of biogenic amine receptors including biosynthesis of their ligands, receptor coupling to intracellular effectors, some aspects of their distribution and signalling behaviour in vivo, and information on functional consequences of receptor malfunction. The interested reader may find additional information on these proteins in text books [1-3] and original literature [4-17] listed at the end of this tutorial.

2 G-protein-coupled receptors

Multicellular organisms contain many different types of cells that vary in size, shape, and function. The human body, for example, consists of about 200 cell types that collectively add up to more than 10^{14} cells, from which approximately 10^{11} *neurons* constitute the brain. Neurons equip the organism with an efficient system for fast information processing, both, in the central and the peripheral nervous system (see also chapter **D.3**). To convey information to target cells, neurons generate *action potentials*. These are small electrical discharges of the membrane potential that originate from the activity of ion channel proteins localized in the cell membrane (see also chapter **B.4**). The action potential reaches the axon of a neuron with velocities of up to 100 m/sec. When the action potential reaches the axon's ending, the *presynaptic terminal* region, it stimulates the release of chemical *transmitters*. Transmitters diffuse to the *postsynaptic terminal* region of a target cell and activate specific proteins that generate an electrical response or cause transient changes in intracellular *second messenger* concentrations of that cell. Changes in second messenger concentrations typically result from activation of GPCRs.

2.1 General properties of G-protein-coupled receptors

The cognate feature of all GPCRs is their highly conserved transmembrane topography. Based on hydrophobicity profile analyses of the amino acid sequences, GPCRs possess seven segments that completely span the *plasma membrane*. These segments are called

transmembrane-segments or -domains (TM; **Fig. 1**). Considerable efforts have been made to determine the protein structure at the atomic level. Up to now, more than 30 unique receptors have been crystallized and their structure has been solved by X-ray or cryo electron microscopy. The structural data unequivocally confirmed the predicted membrane topography with seven TM segments. Due to this highly conserved feature, GPCRs are also called seven transmembrane receptors. The crystal structure data are widely used as templates to predict the structure of other GPCRs that share at least a sequence homology of \geq 30% with the template. Following this approach, the number of available structural data of GPCRs is rapidly increasing.



Fig. 1: Schematic drawing of a GPCR's transmembrane topography. The protein spans the plasma membrane seven times (TM1 - TM7). The amino terminus $(NH_2; start of protein)$ is located extracellularly and the carboxy terminus (COOH; end of protein) is located intracellularly. Transmembrane regions are connected by alternating extracellular (EL) and intracellular (IL) loops. Some GPCRs can be post-translationally modified with fatty acids at cysteine (C) residues that tether the carboxy-terminal loop to the plasma membrane. GPCRs bind to their ligands either with residues located in the amino-terminal loop or with residues located in a 'binding pocket' formed by the TM's in the plane of the membrane. Ligand-binding leads to a conformational change of the GPCR that is conveyed via its ILs to GTP-binding (G) proteins (arrows). G-proteins transduce GPCR activation to downstream effectors.

All GPCRs are classical type II membrane proteins: the amino (N)-terminus is located on the extracellular side of the membrane whereas the carboxy (C)-terminus is located in the *cytosol*. As depicted in **Fig. 1**, the TM segments are linked by three extracellular loops (EL) that alternate with three intracellular loops (IL). Cysteine residues (C) in the cytoplasmic tail of the proteins can be modified by fatty acids and thereby tether the carboxy-terminal loop to the plasma membrane. Many GPCRs also undergo post-translational N-linked *glycosylation*, a modification of the protein with sugar residues covalently bound to asparagine residues in the N-terminal loop.



Fig. 2: Intracellular signalling pathways initiated by GPCR activation. Binding of a ligand to a GPCR induces a conformational change of the receptor that is transmitted intracellularly to a GTP-binding (G)-protein. Amino-acid residues in the primary structure of the GPCR determine the coupling specificity and efficacy to a G-protein. If an inhibitory G-protein (G_{ai}) is activated, production of cyclic adenosine monophosphate (cAMP) is blocked. If a stimulatory G-protein (G_{as}) is activated, production of cAMP is induced. If a $G_{aq/o}$ -protein is engaged, phospholipase C (PLC) is activated. This enzyme hydrolyses phosphatidylinositol-4,5-bisphosphate (PIP₂) to diacylglycerol (DAG) and inositol-1,4,5-trisphosphate (IP₃). This messenger freely diffuses in the cytosol and binds to IP₃-gated ion channels (IP₃R) located in the membrane of the endoplasmic reticulum (ER). Upon IP₃-binding the channels open and Ca²⁺ ions flow from the ER into the cytosol.

2.2 Signalling properties of G-protein-coupled receptors

Depending on the receptor class (see **2.3**), ligands can interact with a GPCR at different locations. Some receptors contain rather large N-terminal loops, e.g. glutamate receptors, adhesion receptors or peptidergic receptors. In this case, the ligand binding site is located in the N-terminal loop of the GPCR. In contrast, for most receptors belonging to the rhodopsin-like GPCR family ligand-binding takes place in a binding pocket formed by the TM regions in the plane of the membrane. For the branch of aminergic receptors, which will be discussed in detail below, highly conserved residues participating in ligand binding have been uncovered. These residues were experimentally proven by mutagenesis experiments followed by expression and pharmacological testing the receptor variants.

The physical interaction between a ligand and its receptor induces a conformational change in the GPCR that is transferred intracellularly to trimeric GTP-binding (G)-proteins. G-proteins are anchored in the plasma membrane by fatty acid residues added after *translation* of the protein. Residues in close vicinity to the plasma membrane of IL2, 3, and 4 of GPCRs determine the specificity and efficacy of G-protein activation and thus, the intracellular signalling pathway that is addressed (see **Fig. 1 & 2**).

The activated GPCR acts as a *nucleotide exchange factor* on its interacting G-protein. At rest, G-proteins contain one α , one β , and one γ subunit. The α -subunit is bound to one molecule of guanosine-di-phosphate (*GDP*). Interaction with the activated GPCR induces an exchange reaction in the G_{α} subunit: GDP is released and a molecule of *GTP* is bound. Thereby,

activation of the GPCR is transmitted to the G-protein. The GTP-bound G_{α} subunit dissociates from the $G_{\beta/\gamma}$ subunits and interacts with downstream effector enzymes (Fig. 2).

The most common forms of GPCR-induced signalling lead to transient changes of cellular cAMP or Ca²⁺ concentrations. As previously stated, amino-acid residues in IL2, 3, and 4 determine the binding specificity between a GPCR and a G-protein. Activation of an inhibitory G-protein ($G_{\alpha i}$) leads to an inhibition of *adenvlvl cyclases*. These enzymes synthesize cAMP from ATP. Thus, inhibiting their enzymatic activity leads to a reduction of the intracellular cAMP concentration. Because cAMP is an important regulator of kinases, enzymes transferring phosphate groups to proteins, cellular processes that require cAMPdependent phosphorylation are impaired. In contrast, when a stimulatory G-protein (G_{us}) is activated, adenylyl cyclase activity is enhanced. This results in production of cAMP and may lead to increased levels of protein phosphorylation. Finally, if a $G_{\alpha\alpha/o}$ -protein ($G_{\alpha\alpha/o}$) is activated, the effector enzyme phospholipase C (PLC) is activated. This enzyme hydrolyses phosphatidylinositol-4.5-bisphosphate (PIP_2) anchored in the plasma membrane. Diacylglycerol (DAG) remains membrane-bound whereas inositol-1,4,5-trisphosphate (IP₃) diffuses in the cytosol (Fig. 2). When IP_3 -gated ion channels in the membrane of the endoplasmic reticulum (ER) open upon IP₃ binding, Ca^{2+} ions flow from the ER into the cvtosol. The GPCR-induced increase of the intracellular Ca²⁺ concentration can mediate different cellular responses including Ca²⁺-dependent phosphorylation of proteins, opening of ion channels, or Ca^{2+} -dependent mechanisms of gene activation.

Although GPCR-mediated cellular cAMP and/or Ca^{2+} signals are the dominating events, additional signalling pathways do exist as well. A variety of cellular kinases can be activated leading to phosphorylation of proteins that control or modulate gene *transcription*, cell proliferation or cellular metabolism. Furthermore, certain potassium and calcium channels have been found to be directly modulated by binding to G-protein subunits in response to GPCR activation.

It is important to note that GPCR-mediated signalling typically results in a high amplification rate of the primary signal (= binding of one ligand to one receptor) because, one activated GPCR stimulates several hundred G-proteins, which in turn activate effector enzymes that convert hundreds to thousands of substrate molecules per second.

2.3 Classification of G-protein-coupled receptor families

Combined data from cloned and functionally expressed GPCR-encoding genes and from analyses of completely sequenced *genomes* suggest that vertebrates may contain up to 2000 genes coding for GPCRs. Thus, GPCRs form the largest family of membrane receptors identified so far. Notably, up to 1000 genes code for receptors engaged in detecting olfactory cues. It is worth mentioning that about 40% of prescribed drugs either directly target GPCRs or their downstream effectors. This classifies GPCRs as lead structures for people's health as well as molecules of high economic impact.

Quite as large as the GPCR gene family is the spectrum of potential ligands binding to and activating the receptors. A systematic classification system was developed based on a *phylogenetic* analysis of human GPCR sequences. In this scheme, GPCRs assemble in five distinct families: glutamate (G), rhodopsin-like (R), adhesion (A), frizzled/taste (F), and

secretin (S) receptors referred to as the **GRAFS** classification system. The predominant number of receptors (>90%) are present in the rhodopsin-like family that also includes olfactory receptors.

The phylogenetic tree depicted in **Fig. 3** nicely illustrates the GRAFS classification system. In the tree, the Glutamate receptor family is represented by 15 members: eight metabotropic glutamate receptors (GRM), two GABA receptors, a Ca^{2+} -sensing receptor (CASR), and four taste receptors (TAS). These proteins share a rather long N-terminal loop consisting of 280 – 580 amino-acid residues. For glutamate receptors it is assumed that this loop forms two distinct domains separated by a glutamate-binding cavity. The structure resembles a 'Venus fly trap' which closes upon binding of glutamate.



Fig. 3: Phylogenetic relationship between unique GPCRs identified in the human genome. For calculating the tree, sequence similarity was determined for the membrane spanning regions using the maximum parsimony method on 1000 replicas. The tree consists of five main branches representing the 'GRAFS' classification system (<u>Glutamate; Rhodopsin-like;</u> <u>Adhesion; Frizzled/TAS2; Secretin</u>). Note that members of the rhodopsin-like GPCR family are not depicted for clarity. (Figure modified from [9])

The Rhodopsin-like receptor family contains a total of 702 sequences. The sequences can be assigned to four sub-branches. The first branch (89 GPCRs) contains, e.g., opsin receptors, biogenic amine receptors, prostaglandin receptors, cannabinoid receptors, and adenosine-binding receptors. The second branch (36 GPCRs) contains receptors binding to peptidergic ligands. The third branch (59 GPCRs) contains 42 chemokine receptors, 15 somatostatin-, opioid-, and galanin-binding (SOG) receptors, and two melanin concentrating hormone (MCH) receptors. The fourth branch (518 GPCRs) contains 460 olfactory receptors. In

addition, eight glycoprotein receptors, 42 purine-binding receptors, and eight MAS oncogene receptors, which bind to angiotensin, also assemble in this branch.

The Adhesion receptor family is represented by 24 members. These receptors contain functional domains in their N-terminal loops, e.g., EGF-like repeats, mucin-like regions, as well as cysteine-rich motifs. The N-termini of these proteins are highly variable in length, ranging from 200 to 2800 amino acids and are thought to participate in cell adhesion. The Frizzled/TAS2 receptor family contains a total of 24 GPCRs. As depicted in **Fig. 3**, this family contains two distinct branches. One branch contains the frizzled receptors (FZD) that control cell fate, proliferation, and polarity during development. They are activated by glycoproteins termed *Wnt*. Binding of the ligand occurs at the N-terminal loop consisting of ca. 200 amino-acid residues. The loop contains conserved cysteine residues that might participate in ligand-binding. The second branch of the family contains taste 2 (TAS2) receptors. The 13 members are assumed to function as bitter taste receptors on the tongue.

Finally, the Secretin receptor family harbors 15 GPCRs. These receptors bind to large peptides and often act in a paracrine manner. The amino-terminal loops contain 60 to 80 amino-acid residues including one highly conserved cysteine bridge.

3. Biogenic amine receptors

Biogenic amine receptors are an important group of membrane proteins that belong to the rhodopsin-like receptor family (see **Fig. 3**). They have crucial roles in cellular signalling in both, vertebrates and invertebrates. Activation of these receptors takes place by binding to small monoaminergic compounds. These biogenic amines are biochemically synthesized from different amino acids and will be introduced in the following chapters.

3.1 Biosynthesis of biogenic amines

In both, vertebrates and invertebrates, the group of biogenic amines consists of five members (see **Tab.1**). In addition to molecules shared by both phylogenetic groups (*dopamine*, *histamine*, *serotonin*), some seem to be synthesized preferentially in either vertebrates (*norepinephrine*, *epinephrine*) or invertebrates (*tyramine*, *octopamine*). Biogenic amines are synthesized from three different amino acids. The biosynthetic pathways are described in the following sections.

vertebrates	invertebrates
dopamine	dopamine
histamine	histamine
serotonin	serotonin
norepinephrine = noradrenaline	(norepinephrine)
epinephrine = adrenaline	tyramine / octopamine

Table 1. Presence of biogenic amines in vertebrates and invertebrates

3.1.1 Catecholamines and phenolamines are derived from tyrosine

In vertebrates, the amino acid tyrosine gives rise to the *catecholamines* dopamine, norepinephrine (=noradrenaline), and epinephrine (=adrenaline) (see **Fig. 4**). Biosynthesis starts with hydroxylation in the *meta*-position of tyrosine to 3,4-dihydroxy-L-phenylalanine (L-DOPA). The reaction is catalyzed by tyrosine hydroxylase, the rate-limiting enzyme in catecholamine synthesis. In a second step L-DOPA is decarboxylated to dopamine. The conversion is mediated by the enzyme DOPA decarboxylase (DDC). In certain vertebrate neurons dopamine can be further metabolized. In these cells, the enzyme dopamine β -hydroxylase catalyzes the addition of a hydroxyl group to the β -carbon on the side chain of dopamine, resulting in the formation of norepinephrine. Finally, norepinephrine can be metabolized to epinephrine. This is achieved by the enzyme phenylethanolamine N-methyl-transferase which adds a methyl group to the nitrogen of norepinephrine.



Fig. 4: Biosynthesis of catecholamines and phenolamines from tyrosine. The biogenic amines dopamine, norepinephrine, and epinephrine are synthesized from the amino acid tyrosine by consecutive enzymatic reactions. In the first step, tyrosine hydroxylase adds a hydroxyl group (OH) to the aromatic ring. Then DOPA decarboxylase cleaves off the COOH-group from L-DOPA generating dopamine. Hydroxylation of dopamine by dopamine β -hydroxylase results in norepinephrine production. This biogenic amine can be converted to epinephrine by the enzyme phenylethanolamine N-methyltransferase. Invertebrates convert tyrosine to dopamine as shown for catecholamine biosynthesis. In addition, invertebrates use tyrosine decarboxylase to produce tyramine that can be converted to octopamine by tyramine β -hydroxylase as depicted for phenolamine biosynthesis.

While the pathway to synthesize dopamine is identical in invertebrates, norepinephrine and epinephrine have not been unequivocally identified in the fruit fly *Drosophila melanogaster*. However, low concentrations have been detected in some other insect species.

Invertebrates use another biochemical pathway to convert tyrosine into the phenolamines tyramine and octopamine (see **Fig. 4**). In a first step, tyrosine is decarboxylated to tyramine by tyrosine decarboxylase. Similar to the conversion of dopamine to norepinephrine, tyramine is hydroxylated on the β -carbon of the side chain. This reaction is catalyzed by tyramine β -hydroxylase and generates octopamine. As octopamine and norepinephrine are structurally very similar, it has been suggested that the adrenergic system (norepinephrine/epinephrine) of vertebrates and the tyraminergic/octopaminergic system of invertebrates are also functionally similar.

3.1.2. Serotonin is derived from tryptophan

In both, vertebrates and invertebrates, identical biochemical pathways exist to synthesize the indolamine 5-hydroxytryptamine (5-HT, serotonin) from L-tryptophan (see **Fig. 5**). In the first and rate-limiting reaction, a hydroxyl group is added to the indole ring in the 5'-position by the enzyme tryptophan hydroxylase. In the second reaction, 5-hydroxytryptophan is decarboxylated by DDC to serotonin. Since DDC also participates in the decarboxylation of L-DOPA to dopamine (see **Fig. 4**), a defect or loss of function of this enzyme will simultaneously result in a severely impaired production of both, dopamine and serotonin.



Fig. 5: Bioynthesis of serotonin from tryptophan. The biogenic amine serotonin (= 5-OH-tryptamin) is produced from the amino acid tryptophan by two enzymatic steps. In the first reaction tryptophan hydroxylase adds a hydroxyl group (OH) to the aromatic ring. Then the enzyme DOPA decarboxylase cleaves off the COOH-group from 5-OH-tryptophan leading to serotonin production.

3.1.3. Histamine is derived from histidine

A single decarboxylation step converts L-histidine to histamine (see **Fig. 6**). The reaction is mediated by the enzyme histidine decarboxylase. Histamine is considered as one of the most important mediators of allergy and inflammation. In the vertebrate CNS, however, histamine is synthesized from a small population of neurons located in the posterior hypothalamus. These neurons project to most cerebral areas and have been implicated in hormonal secretion, cardiovascular control, thermoregulation, and memory functions.



Fig. 6: Biosynthesis of histamine from histidine. The biogenic amine histamine is produced when the enzyme histidine decarboxylase cleaves off the COOH-group from histidine.

To protect an organism against excessive biogenic amine activity, certain inactivation mechanisms exist. Monoamine oxidase (MAO) oxidizes catecholamines to their corresponding aldehydes. The aldehydes can be further converted by aldehyde reductase or aldehyde dehydrogenase to alcohol or acid derivatives, respectively, resulting in the complete loss of a catecholamine's biological activity. While MAO is predominantly active inside the cell, the enzyme catechol-O-methyltransferase (COMT) inactivates catecholamines in the extracellular space. The enzyme transfers a methyl group (CH₃) to the *meta* hydroxy group on the catecholamine ring which impairs binding of the modified biogenic amine to the respective GPCR.

Serotonin inactivation is achieved by MAO-dependent oxidation of the biogenic amine to 5-OH-indoleacetaldehyde. The aldehyde is further converted to 5-OH-indole acetic acid by aldehyde dehydrogenase. An alternative pathway exists in the pineal gland. Here, serotonin is converted to melatonin by serotonin-N-acetyltransferase and hydroxyindole-O-methyltransferase. Although melatonin has negligible effects on serotonin receptors, it possesses a pronounced biological activity. Melatonin is involved in the synchronization of circadian rhythms, including sleep and wakefulness or blood pressure regulation, most likely by interacting with specific receptors.

Histamine can be inactivated by MAO as well as by histamine methyltransferase: The latter adds a methyl group to the tele-nitrogen (NH; see **Fig. 6**) of the biogenic amine.

3.2 Biogenic amine receptors

3.2.1 Histamine receptors

Four unique *genes* have been cloned from vertebrates encoding functional histamine receptors (H1 - H4). Histamine H_1 receptors are expressed throughout the body. High concentrations of this receptor subtype were found in the hypothalamus, aminergic and cholinergic brainstem nuclei, as well as in the thalamus and cortex. These brain regions are important for modulating the internal clock and are a main target for clinical drugs. Histamine binding to H1 receptors typically results in a behavioural status of increased wakefulness and alertness. Thus, when a patient takes antihistaminergic drugs, this may cause drowsiness because the functional activity of H1 receptors causes hives, broncho-constriction, motion sickness, and smooth muscle relaxation. It has also been found that activation of these receptors triggers the symptoms of hayfever and other allergies.

Histamine H2 receptors are also widely expressed in the body. For this receptor subtype high densities have been uncovered in the basal ganglia of the central brain, in the amygdala, the hippocampus and cortex. This receptor type is thought to regulate cellular processes causing neuronal plasticity. Furthermore, histamine H2 receptors were found on parietal cells located in the stomach lining where they control gastric acid secretion, which - in excess - can result in *gastroenteritis*. Histamine H2 receptors are also present on white blood cells (*neutrophils*). Activation of these receptors can inhibit antibody and cytokine production.

receptor	G-protein coupling	effect
H1	G _{q/11}	$IP_3/Ca^{2+}\uparrow$
H2	Gs	cAMP ↑
Н3	G _{i/o}	cAMP↓
H4	G _{i/o}	cAMP↓

 Table 2. Cellular signalling by histamine receptors

High levels of histamine H3 receptor expression were found in the anterior cerebral cortex, the cerebellum, substantia nigra, and the brainstem. In these brain regions, the receptors modulate axonal and synaptic plasticity. An interesting feature of H3 receptors has been uncovered *in vivo*: they are active even in the absence of histamine, a phenomenon called, constitutive activity. The receptors were also found to impair histamine biosynthesis in the body: the more H3 receptors are activated, the less histamine is produced.

The youngest member of the histamine receptor family, the H4 receptor, was molecularly cloned in 2001. It shares cellular signalling properties with histamine H3 receptors (see **Table 2**). Histamine H4 receptors were identified in many cells of hematopoietic origin (blood cells), as well as the spleen, lung, liver, gut, and neurons. Functionally, H4 receptors regulate the levels of white-blood cell release from bone marrow.

3.2.2 Serotonin receptors

The family of serotonin (5-HT)-activated GPCRs contains 13 members (see **Table 3**). Based on sequence similarity and functional signalling properties, 5-HT receptors can be subdivided into six groups. Generally, the genes encoding 5-HT receptors are widely expressed across the mammalian brain.

receptor	G-protein coupling	effect
5-HT _{1A/1B/1D/1E/1F}	G _{i/o}	cAMP↓
5-HT _{2A/2B/2C}	G _{q/11}	$IP_3/Ca^{2+}\uparrow$
5-HT ₄	Gs	cAMP ↑
5-HT _{5A/(5B)}	G _{i/o}	cAMP↓
5-HT ₆	Gs	cAMP ↑
5-HT ₇	Gs	cAMP ↑

Table 3. Cellular signalling by serotonin (5-HT) receptors

The group of 5-HT₁ receptors contains six members. 5-HT_{1A} receptors are expressed in the limbic system and mesencephalic raphe nuclei where they are located on the surface of postsynaptic neurons. Activation of 5-HT_{1A} receptors primarily leads to an inhibition of adenylyl cyclase activity and results in a reduction of intracellular cAMP. The presence of 5-HT_{1A} receptors in high density in the limbic system strongly suggests that effects of serotonin or serotonergic drugs on emotional states and behavior could be mediated by these receptors. In the neocortex, 5-HT_{1A} receptors may also be involved in cognitive or integrative functions. In the dorsal and median raphe nuclei, the 5-HT_{1A} receptors were found to act as *somatodendritic* autoreceptors and thus, to modulate the activity of serotonergic neurons in a negative feedback loop. Activation of the receptors causes a decrease in the firing rate of serotonergic neurons and a reduction in the release of serotonin from serotonergic terminals. The effect most likely results from 5-HT_{1A}-dependent activation of G_{β/γ} subunits that can bind to potassium channels. Upon interaction with G_{β/γ} subunits these channels open and cause neuronal hyperpolarization which impairs neurotransmitter release.

Like 5-HT_{1A} receptors, 5-HT_{1B} and 5-HT_{1D} receptors mediate their cellular effects by inhibiting adenylyl cyclase activity. The receptors are expressed in basal ganglia of the substantia nigra, in the limbic system (amygdala, hippocampus), as well as on vascular smooth muscles. There is experimental evidence that the receptors are located on presynaptic terminals of serotonergic neurons where they modulate the release of their own ligand, serotonin. The receptors also have been identified postsynaptically. Here, they are thought to modulate the release of other neurotransmitters, such as acetylcholine in the hippocampus or dopamine in the prefrontal cortex. The presence of these receptors in high density in the basal ganglia raises the possibility that they may play a role in diseases which involve these structures, such as Parkinson's disease.

The 5-HT_{1E} receptor has been identified in the frontal cortex and the limbic system (amygdala, hippocampus). The function of the 5-HT_{1E} receptor in intact tissue is not known. However, when expressed after transfection of cultivated cells, the receptor inhibits adenylyl cyclase activity. With 64% sequence homology, the amino-acid sequence of 5-HT_{1E} shares a high degree of similarity with the 5-HT_{1D} receptor.

The 5-HT_{1F} receptor is expressed in the cortex, hippocampus, and raphe nuclei. The receptor sequence is very similar to the 5-HT_{1E} receptor (~61% homology). Activation of 5-HT_{1F} receptors causes inhibition of adenylyl cyclase activity. *In vivo*, 5-HT_{1F} receptor activation inhibits neurogenic dural inflammation, a status frequently occurring as a part of migraine-associated pain pathophysiology.

The group of 5-HT₂ receptors contains three members (see **Table 3**). The receptors stimulate hydrolysis of PIP₂ by phospholipase C into DAG and IP₃, thereby causing release of Ca²⁺ from intracellular stores (see **Fig. 2**). High densities of $5HT_{2A}$ receptors were found in many cortical areas, especially in the frontal cortex. The receptors are also present in the limbic system and in a region which is connected to the visual cortex (claustrum). 5-HT_{2A} receptors in the cortex are located postsynaptically. Activation of 5-HT_{2A} receptors mediates neuronal depolarization, a result of the closing of potassium channels. In pharmacological studies it was uncovered that hallucinogenic amphetamine derivatives, like 2,5-Dimethoxy-4-bromoamphetamine (DOB), act as agonists (activators) on this receptor sub-type.

The 5-HT_{2B} receptor was among the first of the serotonin receptors to be pharmacologically characterized. In humans, 5-HT_{2B} receptor mRNA has been found, e.g., in cerebellum, cerebral cortex, amygdala, substantia nigra, caudate, and thalamus. It is assumed that 5-HT_{2B} receptors play an important role in anxiety. Furthermore, the receptors are present in the stomach fundus where their activation results in the contraction of fundus smooth muscle. As shown in **Table 3**, 5-HT_{2B} receptors cause inositol lipid hydrolysis followed by an increase of intracellular Ca²⁺.

Finally, 5-HT_{2C} receptors are present in high density on epithelial cells of the *choroid plexus*. It has been proposed that activation of these receptors could regulate the composition and volume of the cerebrospinal fluid. 5-HT_{2C} receptors are also found throughout the brain, especially in the limbic system, as well as in regions associated with motor behavior, including the substantia nigra and globus pallidus. In these regions, the density of 5-HT_{2C} receptors is much lower than in the choroid plexus. Activation of neural 5-HT_{2C} receptors may result in hypolocomotion, hyperthermia, and anxiety.

Serotonin 5-HT₄ receptors are located primarily in the striatum, substantia nigra, hippocampus, and olfactory tubercle. The receptor exerts its activity by stimulating adenylyl cyclase. In the striatum, activated 5-HT₄ receptors indirectly mediate an enhancement of dopamine release, although these receptors seem to be absent from striatal dopaminergic terminals. Furthermore, 5-HT₄ receptors located on neurons of the myenteric plexus of the ileum, smooth muscle cells, and secretory cells can evoke secretion and the peristaltic reflex.

The subgroup of 5-HT₅ receptors consists of two members. Stimulation of the 5-HT_{5A} receptor leads to an inhibition of adenylyl cyclase activity. 5-HT_{5A} receptors are expressed in the cerebral cortex, limbic system, cerebellum, and thalamic nuclei. Although 5-HT_{5A} and 5-HT_{5B} receptors share ~77% sequence similarity to each other, functional expression of the human *5-ht_{5b}* gene has not been achieved, so far. However, transcripts of the 5-HT_{5B} receptor have been detected in the limbic system (hippocampus), habenula, and dorsal raphe nucleus. In addition to neuronal expression, 5-HT_{5A} receptors are also located on *astrocytes*. Relatively little is known about the functional role of 5-HT_{5A} receptors, but studies on transgenic mice suggest that the receptor is engaged in controlling animal's locomotion.

Serotonin 5-HT₆ receptors are located in the striatum, nucleus accumbens and hippocampus. Activation of these receptors stimulates adenylyl cyclase activity. Physiologically, 5-HT₆ receptors modulate cholinergic and dopaminergic neurotransmission, and were associated with spatial learning and memory.

Another member of the serotonin GPCR family is the 5-HT_7 receptor. Like 5-HT_6 receptors, 5-HT_7 receptors are also positively coupled to adenylyl cyclase activity (see **Table 3**). 5-HT_7 receptors are located in the thalamus, hypothalamus, and hippocampus, as well as in the periphery. The 5-HT_7 receptor is involved in thermoregulation, circadian rhythm, learning, and memory. Furthermore, 5-HT_7 receptors may be involved in mood regulation. Therefore, these receptors could serve as targets in treating mood-associated diseases, e.g., depression.

In summary 5-HT receptors have important physiological functions including thermoregulation, appetite control, regulation of sleep, and regulation of mood. They also play a role in aggression, sexual behavior, and cardiovascular disorders. Various neuropsychiatric disorders like anxiety, depression, schizophrenia, migraine, and drug abuse are controlled and/or modulated by serotonin and serotonin receptor activation. In this respect, 5-HT receptors are the target of a variety of pharmaceutical drugs including antidepressants, antipsychotics, antiemetics, antimigraine drugs, or hallucinogens. Ongoing research is devoted to develop and/or identify selective agents targeting specific receptor sub-type(s) to treat or alleviate the disorders symptoms.

3.2.3 Catecholamine receptors

 $\beta_{1/2/3}$ -adrenergic

The family of catecholamine receptors includes GPCRs activated by dopamine, norepinephrine and epinephrine. Two major types of dopamine receptors exist and are classified according to their positive (D₁/D₅) or negative (D₂/D₃/D₄) coupling to adenylyl cyclase (see **Table 4**). Norepinephrine and epinephrine bind to adrenergic receptors which are classified as α - and β -adrenergic receptors (see **Table 4**). Two subclasses of α -adrenergic receptors are known: α_1 - and α_2 -receptors. Both subclasses contain three members, α_{1A-C} and α_{2A-C} . Whereas α_1 -receptors cause inositol lipid hydrolysis followed by an increase of intracellular Ca²⁺, α_2 -receptors inhibit adenylyl cyclase activity. Furthermore, three β -adrenergic receptor subtypes are known (β_{1-3}) which are all positively coupled to adenylyl cyclase.

cAMP ↑

receptor	G-protein coupling	effect
dopamine D _{1/5}	Gs	cAMP ↑
dopamine D _{2/3/4}	G _{i/o}	cAMP↓
$\alpha_{1A/B/C}$ -adrenergic	G _{q/11}	$IP_3/Ca^{2+}\uparrow$
$\alpha_{2A/B/C}$ -adrenergic	G _{i/o}	cAMP↓

Table 4. Cellular signalling by catecholamine receptors

Gs

As mentioned in chapter 2.2, binding of a ligand to these receptors takes place in a binding pocket formed by the TM regions in the plane of the membrane. Specific residues located in different TM-segments are highly conserved in all catecholamine and phenolamine receptors (see Fig. 7). An aspartic acid residue (D) in TM3, serine residues (S) in TM5, and a phenylalanine residue (F) in TM6 were identified to determine the ligand binding properties of biogenic amine receptors. The main energetic contribution to ligand - receptor interaction originates from ionic bonding between the carboxyl-group of the aspartic acid residue in TM3 and the protonated amino group (NH₃⁺) of the biogenic amine. In TM5, serine residues arranged in an "S-X-X-S" motif (X = any, non-charged amino acid) form hydrogen bonds with hydroxyl residues on the catechol- or phenyl-ring of catecholamines or phenylamines, respectively. The aromatic amino-acid phenylalanine in TM6 forms π - π electron bonds with the aromatic ring of the biogenic amine, thereby stabilizing the ligand - receptor interaction.



Fig. 7: Schematic drawing of a GPCR with its ligand binding residues. Amino acids in TM segments participating in binding to catecholamines or phenolamines are indicated. Interaction of a ligand to an aspartic acid residue (D) in TM3, serine residues (S) in TM5, and a phenylalanine residue (F) in TM6 result in a high affinity binding status.

3.2.3.1 Dopamine receptors

Dopamine receptors contribute to many neurological processes, e.g., motivation, reward, memory, learning, motor control, and modulation of neuroendocrine signalling. Abnormal signalling of the dopaminergic system can cause a variety of neuropsychiatric disorders including Tourette's syndrome, Parkinson's disease, schizophrenia, or attention-deficit hyperactivity disorder (ADHD). Thus, dopamine receptors serve as drug targets for antipsychotics and psychostimulants.

The dopamine D_1 receptor family consists of the D_1 - and D_5 receptors. The D_1 receptor is the most abundant dopamine receptor expressed in the brain. High levels are found in the cortex,

striatum, and limbic system. D_1 receptors are exclusively present on postsynaptic neurons. Recent data suggest that activated D_1 receptors have synergistic effects on D_2 receptormediated motor responses. This led to the development of receptor agonists for the treatment of Parkinson's disease.

Like D_1 receptors also the D_5 receptors are positively coupled to adenylyl cyclase and cause a rise in intracellular cAMP upon activation. D_5 receptors are located in the cortex and the limbic system of the brain. Based on their distribution pattern it is assumed that D_5 receptors play a role in modulating behaviour, emotion, and long-term memory. It has also been uncovered that D_5 receptors have a 10-fold higher affinity for dopamine than D_1 receptors. Recently it has been suggested that D_5 receptors play a role in blood pressure regulation, thus D_5 receptor-specific drugs could be promising candidates for treating hypertension.

Members of the dopamine D_2 receptor family ($D_{2/3/4}$) share the functional property to inhibit adenylyl cyclase activity. D_2 receptors are present in many brain regions, with the highest levels expressed in the basal ganglia of the striatum. The basal ganglia participate in motor control and learning. Together with D_1 and D_3 receptors, D_2 receptors are the primary determinants of locomotor control. The roles of D_2 and D_3 receptors, compared to D_1 receptors, are more complex because both receptor subtypes are expressed pre- and postsynaptically. In the presynaptic terminus they act as autoreceptors, providing an important negative feedback mechanism that adjusts neuronal firing rate, synthesis, and release of neurotransmitter in response to changes in extracellular dopamine levels. Activation of presynaptic D_2 autoreceptors generally causes a decrease in dopamine release that results in decreased locomotor activity, whereas activation of postsynaptic receptors (e.g., D_1 receptors) stimulates locomotion. In addition to effects on motor control, D_2 receptors are an essential target for antipsychotic drugs applied in treating schizophrenia and other idiopathic psychotic disorders.

High levels of dopamine D_3 receptors are present in two regions of the limbic system, i.e., the islands of Calleja and the nucleus accumbens. Both regions are involved in the reinforcing effects of pleasurable activities and it has been shown that D_3 receptors are associated with reward and reinforcement mechanisms. Altering or manipulating dopamine receptor function can result in a significant modulation of the responses to natural rewards as well as addictive drugs. Notably, certain D_3 receptor ligands were found to antagonize reinforcing behaviors of cocaine and alleviating DOPA-induced dyskinesia in patients with Parkinson's disease. As described for D_2 receptors, there is evidence that D_3 autoreceptors also contribute to the presynaptic regulation of dopamine release and thereby complement D_2 autoreceptor's function.

The dopamine D_4 receptor is expressed in the cortex, amygdala, hypothalamus, and pituitary. The receptor is involved in exploratory behavior and motor coordination. Mutations of the receptor-encoding gene can cause various behavioral phenotypes, including attention-deficit hyperactivity disorder (ADHD), autonomic nervous system dysfunction, and schizophrenia. Receptor variants have been uncovered harboring variable numbers of tandem repeats in the coding sequence of the gene's *exon* 3. One of these variants has a greater incidence of causing ADHD. In addition, the D_4 receptor is a target for drugs that are used in treating schizophrenia or Parkinson's disease.

3.2.3.2 Norepinephrine and epinephrine receptors

Norepinephrine and epinephrine, which are almost exclusively present in vertebrates, both can activate α - and β -adrenergic receptors. While β -adrenergic receptors are a target of " β -blockers" to treat hypertension or congestive heart failure, the role of α -adrenergic receptors is less clear. Yet, several anti-depressants and anti-psychotic drugs block the activity of α_1 -adrenergic receptors.

Molecular cloning and analysis of the deduced amino acid sequences revealed that there are three subtypes of α_1 -adrenergic receptors ($\alpha_{1A/B/C}$) and also three subtypes of α_2 -adrenergic receptors ($\alpha_{2A/B/C}$). The receptors have been identified in the brain and in peripheral tissues. Activation of α_1 -adrenergic receptors causes hydrolysis of PIP₂ into DAG and IP₃ by phospholipase C followed by Ca²⁺ release from intracellular stores. The receptors have been localized both, pre- and postsynaptically. It is assumed that presynaptic receptors act as autoreceptors as previously described for dopamine receptors (see 3.2.3.1). Some agonists acting selectively on α_1 -adrenergic receptors are used as nasal decongestants. Antagonists acting on these receptors are used to treat hypertension and benign prostatic hyperplasia. The subgroup of α_2 -adrenergic receptors is activated by epinephrine and with lower potency by norepinephrine. The receptors are negatively coupled to adenylyl cyclase. Some agonists of α_2 -adrenergic receptors can be used to control hypotension and bradycardia, can induce hypnotic effects and analgesia, and modulate seizure activity as well as platelet aggregation. In contrast to agonists, antagonists acting on these receptors can act as anti-depressants. There is evidence that the α_{2B} subtype participates in neurotransmission in the spinal cord, whereas the α_{2C} subtype regulates catecholamine release from adrenal chromaffin cells.

As with α_1 - and α_1 -adrenergic receptors, three subtypes of β -adrenergic receptors exist. β_1 adrenergic receptors are mainly expressed in the heart and in the cerebral cortex. In the cerebellum and the lung β_2 -adrenergic receptors predominate, but in many tissues β_1 - and β_2 adrenergic receptors are co-expressed and fulfill the same physiologic function. A third subtype of β -adrenergic receptor (β_3) has been identified whose pharmacological properties are clearly distinct from β_1 -/ β_2 -adrenergic receptors. β_3 -adrenergic receptors are mainly expressed in adipose tissue. All three receptors share the property to activate adenylyl cyclase upon binding to norepinephrine or epinephrine. Physiologically, activation of β -adrenergic receptors increases heart rate, stimulates cardiac contraction, and dilates the bronchi as well as blood vessels. Agonists acting on β_1 -adrenergic receptors are used to acutely treat cardiogenic shock. Antagonists acting on β_1 -adrenergic receptors are applied to treat hypertension, cardiac arrhythmias and cardiac failure. Agonists of β_2 -adrenergic receptors are used to treat respiratory disorders as well as to overcome bronchial construction. The β_3 -adrenergic receptor has been linked to hereditary obesity, regulation of lipid metabolism and thermogenesis, and the development of diabetes. Some β_3 -receptor selective agonists evoke anti-stress responses and cause nonshivering thermogenesis in animal studies.

In summary, catecholamine receptors have important neurological and physiological functions. Dopamine receptors participate in neural processes, including motivation, learning, memory, reward, motor control, as well as neuroendocrine signalling. Activation of α -adrenergic receptors results in constriction of blood vessels, relaxation of intestinal muscles, and dilation of the pupils. Activation of β -adrenergic receptors stimulates cardiac contraction, heart rate, but also thermogenesis. Therapeutically, patients suffering from tachycardia, high

blood pressure, or angina pectoris are often treated with drugs blocking β -adrenergic receptor activity.

3.2.4 Phenolamine receptors

Low levels of tyramine, octopamine, tryptamine, and β -phenylethylamine have been identified in mammals. These substances account for less than 1% of all biogenic amines present in the mammalian brain. Because of their low abundance, they are called "trace amines". Interestingly, the level of trace amines is altered in various human disorders including, e.g., depression, hepatic encephalopathy, hypertension, Parkinson's disease, and schizophrenia. In contrast to vertebrates, the phenolamines tyramine and octopamine have important physiologic functions in invertebrates where they act as neurotransmitters, neuromodulators, and neurohormones. At present, two types of tyramine receptors (TAR) have been characterized which are classified according to their negative (TAR1) or positive (TAR2) coupling to adenylyl cyclase (see **Table 5**). The other important phenolamine of invertebrates, octopamine, binds to OCT α - or OCT β -type receptors. The receptors are classified according to their functional relationship to adrenergic receptors (see **Table 4**). It should be noted that TAR- and OCT-receptors can be activated by both phenolamines, but the receptor families differ in their preferential binding to either tyramine or octopamine.

receptor	G-protein coupling	effect
Tyramine (TAR1)	G _{i/o}	cAMP↓
Tyramine (TAR2)	Gs	cAMP ↑
Octopamine (OCTaR)	G _{q/11}	$IP_3/Ca^{2+}\uparrow$
Octopamine (OCTβR)	Gs	cAMP ↑

Table 5. Cellular signalling by phenolamine receptors

3.2.4.1 Tyramine receptors

The first tyramine receptor was molecularly cloned from the fruit fly, *Drosophila melanogaster*. Although additional receptors were isolated from different invertebrate species and/or uncovered in databases containing completely sequenced genomes, only a small number of these genes were functionally expressed and pharmacologically characterized. The majority of these receptors inhibit cAMP production in response to tyramine application. A human and a rat trace amine receptor were also identified that are more potently activated by tyramine than by octopamine. These mammalian receptors cause stimulation of adenylyl cyclase upon ligand activation. This finding was surprising because it was generally assumed that tyramine receptors are negatively coupled to adenylyl cyclase, at least in invertebrates. This hypothesis currently changed with the characterization of a TAR2 receptor from the honeybee, *Apis mellifera*. The receptor specifically activates adenylyl cyclase in response to tyramine and thus, a second family of tyramine receptors can be defined whose members stimulate cAMP production.

Tyramine (and octopamine) is present in high concentrations in the CNS and periphery of invertebrates. For a long time, tyramine was considered only as a biochemical precursor of octopamine (see **Fig. 4**) and might not have significant neuroactive functions. This assessment changed once tyramine-specific receptors had been characterized. Nevertheless, knowledge about the physiological role of tyramine is scarce. So far, tyramine has been shown to alter the behavioral sensitivity of *Drosophila melanogaster* to cocaine and causes an increase in chloride conductance across the fly's *Malpighian tubule*. In addition, tyramine is considered to act as a neuromodulator in the olfactory system and at the neuromuscular junction of this insect. In the cockroach, tyramine stimulates *trehalose* metabolism in isolated fat bodies, inhibits the contraction of locust visceral muscles and, after repeated injections, it reduces locust viability.

3.2.4.2 Octopamine receptors

The physiological role of octopamine has been studied *in vivo* extensively. In intact organ preparations, in membrane preparations of various tissues, and in various insect cell lines, octopamine application evoked either increases in intracellular Ca^{2+} and/or cAMP concentrations. In insects, octopamine is referred to as a "flight or fight" hormone or a "sympathetic" circulation hormone because, like the adrenergic system of vertebrates, the octopaminergic system seems to adapt animals to energy-demanding situations. In addition, octopamine influences the response characteristics of sensory organs.

Octopamine plays an important role during insect flight. During long flight periods, the flight muscles switch to lipid metabolism as the energy source. It is assumed that a decrease in octopaminergic neurotransmission switches off the glycolytic pathway and thereby turns on lipid metabolism. In a variety of behavioral tests, octopamine has been shown to reduce response thresholds and habituation rates of feeding responses in honeybees and flies. Visual responses in bees and locusts are affected, as well. Other behaviors that are induced or modulated by octopamine are pharyngeal pumping, locomotion, and egg-laying in *Caenorhabditis elegans*, firefly flashing, as well as nestmate recognition and the onset of foraging in honeybees.

In recent years, a large number of octopamine receptors have been molecularly cloned and pharmacologically characterized from various invertebrates including fruit flies, honeybees, moths, sea slugs (*Alplysia californica*) or pond snails (*Lymnea stagnalis*). The proteins possess signalling properties that largely confirm what is known from native receptors. The proteins can be assigned to two distinct groups according to sequence similarity and the intracellular signalling pathway that is activated (see **Table 5**).

We uncovered a fascinating signalling behavior of an octopamine receptor from *Drosophila melanogaster*. Heterologously expressed DmOCT α 1B receptors cause Ca²⁺ oscillations in the continuous presence of octopamine (see **Fig. 8**). Using a series of pharmacological experiments combined with mutagenesis and functional expression of the mutant receptors, we identified a single amino-acid residue in the third intracellular loop of the receptor (threonine; T₃₅₂) that undergoes transient phosphorylation/dephosphorylation cycles resulting in termination and re-activation of the receptor's signalling capabilities.





Fig. 8: An octopamine receptor (DmOCTaR1B) causes oscillation of the intracellular Ca^{2+} concentration. The receptor was expressed in a cell line. Cells were loaded with a Ca^{2+} -sensitive fluorescent dye to monitor changes in intracellular Ca^{2+} concentrations. A) Without ligand (ES; extracellular solution) no change in fluorescence is registered. When octopamine is applied, an oscillatory change of the Ca^{2+} -dependent fluorescence is observed. B) Octopamine induced Ca^{2+} oscillations, like in A), were blocked by co-application of a PKC inhibitor (PKC block). Note the elevated level of Ca^{2+} -dependent fluorescence in the presence of the inhibitor. C) The phosphorylation site in the third intracellular loop is depicted. Phosphorylation (+P) causes uncoupling of the GPCR from its signalling pathway and results in a decrease of the fluorescence to basal values. When the phosphate group is cleaved off by a phosphatase (-P; dephosphorylation), the receptor is reactivated and can cause another Ca^{2+} -dependent fluorescence signal.

Octopamine binding to DmOCT α 1B results in IP₃-mediated Ca²⁺ release from intracellular stores. The increase in Ca²⁺ activates protein kinase C (PKC) which phosphorylates a threonine residue (T₃₅₂) in the third IL of the receptor. The phosphorylated receptor is desensitized, i.e., the protein is impaired to further stimulate the G_q-protein. From the desensitized state, the receptor is relieved by dephosphorylation mediated by cell endogenous phosphatases. As octopamine is still present in the external solution, the receptors are reactivated and can cause another rise of the intracellular Ca²⁺ concentration. Receptor desensitization and resensitization are efficient biochemical regulatory circuits to prevent overshooting cellular reactions that may have detrimental effects on cell survival or propagation.

Like β -adrenergic receptors, OCT β -type receptors are positively coupled to adenylyl cyclase. Three receptor subtypes have been characterized from the fruit fly, four from the honeybee, and also four from the bay barnacle. In insects, the receptors are expressed in brain regions participating in sensory information processing as well as in learning and memory. Many behavioral reactions have been attributed to the signalling action of octopamine and especially to octopamine-evoked changes in cellular cAMP levels. Activation of octopamine receptors has been shown to modulate the responsiveness of sensory receptors, interneurons, and *motor neurons* and thus to affect complex behavioral responses.

4. Summary

This chapter aimed to introduce some aspects of the molecular, structural, and functional signalling properties of GPCRs. With a total number of ≥ 2000 members, these membrane proteins are encoded by the largest gene family present in mammalian genomes. The receptors are activated by a broad range of ligands including photons, odorants, flavors, peptides, amino acids, as well as derivatives of amino acids like biogenic amines. The receptors share a common transmembrane topography, characterized by seven membranespanning segments. Ligand binding to a receptor induces a conformational change that is transduced to trimeric G-proteins. Subunits of these G-proteins activate downstream effectors which cause transient changes in intracellular concentrations of ions and cyclic nucleotides. Depending on the affected cellular pathway, specific responses to the original signal are evoked. GPCR-mediated signalling provides cells with a mechanism resulting in a high amplification rate of the primary stimulus. In this tutorial, emphasis has been put to cover biogenic amine receptors in more detail. The role of histamine receptors in triggering symptoms of havfever or allergies has been discussed. Serotonin receptors have been introduced which participate in thermoregulation, appetite control, regulation of sleep, mood, or aggression. Neuropsychiatric disorders like anxiety, depression, schizophrenia, or migraine can be caused by mal-functions of the serotonergic system. Dopamine receptors participate in neurological processes, like motivation, reward, motor control, and modulation of neuroendocrine signalling. Abnormal signalling of the dopaminergic system can cause Tourette's syndrome, Parkinson's disease, or ADHD. Adrenergic receptors are well known for their role in treating hypertension or congestive heart failure. Finally, tyramine and octopamine receptors, which are almost exclusively expressed in invertebrates, are known to adapt animals to energy-demanding situations, but also to have modulatory functions in tuning the sensitivity of sensory receptors and to affect complex behavioral responses. In mammals, GPCRs are important drug targets with pharmaceuticals covering antipsychotics, psychostimulants, antidepressants, antiemetics, antimigraine drugs, or hallucinogens. In order to treat and alleviate symptoms of disorders resulting from GPCR mal-function it is still necessary to develop or improve drugs targeting receptor sub-type(s) specifically. With the availability of 3D structural data from an increasing number of GPCRs this challenge can now be addressed by rational drug design, i.e., fitting a potential ligand to a given GPCR structure.

References

[1] Berg, J.M., Stryer, L., Tymoczko, J.L. and Gatto, G.J. "*Biochemistry*" 8th Ed. (W.H. Freeman & Co., New York, 2015). See especially chapter 14 and 33.

[2] Brady, S.T., Siegel, G.J., Albers, R.W. and Price, D.L. "*Basic Neurochemistry*" 8th Ed. (Elsevier Ltd, Oxford, 2012). See especially part II and III.

[3] Kandel, E.R., Schwartz, J.H., Jessell, T.M., Siegelbaum, S.A., and Hudspeth, A.J. "*Principles of Neural Science*" 5th Ed. (Mcgraw-Hill Pub. Comp., New York, 2012)

[4] Alexander, S.P.H. et al. (2013) The concise guide to pharmacology 2013/14: G proteincoupled receptors. British J. Pharmacol. <u>170</u>, 1459-1581. doi: 10.1111/bph.12444/full

[5] Balfanz, S., Strünker, T., Frings, S. and Baumann, A. (2005) A family of octopamine receptors that specifically induce cyclic AMP production or Ca²⁺ release in *Drosophila melanogaster*. J. Neurochem. <u>93</u>, 440-451. doi: 10.1111/j.1471-4159.2005.03034.x

[6] Balfanz, S., Jordan, N., Langenstück, T., Breuer, J., Bergmeier, V. and Baumann, A.
(2014) Molecular, pharmacological, and signaling properties of octopamine receptors from honeybee (*Apis mellifera*) brain. J. Neurochem. <u>129</u>, 284-296. doi: 10.1111/jnc.12619

[7] Beaulieu, J.-M. and Gainetdinov, R.R. (2011) The physiology, signaling, and pharmacology of dopamine receptors. Pharmacol. Rev. <u>63</u>, 182-217. doi: 10.1124/ pr.110.002642

[8] Blenau, W. and Baumann, A. (2016) Octopaminergic and tyraminergic signaling in the honeybee (*Apis mellifera*) brain: Behavioral, pharmacological, and molecular aspects. In: Farooqui, T. and Farooqui, A.A. (Eds), Trace Amines and Neurological Disorders. Oxford: Academic Press, 203-220.

[9] Fredriksson, R., Lagerström, M.C., Lundin, L.-G. and Schiöth, H.B. (2003) The Gprotein-coupled receptors in the human genome form five main families. Phylogenetic analysis, paralogon groups, and fingerprints. Mol. Pharmacol. <u>63</u>, 1256-1272. doi: 10.1124/mol.63.6.1256

[10] Hoff, M., Balfanz, S., Ehling, P., Gensch, T. and Baumann, A. (2011) A single amino acid residue controls Ca²⁺ signaling by an octopamine receptor from *Drosophila melanogaster*. FASEB J. <u>25</u>, 2484-2491. doi: 10.1096/fj.11-180703

[11] Huang, Y. and Thathiah, A. (2015) Regulation of neuronal communication by G proteincoupled receptors. FEBS Lett. <u>589</u>, 1607-1619. doi: 10.1016/j.febslet.2015.05.007 [12] Lee, S.-M., Booe, J.M. and Pioszak, A.A. (2015) Structural insights into ligand recognition and selectivity for class A, B, and C GPCRs. Eur. J. Pharmacol. <u>15</u>, 196-205. doi: 10.1016/j.ejphar2015.05.013

[13] Limbird, L.E. (2011) Historical perspective for understanding of adrenergic receptors. Current Topics in Membranes <u>67</u>, 1-17. doi: 10.1016/B978-0-12-384921-2.00001-X

[14] McCorvy, J.D. and Roth, B.L. (2015) Structure and function of serotonin G protein coupled receptors. Pharmacol. Ther. 150, 129-142. doi:10.1016/j.pharmthera.2015.01.009

[15] Panula, P. and Nuutinen, S. (2013) The histaminergic network in the brain: basic organization and role in disease. Nat. Rev. Neurosci. <u>14</u>, 472–487. doi: 10.1038/nrn3526

[16] Reim, T., Balfanz, S., Baumann, A., Blenau, W., Thamm, M. and Scheiner, R. (2017) AmTAR2: functional characterization of a honeybee tyramine receptor stimulating adenylyl cvclase activity. Insect Biochem. Mol. Biol. 80, 91-100. doi:10.1016/j.ibmb.2016.12.004

[17] Stevens, R.S., Cherezov, V., Katritch, V., Abagyan, R., Kuhn, P., Rosen, H. and Wüthrich, K. (2013) GPCR Network: a large-scale collaboration on GPCR structure and function. Nat. Rev. Drug Discov. <u>12</u>, 25-34. doi:10.1038/nrd385

Appendices

For critical reading and comments on the manuscript I'd like to acknowledge M. Deutsch and S. Balfanz (both ICS-4).

A Glossary

action potential	small and transient electrical discharge of the membrane potential		
adenylyl cyclase	enzyme converting ATP to cAMP		
amino acid	monomeric unit of a protein		
amino terminus	start of a protein characterized by the free amine group $(-NH_2)$		
astrocyte	astrocytes ('star-like cells') are the most numerous and diverse		
5	neuroglial cells in the central nervous system		
ATP	adenosine triphosphate		
axon	long process of a neuron that conducts nerve impulses (= action		
	potentials)		
cAMP	cyclic adenosine 3',5' monophosphate		
carboxy terminus	end of a protein, terminated by a free carboxyl group (-COOH).		
catecholamine	contains a catechol group (benzene ring with two hydroxyl		
	groups) and an attached amine (NH ₂) group		
choroid plexus	choroid plexus is a plexus of cells that produces the		
	cerebrospinal fluid in the ventricles of the brain		
cytosol	portion of a cell's content outside the nucleus		
dopamine	neurotransmitter synthesized from the amino acid tyrosine		
endoplasmic reticulum	system of double membranes in the cytoplasm of eukaryotic		
	cells		
epinephrine	syn. adrenalin; neurotransmitter synthesized from the amino acid tyrosine		
exon	part of gene coding for protein sequence		
gastroenteritis	inflammation of the lining of the intestines caused by a virus,		
	bacteria or parasites		
GDP	guanosine diphosphate		
gene	chromosomal segment that codes for a functional protein or		
	RNA		
genome	all the genetic information encoded in a cell		
glycosylation	post-translational modification of a protein, sugar residues are		
	covalently bound to, e.g., asparagine residues		
GPCR	GTP-binding protein coupled receptor serves to register external		
	signals or ligands		
GTP	guanosine triphosphate		
histamine	neurotransmitter synthesized from the amino acid histidine		
Malpighian tubule	excretory and osmoregulatory system found in insects and other		
	invertebrates		
motor neuron	neuron conducting electrical impulses from the spinal cord to muscles		
neuron	cell in the central or peripheral nervous system		
neutrophils	most common type of white blood cell		

norepinephrine	syn. noradrenalin; neurotransmitter synthesized from the amino acid tyrosine		
nucleotide exchange factor	proteins that stimulate the exchange of nucleoside diphosphates for nucleoside triphosphates bound to other proteins		
octopamine	invertebrate-specific neurotransmitter synthesized from the amino acid tyrosine		
phylogenetics	study of the evolutionary history and relationship among individuals or groups of organisms deduced from DNA or protein sequences		
plasma membrane	physical barrier that surrounds the cell surface and encloses the cytoplasm		
presynaptic terminal	distal termination of an axon specialized for signal transmission to target cells		
postsynaptic terminal	specialized region of a cell equipped for transmitter perception recognition released from presynaptic terminus		
RNA	ribonucleic acid		
second messenger	intracellular signalling molecules triggering physiological changes of a cell		
serotonin	neurotransmitter synthesized from the amino acid tryptophan		
somatodendritic	region of a neuron that includes the cell body and dendrite(s), but excludes the axon.		
transcription	enzymatic process whereby the genetic information contained in DNA is used to specify a messenger RNA molecule		
translation	process of converting genetic information contained in a mRNA molecule into a protein taking place at the ribosome		
transmitter	small organic compound released from presynaptic terminus of an axon in response to an action potential		
trehalose	is a natural disaccharide formed by an glucoside bond between two α -glucose units		
tyramine	invertebrate neurotransmitter synthesized from the amino acid tyrosine		
Wnt	abbreviation derived from wingless-gene and Int-1 proto oncogene which both are involved in embryonic development		

B7 Amyloid aggregation

Birgit Strodel Structural Biochemistry Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	1 Introduction		2
2	2 Amyloid oligomers		3
3 Functional amyloids			5
4	4 Modeling amyloid aggregation		8
	4.1 Molecular dynamics simulations of protein aggregation		9
	4.2 Enhancing the sampling of amyloid aggregation simulations		9
	4.3 Kinetics of protein aggregation determined by atomistic simulation	ns	10
5	5 Outlook		11

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Amyloids are fibrils made by proteins in a β -strand conformation that lay perpendicular to the fibril axis [1, 2]. They are normally formed by two or more β -sheets which aggregate laterally to form the final fibrils. This structure is called cross- β because of its characteristic X-ray diffraction pattern with a reflection at 4.8 Å corresponding to the distance between β -strands in a sheet, and a second reflection at ~10 Å corresponding to the distance between β -sheets. The essential characteristics of amyloids are summarized in figure 1. It has been hypothesized that most proteins can form amyloids in the right environment [3], but most proteins in physiological conditions have a high β -sheet propensity and their monomers are usually disordered.

Amyloids have been extensively studied because of their association with various diseases. In particular, the deposits found in the brains of patients with Alzheimer's or Parkinson's disease are composed of amyloids. A list of some of the diseases related to amyloids can be seen in table 1. Even though amyloids are related to many diseases, it is still unclear what role they play in these diseases. Notably, it has been observed that the concentration of fibrils does not correlate with the stage of the disease [5]. However, the concentration of soluble oligomers, made up of only a few peptides, have a stronger correlation with the stage of the disease. This



Fig. 1: (A) Amyloids have a cross- β diffraction pattern with a reflection at 4.8 Å caused by the interstrand spacing and a second diffraction around 10 Å caused by the intersheet spacing. In (B) and (C), the structure of GNNQQNY microcrystals from Sup35, the first amyloid structure resolved, is shown [4]. Amyloids are stable because of the hydrogen bonds between the backbone atoms of the peptides. This figure is reproduced with permission of Elsevier.

Annulus

Protein	Disease	Reference	
β -amyloid	Alzheimer's disease	Benilova et al. [7]	
α -synuclein	Parkinson's disease	Irwin et al. [8]	
IAPP (amylin)	Diabetes mellitus type 2	Westermack et al. [9]	
Huntingtin	Huntington's disease	Perutz [10]	
Prion protein PrPSc	Transmissible spongiform encephalopathy	Aguzzi and Calella [11]	
"Natively unstructured" Partially structured monomer			

 Table 1: Disease-related amyloids.



Paranucleus

has led to the hypothesis that it is not amyloid fibrils but soluble oligomers, which are the toxic species [6]. It has also provoked the extensive study of amyloid aggregation, which was found to be a very complex process with many intermediate species as depicted in figure 2 for the aggregation mechanism proposed for the amyloid- β peptide (A β) associated with the development of Alzheimer's disease.

2 Amyloid oligomers

Dodecame

Amylospheroid

The aggregation of proteins into amyloids must start with the initial formation of small nuclei, which later expand into fibrils. Because of their transient nature, these initial aggregates are hard to study experimentally and many open questions remain about their formation [13–15]. In the case of Alzheimer's disease, extensive experimental evidence has been generated during the past two decades, which implicates $A\beta$ oligomers in the development of the disease rather than the aggregation end stage, extremely stable and quite inert amyloid fibrils [14, 16–20]. The results of the various studies suggest that many oligomers different in terms of sizes and structures can be formed, as shown for $A\beta$ in figure 3. More details about amyloid oligomers can be found in a recent review by Strodel and co-workers [19].

Oligomerization seems to be a crossing point between the formation of toxic oligomers, and an intermediate or 'seed' on the pathway to the formation of relatively inert fibrils. Moreover, *in vitro* studies have shown that $A\beta$ can assemble into β -sheet-rich oligomers that are on the path towards the assembly of ordered amyloid fibrils [30]. It had been previously shown that the β -sheet secondary structure is the requisite conformation for soluble $A\beta$ to form mature

10 nm

Fibril

Protofibril



Fig. 3: Different $A\beta$ species ranging from dimers to fibrils are presented. Dimers are flexible and disordered involving β -sheets at different positions and varying degree as shown in (A) to (C) for dimers obtained from simulations in [21], [22] and [23], respectively. The oligomers become more ordered as the oligomer size increases, as shown here for two hexameric structures obtained from simulations (D) [24] and from solid-state NMR (E) [25]. In both models the individual peptides are composed of three β -strands connected by turns, which is similar to the structure obtained from solid-state NMR for disc-shaped pentamers (F) [26]. The dimer model in (G) was found to be the building block of $A\beta$ (pre)globulomers [27]. Here, the C-terminal β -strands already form an in-register parallel β -sheet as observed for the $A\beta$ fibril with PDB ID 2BEG in (H) [28]. Recently, also an S-shaped fibril structure was found for $A\beta$ fibrils (PDB ID 2MXU), which is shown in (I) [29] In some panels not all of the N-terminal residues are shown as they are mostly disordered. The peptides are shown in Cartoon and the names of the first, last and structure-characterizing residues are denoted. This figure has been reproduced from [19].
fibrils [31]. It was also shown that disc-shaped, neurotoxic, pentameric oligomers prepared at 25 °C do not have the β -sheet structure characteristic of fibrils (see figure 3F). However when the temperature was increased to 37 °C the structures converted to a β -sheet-rich state and the oligomers demonstrated a significantly decreased toxicity [26]. It becomes apparent that oligomers can be characterized by their fibrillogenic propensities, which can be distinguished by antibodies. The prefibrillar, annular oligomers can be recognized by the A11 polyclonal antibody, whereas the O4 antibody recognizes the more fibrillar oligomer [32]. Neuronal lipids have been reported to stabilize the fibril and cause it to release toxic oligomers that are biophysically and biochemically similar to, and of similar toxicity as the oligomers used to form the fibrils [33]. Moreover, Knowles and co-workers have found that once a small but critical concentration of amyloid fibrils has accumulated, the toxic oligometric species are predominantly formed from monomeric peptide molecules through a fibril-catalyzed secondary nucleation reaction, rather than through a classical mechanism of homogeneous primary nucleation [34]. Amyloid oligomers have demonstrated toxicity through receptor-mediated, membrane, and intracellular mechanisms. Oligometric A β has been shown to cause cell death by interacting with the N-methyl-D-aspartate receptor (NMDAR), by inducing nerve growth factor (NGF), the p75 neurotrophin (p75NTR), part of the tumor necrosis factor superfamily by receptor-mediated mechanisms. It has been suggested that $A\beta$ oligomers also interact with the cell membrane and cause cell-death by disrupting Ca²⁺ homeostasis, in what has been called the channel hypothesis. A β oligomers have also been shown to move into the cell through the scavenger receptor for advanced glycation end products in neurons and microglia. Moreover, the formation of amyloid oligomers has been associated with the production of free radicals, thus it has been proposed that these polypeptides demonstrate toxicity via oxidative stress [35]. Another important determinant for the cytotoxicity of amyloid oligomers is the surface hydrophobicity as for different amyloid proteins or peptides it has been shown that the toxicity of their oligomers increases with increasing hydrophobicity [36–38]. Moreover, it seems that smaller amyloid oligomers are usually more toxic than larger oligomers [12, 37, 39, 40].

3 Functional amyloids

Most amyloids are associated to diseases. However, the last few years have seen the discovery of amyloids which have normal physiological roles. These amyloids are called *functional amyloids*, and their discovery has changed our understanding of the amyloid fold. Today, we know that functional amyloids can be found in almost every kingdom (see table 2). The functions of these amyloids are diverse and include storage of peptides, scaffold for melanin formation, adhesives, or biofilms.

The first functional amyloid was found in the fungus *Podospora anserina* in 1997 and is made of the HET-s protein [41]. *Podospora anserina* is a filamentous fungus which can have heterokaryon cells, i.e., cells with two different nuclei. Such condition occurs when two cells undergo vegetative cell fusion, which can happen within the same individual but also between different individuals. However, to limit infections, such cell fusion must be tightly controlled. In particular, cells must check if the two nuclei are incompatible. To check such incompatibility, fungi have certain loci to control the viability of the heterokaryon cell. One such loci is the HET-s protein. If two different antagonistic alleles are expressed in the same cytoplasm, the cell is killed. Different alleles code for proteins which differ in 14 amino acids, but only one difference is enough to kill the cell. Coustou et al. [41] showed that Het-s can form heterooligomers

Protein	Organism	Function	Reference	
HET-s	Podospora anserina	Heterokaryon	Costou et al. [41]	
	(fungus)	incompatibility		
Curli	Escherichia coli	Extracellular adhesion,	Chapman et al. [42]	
	(bacterium)	invasion and		
		biofilm formation		
Pmel17	Homo sapiens	Melanin scaffold	Fowler et al. [43]	
Peptides hormones	Mammalia and	Storage and release	Maji et al. [44]	
	Amphibia	of hormones		
RIP1/RIP3	Homo sapiens	Signalling complex for	Li et al. [45]	
		programmed necrosis		
CPEB	Aplysia (mollusca)	Long-term memory	Si et al. [46]	
۸	1	D	J	

Table 2: Functional amyloids in biology.



Fig. 4: *Structure of the Het-s amyloid fibrils from (A) side-view and (B) top-view [47]. Each monomer forms two windings and is represented by a different color.*

which have an amyloid structure. These amyloids formed by proteins from different alleles cause the death of the cell. Wasmer et al. [47] determined the structure for the section 218–289 of the Het-s protein (see figure 4). As other amyloids, Het-s has a very high β -sheet content as it forms a β -solenoid, where each protein has two windings. This organization is in concordance with the fact that HET-s is not polymorphic as many aberrant amyloids are. Furthermore, it exemplifies the fact that HET-s has evolved for its role, in opposition to aberrant amyloids which are caused by misfolded proteins.

Functional amyloids were also discovered in bacteria such as *Escherichia coli* [42]. *Escherichia coli* produces different proteinaceous filaments which are used extracellularly to promote colonization, enter into host cells and to organize communities such as biofilms. One of these filaments used to colonize inert surfaces are composed of fibrils called curli. These fibrils have all the properties of amyloid fibrils. Curli fibrils also have a complex network of protein interactions for their formation. In particular, the protein CsgB is needed to nucleate the aggregation of the protein CsgA, which is the main component of the curli fibrils. Other proteins such as



Fig. 5: *Pmel17 forms amyloids inside melanosomes. Later, melanin precursors are activated through tyrosinases. Pmel17 amyloids template the aggregation of the activated melanin precursors [43]. This figure is reproduced from [43].*

CsgE, CsgF and CsgG are also involved. Such complex mechanism also shows the influence of evolution in the mechanism of functional amyloid formation.

The first functional amyloids found in mammals was Pmel17 [43]. Pmel17 is a protein that is essential for the production of melanosomes. Melanosomes are large organelles which store melanin, a pigment which is found in the skin and is responsible for the protection from UV rays. Melanin is formed by the aggregation of melanin precursors, which first must be activated by tyrosinases. Pmel17 forms amyloid fibrils which serve as a scaffold for such aggregation. A sketch of the process can be seen in figure 5. Considering that melanin precursors are toxic, melanin formation must be strictly regulated and Pmel17 plays a role in this regulation mechanism. Furthermore, amyloids may even be toxic themselves and they also need to be highly regulated. In fact, Pmel17 is trafficked to melanosomes as a transmembrane protein that cannot aggregate. Later, the M α section of the protein, which is the section that aggregates, is secreted through proteolysis. This section then aggregates and forms amyloids. This process is very similar to the pathway in which Alzheimer's β -amyloid is created. One of the most interesting characteristics of Pmel17 aggregation is that it is much faster than the aggregation of toxic amyloid such as β -amyloid or α -synuclein. This is one of the possible mechanisms for amyloids to avoid being in the oligomer state and diminish their toxicity.

Furthermore, a number of peptide hormones have been found to be stored as functional amyloids [44]. These hormones are stored in the Golgi apparatus and later transported in secretory granules to the extracellular space where these amyloid disaggregate into monomers, which are the active species. This process is illustrated in figure 6. In this case, the amyloid fold is an excellent fold for such a role as hormones can be stored in a high concentration. Storing peptides in such a manner will help for a rapid secretion, as the synthesis of the peptide would take much longer. The amyloid fold can also explain how secretory granules are composed of single peptide species, as amyloid formation is sequence specific. Maji et al. [44] observed that the hormones aggregate only under certain environmental conditions, such as at a specific pH or in the presence of helper molecules such as glycosaminoglycans (GAGs). Such delicate environmental conditions could be relevant to how these functional amyloids avoid being toxic. The last functional amyloid discovered in humans was a heterooligomeric amyloid signalling

complex formed by RIP1 and RIP3 kinases [45]. These kinases are required for programmed necrosis and function in a feed forward mechanism where kinase activation and RIP1/RIP3



Fig. 6: Sketch of how hormone peptides are stored as functional amyloids. Hormones are stored as amyloids in the Golgi apparatus. Later, they are transported in secretory granules to the extracellular space. When outside the cell, the amyloids dissolve and hormone monomers become active [44]. This figure is reproduced from the Supporting Material of [44].

aggregation reinforce each other. When RIP1 and RIP3 are not phosphorylated they do not aggregate but when phosphorylated they form amyloids. Interestingly, when aberrant amyloids, such as τ -protein and α -synuclein, are phosphorylated they aggregate faster too. As these kinases are signalling for programmed necrosis, it has been hypothesised that these amyloids are the ones that kill the cell.

Finally, one of the most fascinating functionalities in which amyloids may play a role is longterm memory [46,48]. Long-term memory, unlike short-term memory, needs the creation of new synapse pathways between neurons. The creation of new synapses is mediated by serotonin, a monoamine neurotransmitter. However, these changes must be made permanent and how this happens has been challenging to understand. Considering that amyloid fibrils are much more stable than most proteins, it has been suggested that they play a role in long-term memory. In particular, Si et al. [46, 48] showed that the cytoplasmic polyadenylation element binding protein (CPEB), which is related to long-term memory in *Aplysia* (a sea slug usually studied by neurobiologists because it has only around 20,000 neurons), forms amyloids. In particular, CPEB appears to be in two different states: a first one which is not active and does not aggregate, and a second one in which it is active and aggregates into a prion-like amyloid. In the prion-like state, it stimulates other non-active CPEB to become active and aggregate, and later bind to RNA. Serotonin increases the production of CPEB. There are still many open questions in how amyloids relate to long-term memory. Furthermore, it would be interesting to observe how this translates into the human brain, which is much more complex than the one from *Aplysia*.

4 Modeling amyloid aggregation

One of the grand challenges of biophysics and biochemistry is to understand the principles that govern protein aggregation. Protein aggregation is a highly complex process that is sensitive to initial conditions, operates on a huge range of timescales and its products range from dimers to macroscopic fibrils. Molecular simulations can help understand protein aggregation. In particu-

lar, atomistic simulations can be used to study the initial formation of toxic oligomers which are hard to characterize experimentally, and to understand the difference in aggregation behaviour between different amyloidogenic proteins. Here, we review the latest atomistic simulations of protein aggregation.

4.1 Molecular dynamics simulations of protein aggregation

Molecular simulations have helped tremendously in understanding the underlying physics and chemistry behind protein dynamics [49]. The standard method for simulating proteins is molecular dynamics (MD) simulations , in which the positions and velocities of atoms are calculated using classical mechanics. Proteins and solvent molecules are usually represented atomistically. This detailed representation leads to extremely expensive simulations. Therefore, for processes such as protein aggregation, coarse-grained models, in which many atoms are simulated as a single bead, have been extensively used [50–54]. However, with the improved parallelization of widely used MD codes [55, 56], the advent of special purpose parallel architectures [57] and the implementation of various MD codes into graphical processing units (GPUs) [58, 59], simulations which were impossible a few years ago have become a reality. Thus, it is now possible to study the aggregation of amyloidogenic peptides by means of explicit solvent atomistic simulations.

The most straightforward way of studying protein aggregation is by placing a number of peptides in a solvated box and allowing them to aggregate. The first study of protein aggregation using an explicit water all-atom force field was performed by Klimov and Thirumalai [60]. They studied the aggregation of a fragment of Alzheimer's A β , A β_{16-22} and used a harmonic constraint between the center of the water box and the oligomer center of mass in order to speed up the aggregation process. They observed that $A\beta_{16-22}$ forms antiparallel β -sheets, after passing through an α -helical intermediate. Later, many similar studies were performed in which different amyloidogenic peptides were allowed to aggregate without constraints. In most of these studies [61-64], it is observed that the peptides first collapse to form only partially ordered aggregates and later, these aggregate evolve into ordered β -sheet structures. A particularly interesting study is the one by Matthes et al. [62] in which they investigated the driving forces behind protein aggregation for various peptides. They observed a two-step mechanism in which the peptides first aggregate into partially ordered aggregates driven by solvation free energy, and later these aggregates reorder into β -sheets driven by the optimization of interpeptide interactions. Strodel and co-workers have recently developed a method using transition networks to elucidate the aggregation pathways sampled during MD simulations [63]. This approach revealed for A β that the oligomers leading to the type of oligomer distributions observed in experiments originate from compact conformations. Extended oligomers, on the other hand, contribute more to the production of larger aggregates thus driving the aggregation process [65].

4.2 Enhancing the sampling of amyloid aggregation simulations

Even though it is now possible to simulate larger systems for longer times, protein aggregation simulations are hard to converge. Hence, enhanced sampling techniques are normally used to accelerate the convergence of these simulations. One of the most commonly used enhanced sampling algorithms is replica exchange molecular dynamics (REMD) [66]. In this algorithm, one runs multiple replicas of the system under study at different temperatures. Then, based on the Metropolis criterion, one periodically exchanges configurations between replicas at different

temperatures. High temperature replicas are used to enhance the sampling of the free energy surface, whereas low temperature replicas are used to sample the system at a relevant temperature. REMD simulations have been extensively used to study protein aggregation [21, 23, 67–70]. For example, REMD was used to study the aggregation of 16 Alzheimer's $A\beta_{37-42}$ peptides by Nguyen and Derreumaux [70]. They observed that the global free energy minimum is characterized by 2- or 3-stranded β -sheets. However, they also observed a non-negligible amount of 5- and 6-stranded antiparallel sheets. One of the problems of REMD simulations is that, because of the exchanges between replicas, kinetic information is lost and, hence, it is hard to have a detailed understanding of the dynamics of the aggregation process. Also, because the number of replicas depend on the number of atoms of the system, REMD can easily become expensive. The latter problem can be alleviated by using the Hamiltonian REMD (HREMD) method, in which the different replicas have different Hamiltonians but are simulated (in most cases) at a constant temperature [71-73]. The idea behind this approach is that the system is trapped in a local minimum because of strong bonded and non-bonded interactions. Thus, if these interactions are weakened by applying a scaling factor, the energy landscape is sampled more efficiently. Such scaling factor is chosen in a way so that the number of replicas does not depend on the number of solvent molecules. In HREMD different scaling factors are used for the different replicas, including one replica with the unmodified Hamiltonian. Compared to temperature REMD, fewer replicas are required if only the Hamiltonian of the aggregating peptides but not the solvent is modified. Laghaei et al. [74] used the combination of temperature and Hamiltonian REMD to characterize the structure and thermodynamics of the full-length hIAPP dimer.

Another algorithm used to enhance sampling is metadynamics [75]. In this method, a history dependant bias acts in a number of collective variables on the system. The system is then forced to leave the areas of the free energy landscape that have already been visited and, hence, the sampling is enhanced. In the case of studying protein aggregation, bias-exchange metadynamics [76], in which replica exchange and metadynamics are combined, was used [77–80]. Using this methodology, Baftizadeh et al. [77], for example, showed that the aggregation of polyvaline starts with the formation of antiparallel β -sheets. When enough antiparallel sheets are formed, parallel β -sheets begin to appear. The system then falls into a free energy minimum mostly formed by parallel β -sheets.

4.3 Kinetics of protein aggregation determined by atomistic simulations

A number of methods are available for studying rare events which can accurately estimate the kinetics of the processes in question [81], such as transition path sampling (TPS) [82] and forward flux sampling (FFS) [83]. Even though they have not yet been used to study oligomer formation using explicit-solvent atomistic simulations, they have been applied to similar systems. In particular, Schor et al. [84] studied the mechanism of monomer addition to an amyloid fibril of a fragment of the insulin peptide hormone using TPS. TPS calculates the equilibrium path ensemble between two predefined states using a Monte Carlo random walk in trajectory space. Using TPS, Schor et al. [84] proposed a detailed mechanism of the docking of monomers to a growing fibril. In a different study, Luiken and Bolhuis [85] applied FFS to elucidate the aggregation of different amyloidogenic peptides into amyloids represented with a coarse-grained model. In FFS, the initial and final state are separated in terms of an order parameter λ . FFS starts with a normal MD simulation, in which configurations at the initial state are generated. Trial runs are generated from these configurations that reach the next interface defined by λ

or return to the initial state. The algorithm then generates configurations from the following interface and the procedure is repeated until the final state is reached. Luiken and Bolhuis [85] used order parameters dependant on the number of in-register contacts between peptides to study amyloid formation. They found that the nucleation pathway changes from a one-step to a two-step nucleation mechanism with increasing hydrophobicity.

In the last years, it has become popular to construct kinetic models of long-time protein dynamics from multiple MD simulations using so-called Markov state models (MSMs) [86, 87]. MSMs not only use the information of multiple short simulations but also allow to sample the energy landscape more efficiently by adaptive sampling [88]. Kelley et al. [89] used MSMs to analyze atomistic MD simulations of the aggregation of Alzheimer's $A\beta_{21-43}$ into small oligomers (up to tetramers). By analytically considering the diffusion of peptides, they could estimate the aggregation of these peptides at *in vitro* concentrations. Even though this was an important step forward, they did not consider the difference in conformational sampling by oligomers at different peptide concentrations. Perkett and Hagan [90] used MSMs to study the aggregation of virus proteins using a coarse-grained model. They created a theoretical framework in which aggregates are considered as undirected graphs which are then used to construct the MSM. However, neither of these two MSM studies considered the intramolecular dynamics that drive aggregation. Finally, Schor et al. [91] used MSMs to understand the mechanism of monomer addition of transthyretin $TTR_{105-115}$ to a growing $TTR_{105-115}$ fibril. In this case, the intramolecular dynamics was considered and a detailed mechanistic model of the process of monomer addition established.

5 Outlook

The simulation of protein aggregation promoted our understanding of the formation of oligomers by amyloidogenic peptides. Initially, most of these simulations have used coarse-grained models because explicit-solvent atomistic simulations were too expensive to perform. However, with the advance of new technologies and software, atomistic simulations of the aggregation of small amyloidogenic peptides are now possible. These simulations are normally performed at mM concentrations, while the concentration of amyloidogenic peptides *in vivo* is in the pM to nM range and *in vitro* studies are performed at μ M concentrations. Since high peptide concentrations do not give the newly formed oligomers enough time to equilibrate before the attachment of another monomer or oligomer, new approaches or methodologies should be developed in order to enable simulations at much lower concentrations. Brute force simulations at such low concentrations are currently not feasible. Considering that a simulation at mM concentration requires thousands up to one million atoms, the same simulation at nM concentration would require a hundred billion atoms [92], which is orders of magnitude more than the current biggest simulations. Even though eventually such large simulations will probably become possible, most of the simulation time would be wasted on the diffusion of the peptides and the simulation of the solvent. Thus, methods that take into account the low concentration implicitly by means of analytical or multiscale methods [93, 94] are needed. Moreover, in the future we expect more aggregation studies of the entire peptides or even proteins associated with diseases and not only sections of them. These simulations, studying aggregation beyond dimers, are currently only possible using approximations such as an implicit representation of the solvent [64]. Finally, future aggregation simulations should consider the different environments which are essential in modulating amyloid formation, such as membranes [95], transition metal ions [96]

or macromolecular crowding [97]. Such simulations will provide an enhanced understanding of the aggregation of disease-related and functional amyloids, and will help in the design of drugs to fight disease-related amyloids and in the development of novel amyloid-based nanomaterials.

References

- [1] F. Chiti and C. M. Dobson, Annu. Rev. Biochem. 75, 333 (2006).
- [2] J. Greenwald and R. Riek, Structure 18, 1244 (2010).
- [3] L. Goldschmidt, P. K. Teng, R. Riek, and D. Eisenberg, Proc. Natl. Acad. Sci. U.S.A. 107, 3487 (2010).
- [4] D. Eisenberg and M. Jucker, Cell 148, 1188 (2012).
- [5] C. Schmitz, B. P. Rutten, A. Pielen, S. Schäfer, O. Wirths, G. Tremp, C. Czech, V. Blanchard, G. Multhaup, P. Rezaie, H. Korr, H. W. Steinbusch, *et al.*, Am. J. Pathol. **164**, 1495 (2004).
- [6] W. Klein, W. Stine, and D. Teplow, Neurobiol. Aging 25, 569 (2004).
- [7] I. Benilova, E. Karran, and B. De Strooper, Nat. Neurosci. 15, 349 (2012).
- [8] D. J. Irwin, V. M. Y. Lee, and J. Q. Trojanowski, Nat. Rev. Neurosci. 14, 626 (2013).
- [9] P. Westermark, A. Andersson, and G. T. Westermark, Physiol. Rev. 91, 795 (2011).
- [10] M. F. Perutz, Trends Biochem. Sci. 24, 58 (1999).
- [11] A. Aguzzi and A. M. Calella, Physiol. Rev. 89, 1105 (2009).
- [12] D. Teplow, Alzheimers Res. Ther. 5, 39 (2013).
- [13] F. Bemporad and F. Chiti, Chem. Biol. 19, 315 (2012).
- [14] M. Fändrich, J. Mol. Biol. 421, 427 (2012).
- [15] L. Breydo and V. N. Uversky, FEBS Lett. 589, 2640 (2015).
- [16] M. P. Lambert, A. K. Barlow, B. A. Chromy, C. Edwards, R. Freed, M. Liosatos, T. E. Morgan, I. Rozovsky, B. Trommer, K. L. Viola, P. Wals, C. Zhang, *et al.*, Proc. Natl. Acad. Sci. **95**, 6448 (1998).
- [17] T. Huang, D.-S. Yang, N. P. Plaskos, S. Go, C. M. Yip, P. E. Fraser, and A. Chakrabartty, J. Mol. Biol. 297, 73 (2000).
- [18] M. Kirkitadze, G. Bitan, and D. Teplow, J. Neurosci. Res. 69, 567 (2002).
- [19] L. Nagel-Steger, M. C. Owen, and B. Strodel, ChemBioChem 17, 657 (2016).
- [20] S. J. C. Lee, E. Nam, H. J. Lee, M. G. Savelieff, and M. H. Lim, Chem. Soc. Rev. 46, 310 (2017).
- [21] A. Yano, A. Okamoto, K. Nomura, S. Higai, and N. Kurita, Chem. Phys. Lett. 595–596, 242 (2014).
- [22] T. Zhang, J. Zhang, P. Derreumaux, and Y. Mu, J. Phys. Chem. B 117, 3993 (2013).
- [23] B. Tarus, T. T. Tran, J. Nasica-Labouze, F. Sterpone, P. H. Nguyen, and P. Derreumaux, J. Phys. Chem. B 119, 10478 (2015).
- [24] B. Strodel, J. W. L. Lee, C. S. Whittleston, and D. J. Wales, J. Am. Chem. Soc. 132, 13300 (2010).
- [25] C. Lendel, M. Bjerring, A. Dubnovitsky, R. T. Kelly, A. Filippov, O. N. Antzutkin, N. C. Nielsen, and T. Härd, Angew. Chem. Int. Ed. 53, 12756 (2014).
- [26] M. Ahmed, J. Davis, D. Aucoin, T. Sato, S. Ahuja, S. Aimoto, J. I. Elliott, W. E. Van Nostrand, and S. O. Smith, Nat. Struct. Mol. Biol. 17, 561 (2010).
- [27] L. Yu, R. Edalji, J. E. Harlan, T. F. Holzman, A. P. Lopez, B. Labkovsky, H. Hillen,

S. Barghorn, U. Ebert, P. L. Richardson, L. S. L. Miesbauer, D. Bartley, *et al.*, Biochemistry **2009**, 1870 (2009).

- [28] T. Luhrs, C. Ritter, M. Adrian, D. Riek-Loher, B. Bohrmann, H. Dobeli, D. Schubert, and R. Riek, Proc. Natl. Acad. Sci. U.S.A. 102, 17342 (2005).
- [29] Y. Xiao, B. Ma, D. McElheny, S. Parthasarathy, F. Long, M. Hoshi, R. Nussinov, and Y. Ishii, Nat. Struct. Mol. Biol. 22, 499 (2015).
- [30] I. A. Mastrangelo, M. Ahmed, T. Sato, W. Liu, C. Wang, P. Hough, and S. O. Smith, J. Mol. Biol. 358, 106 (2006).
- [31] C. Soto, E. M. Castano, R. A. Kumar, R. C. Beavis, and B. Frangione, Neurosci. Lett. 200, 105 (1995).
- [32] R. Kayed, I. Canto, L. Breydo, S. Rasool, T. Lukacsovich, J. Wu, R. Albay, 3rd, A. Pensalfini, S. Yeung, E. Head, J. L. Marsh, and C. Glabe, Mol. Neurodegener. 5, 57 (2010).
- [33] I. C. Martins, I. Kuperstein, H. Wilkinson, E. Maes, M. Vanbrabant, W. Jonckheere, P. Van Gelder, D. Hartmann, R. D'Hooge, B. De Strooper, J. Schymkowitz, and F. Rousseau, EMBO J. 27, 224 (2008).
- [34] S. I. A. Cohen, S. Linse, L. M. Luheshi, E. Hellstrand, D. A. White, L. Rajah, D. E. Otzen, M. Vendruscolo, C. M. Dobson, and T. P. J. Knowles, Proc. Natl. Acad. Sci. p. 201218402 (2013).
- [35] B. J. Tabner, O. M. El-Agnaf, M. J. German, N. J. Fullwood, and D. Allsop, Biochem. Soc. Trans. 33, 1082 (2005).
- [36] S. Campioni, B. Mannini, M. Zampagni, A. Pensalfini, C. Parrini, E. Evangelisti, A. Relini, M. Stefani, C. M. Dobson, C. Cecchi, and F. Chiti, Nat. Chem. Biol. 6, 140 (2010).
- [37] B. Mannini, E. Mulvihill, C. Sgromo, R. Cascella, R. Khodarahmi, M. Ramazzotti, C. M. Dobson, C. Cecchi, and F. Chiti, ACS Chem. Biol. 9, 2309 (2014).
- [38] A. R. A. Ladiwala, J. Litt, R. S. Kane, D. S. Aucoin, S. O. Smith, S. Ranjan, J. Davis, W. E. V. Nostrand, and P. M. Tessier, J. Biol. Chem. 287, 24765 (2012).
- [39] K. Ono, M. M. Condron, and D. B. Teplow, Proc. Natl. Acad. Sci. 106, 14745 (2009).
- [40] E. Hayden and D. Teplow, Alzheimers Res. Ther. 5, 60 (2013).
- [41] V. Coustou, C. Deleu, S. Saupe, and J. Begueret, Proc. Natl. Acad. Sci. U.S.A. 94, 9773 (1997).
- [42] M. R. Chapman, L. S. Robinson, J. S. Pinkner, R. Roth, J. Heuser, M. Hammar, S. Normark, and S. J. Hultgren, Science 295, 851 (2002).
- [43] D. M. Fowler, A. V. Koulov, C. Alory-Jost, M. S. Marks, W. E. Balch, and J. W. Kelly, PLoS Biol. 4, e6 (2005).
- [44] S. K. Maji, M. H. Perrin, M. R. Sawaya, S. Jessberger, K. Vadodaria, R. A. Rissman, P. S. Singru, K. P. R. Nilsson, R. Simon, D. Schubert, D. Eisenberg, J. Rivier, *et al.*, Science 325, 328 (2009).
- [45] J. Li, T. McQuade, A. B. Siemer, J. Napetschnig, K. Moriwaki, Y.-S. Hsiao, E. Damko, D. Moquin, T. Walz, A. McDermott, F. K.-M. Chan, and H. Wu, Cell 150, 339 (2012).
- [46] K. Si, S. Lindquist, and E. R. Kandel, Cell 115, 879 (2003).
- [47] C. Wasmer, A. Lange, H. Van Melckebeke, A. B. Siemer, R. Riek, and B. H. Meier, Science 319, 1523 (2008).
- [48] K. Si, Y.-B. Choi, E. White-Grindley, A. Majumdar, and E. R. Kandel, Cell 140, 421 (2010).
- [49] R. O. Dror, R. M. Dirks, J. Grossman, H. Xu, and D. E. Shaw, Annu. Rev. Biophys. 41, 429 (2012).
- [50] N. Mousseau, , and P. Derreumaux, Acc. Chem. Res. 38, 885 (2005).

- [51] J. E. Straub and D. Thirumalai, Curr. Opin. Struct. Biol. 20, 187 (2010).
- [52] C. Wu and J.-E. Shea, Curr. Opin. Struct. Biol. 21, 209 (2011).
- [53] A. Morriss-Andrews and J.-E. Shea, J. Phys. Chem. Lett. 5, 1899 (2014).
- [54] A. Morriss-Andrews and J.-E. Shea, Annu. Rev. Phys. Chem. 66, 643 (2015).
- [55] B. Hess, C. Kutzner, D. van der Spoel, and E. Lindahl, J. Chem. Theory Comput. 4, 435 (2008).
- [56] J. C. Phillips, R. Braun, W. Wang, J. Gumbart, E. Tajkhorshid, E. Villa, C. Chipot, R. D. Skeel, L. Kalé, and K. Schulten, J. Comput. Chem. 26, 1781 (2005).
- [57] D. E. Shaw, M. M. Deneroff, R. O. Dror, J. S. Kuskin, R. H. Larson, J. K. Salmon, C. Young, B. Batson, K. J. Bowers, J. C. Chao, M. P. Eastwood, J. Gagliardo, *et al.*, Commun. ACM **51**, 91 (2008).
- [58] M. S. Friedrichs, P. Eastman, V. Vaidyanathan, M. Houston, S. Legrand, A. L. Beberg, D. L. Ensign, C. M. Bruns, and V. S. Pande, J. Comput. Chem. 30, 864 (2009).
- [59] M. J. Harvey, G. Giupponi, and G. D. Fabritiis, J. Chem. Theory Comput. 5, 1632 (2009).
- [60] D. K. Klimov and D. Thirumalai, Structure 11, 295 (2003).
- [61] D. Matthes, V. Gapsys, V. Daebel, and B. L. de Groot, PLoS ONE 6, e19129 (2011).
- [62] D. Matthes, V. Gapsys, and B. L. de Groot, J. Mol. Biol. 421, 390 (2012).
- [63] B. Barz, D. J. Wales, and B. Strodel, J. Phys. Chem. B 118, 1003 (2014).
- [64] B. Barz, O. O. Olubiyi, and B. Strodel, Chem. Commun. 50, 5373 (2014).
- [65] B.Barz, Q. Liao, and B. Strodel, J. Am. Chem. Soc. DOI: 10.1021/jacs.7b10343 (2018).
- [66] Y. Sugita and Y. Okamoto, Chem. Phys. Lett. **314**, 141 (1999).
- [67] G. Bellesia and J.-E. Shea, Biophys. J. 96, 875 (2009).
- [68] E. Rivera, J. Straub, and D. Thirumalai, Biophys. J. 96, 4552 (2009).
- [69] L. Larini and J.-E. Shea, Biophys. J. 103, 576 (2012).
- [70] P. H. Nguyen and P. Derreumaux, J. Phys. Chem. B 117, 5831 (2013).
- [71] S. Jang, S. Shin, and Y. Pak, Phys. Rev. Lett. 91, 058305 (2003).
- [72] P. Liu, B. Kim, R. A. Friesner, and B. J. Berne, Proc. Natl. Acad. Sci. U.S.A. 102, 13749 (2005).
- [73] L. Wang, R. A. Friesner, and B. J. Berne, J. Phys. Chem. B 115, 9431 (2011).
- [74] R. Laghaei, N. Mousseau, and G. Wei, J. Phys. Chem. B 115, 3146 (2011).
- [75] A. Laio and M. Parrinello, Proc. Natl. Acad. Sci. U.S.A. 99, 12562 (2002).
- [76] S. Piana and A. Laio, J. Phys. Chem. B 111, 4553 (2007).
- [77] F. Baftizadeh, X. Biarnés, F. Pietrucci, F. Affinito, and A. Laio, J. Am. Chem. Soc. 134, 3886 (2012).
- [78] F. Baftizadeh, F. Pietrucci, X. Biarnés, and A. Laio, Phys. Rev. Lett. 110, 168103 (2013).
- [79] A. Barducci, M. Bonomi, M. K. Prakash, and M. Parrinello, Proc. Natl. Acad. Sci. U.S.A. 110, E4708 (2013).
- [80] L. E. Buchanan, E. B. Dunkelberger, H. Q. Tran, P.-N. Cheng, C.-C. Chiu, P. Cao, D. P. Raleigh, J. J. de Pablo, J. S. Nowick, and M. T. Zanni, Proc. Natl. Acad. Sci. U.S.A. 110, 19285 (2013).
- [81] K. Klenin, B. Strodel, D. J. Wales, and W. Wenzel, BBA-Proteins Proteom 1814, 977 (2011).
- [82] P. G. Bolhuis, D. Chandler, C. Dellago, and P. L. Geissler, Annu. Rev. Phys. Chem. 53, 291 (2002).
- [83] R. J. Allen, C. Valeriani, and P. R. ten Wolde, J. Phys.: Condens. Matter 21, 463102 (2009).
- [84] M. Schor, J. Vreede, and P. G. Bolhuis, Biophys. J. 103, 1296 (2012).

- [85] J. A. Luiken and P. G. Bolhuis, J. Phys. Chem. B 119, 12568 (2015).
- [86] J.-H. Prinz, H. Wu, M. Sarich, B. Keller, M. Senne, M. Held, J. D. Chodera, C. Schütte, and F. Noé, J. Chem. Phys. 134, 174105 (2011).
- [87] G. Bowman, V. S. Pande, and F. Noé (eds.), An Introduction to Markov State Models and Their Application to Long Timescale Molecular Simulation, Advances in Experimental Medicine and Biology (Springer, 2014).
- [88] G. R. Bowman, D. L. Ensign, and V. S. Pande, J. Chem. Theory Comput. 6, 787 (2010).
- [89] N. W. Kelley, V. Vishal, G. A. Krafft, and V. S. Pande, J. Chem. Phys. 129, 214707 (2008).
- [90] M. R. Perkett and M. F. Hagan, J. Chem. Phys. 140, 214101 (2014).
- [91] M. Schor, A. S. J. S. Mey, F. Noé, and C. E. MacPhee, J. Phys. Chem. Lett. 6, 1076 (2015).
- [92] M. Carballo-Pacheco and B. Strodel, J. Phys. Chem. B 120, 2991 (2016).
- [93] L. W. Votapka and R. E. Amaro, PLoS Comput. Biol. 11, e1004381 (2015).
- [94] S. P. Hirakis, B. W. Boras, L. W. Votapka, R. D. Malmstrom, A. D. McCulloch, and R. E. Amaro, Front. Physiol. 6 (2015).
- [95] M. Bucciantini, S. Rigacci, and M. Stefani, J. Phys. Chem. Lett. 5, 517 (2014).
- [96] J. Nasica-Labouze, P. H. Nguyen, F. Sterpone, O. Berthoumieu, N.-V. Buchete, S. Côté, A. D. Simone, A. J. Doig, P. Faller, A. Garcia, A. Laio, M. S. Li, *et al.*, Chem. Rev. **115**, 3518 (2015).
- [97] R. Ellis and A. Minton, Biol. Chem. 387, 485 (2006).

B8 Optogenetics

V. Gordeliy^{1,2,3}, I. Okhrimenko³ and K. Kovalev^{1,2,3} ¹ Structural Biochemistry, Institute of Complex Systems (ICS-6), Research Centre Juelich, Juelich, Germany ² Institut de Biologie Structurale J.-P. Ebel, Université Grenoble Alpes-CEA-CNRS, Grenoble, France ³ Moscow Institute of Physics and Technology, Dolgoprudniy, Russia

Contents

1	Introd	Introduction	
2	Development of the method		2
	2.1	History of rhodopsins	2
	2.2	History of the optogenetics	3
3	Applie	cations of optogenetic tools	3
	3.1	Rhodopsins as optogenetics tools	3
	3.2	Characterization and functional studies of optogenetic tools	4
	3.3	Structure-based rational design of optogenetic tools	5
4	Persp	ectives and unsolved problems	8
R	eferenc	es	10

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Rhodopsins are probably the most universal biological light-energy transducers and abundant phototrophic mechanisms evolved on Earth. They are found in all the domains of life and in viruses and have a remarkable diversity and potential for biotechnological applications. Channel rhodopsins, H+ and Cl- pumps have become indispensable means of optogenetics and revolutionized neuroscience promising new approaches to the treatment of severe diseases.

However, among 700 identified rhodopsin genes only a few are characterized. This dramatically limits our knowledge of their functions, mechanisms and biotechnological applications. Moreover, high-resolution structures and molecular mechanisms of recently studied new rhodopsins are either not known or limited to non-active states. The amazing scientific and technological potential of rhodopsins is still to be exploited.

2 Development of the method

2.1 History of rhodopsins

Microbial rhodopsins were first discovered in the early 1970s in the archaeal halophiles (saltloving microbes that live in saturated brines)^{1,2}. It was shown that bacteriorhodopsin (BR), a seven alpha helical membrane protein comprising the retinal, is a proton pump used to generate a light-driven proton gradient as a source of energy for the cell. Soon after bacteriorhodopsin, two sensory rhodopsins (also known as slow rhodopsins) were found in the same microbe (*Halobacterium salinarum*)³.

They are light sensors that produce complex phototactic behaviour in the microbes. Then, again in the same microbe a fourth pump was found that functions as an inward directed chloride pump (halorhodopsin) that helps maintaining the proper electrochemical balance⁴. These rhodopsins played a key role in biology and biophysics, in particular, in membrane protein research, extending our understanding of molecular mechanisms of bioenergetics, transmembrane signaling and in the development of new biophysical, biochemical, structural biology (X-ray crystallography and EM) methods, approaches and potential applications⁵. However, at the end of the XXth century not many of us (if any at all) would imagine a new era of the rhodopsin world and their breaking contributions to neuroscience. Not many

era of the rhodopsin world and their breaking contributions to neuroscience. Not many enthusiasts were ready to stay working with these proteins. But nature really plays tricks on us.

A new microbial rhodopsins era started in 2000 when, first of all thanks to metagenomics, rhodopsins were also found in bacteria. An uncultivated marine gammaproteobacterium of the 16S rRNA defined clade SAR86 was confirmed to have a functional proton pump that was named proteorhodopsin⁶. Then due to the rapid increase in bacterial genomes of the first decade of the XXIst century it appeared that there are unexpectedly many rhodopsins and they are present in all the domains of life and even in viruses. It was also unexpected that these 'seemingly the same' proteins can execute diverse functions. Ground breaking discoveries and characterization of channel rhodopsins in 2002⁷ sodium pump in 2013⁸ rhodopsin-guanylyl cyclase in 2014⁹, outward proton pumps xenorhodopsins in 2016¹⁰/2017¹¹ support this conclusion. It is amazing how evolution was using 'minimum' protein modification to engineer rhodopsins with such different functions.

2.2 History of the optogenetics

Application of channel rhodopsin 2 for light induced firing neurons with unprecedented time and special resolution opened a new field – optogenetics. These proteins together with the proton and chloride light-driven pumps have become core tools of optogenetics. Optogenetics – enabling a real-time control of cellular processes, cells, networks and animal behaviour by using light-activated proteins that affect membrane voltage or intracellular signals – is a novel multidisciplinary science and biotechnology field that integrates electrophysiology, structural biology, optics, electronics and genetic engineering. Optogenetic control of nerve cells represents one of the most important technological advancements in neuroscience to date. Optogenetics promises to solve a challenging problem of cellular control by enabling unprecedented temporal and spatial precision in tissues of living animals, which electrical stimulations and pharmacological methods are unable to achieve. It offers a breakthrough in our ability to study the mechanisms of complex processes in neuronal circuits, such as learning and motor function. It also brings a promise of providing an entirely new approach to the restoration of function in blindness or brain degeneration, and to the treatment of a variety of other neurological and mental disorders. In the "Human Brain Project" with the

funding of \notin 1 bln per year for 10 years, initiated by the EU in 2013, optogenetics is one of the key technologies. Similar to that The Brain Research through Advancing Innovative Neurotechnologies[®] (BRAIN) Initiative in the USA is aimed at revolutionizing our understanding of the human brain and it considers optogenetics a key technology to achieve the goals.

3 Applications of optogenetic tools

3.1 Rhodopsins as optogenetics tools

It is evident that rhodopsins, the core tools of optogenetics, will be among key determinants of success of these global initiatives. However, despite their broad impact and unprecedented possibilities for basic and applied research, to date, only a few retinal proteins are studied and available for optogenetics.

Specifically, the major tools are channel rhodopsin 2 (ChR2) from *Chlamydomonas reinhardtii* (a cation channel), Cl- pump halorhodopsin (NphR) from *Natronomonas pharaonis*, and H+ pump archaerhodopsin 3 (Arch3) from *Halorubrum sodomense*. NphRs are used to silence neurons. The cationic ChRs are the only protein able to activate neurons, but it lacks ion selectivity. Numerous protein engineering efforts to modify its ion selectivity have met limited success, but also, as it is a channel, the actual current of any specific ion conductance imposed on the ChR2 by mutation, depends on the actual state of the membrane, i.e., ionic gradients and electrical potential. In the Table 1 you can find the list of ChR variants relevant for optogenetic applications.

ChR variant	Use(s) and special properties
ChR2 H134R	Widely used ChR, increased Na+ conductivity, improved
	retinal binding in <i>Caenorhabditis elegans</i> ^{12,13}
CatCh (ChR2 L132C)	Enhanced Mg2+ and Ca2+ selectivity, large photocurrents, low inactivation ^{14,15}
ChR2 T159C	Large photocurrents, improved retinal binding affinity ^{16,17}

CatCh+ (ChR2 L132C-T159C)	Enhanced Mg2+ and Ca2+ selectivity, large photocurrents, low
	inactivation, provides high light sensitivity to host cells, more
	stable expression than CatCh ^{15,18,19}
ChETA (ChR2 E123T-T159C)	Fast photocycle at the expense of reduced transported charge
	per absorbed photon, reduced voltage dependency of channel
	closing kinetics ^{16,20}
Chronos	Very fast photocycle, large photocurrents, inverse voltage
	dependency of off kinetics, provides high light sensitivity to
	host cells ²¹
Step-function rhodopsins (ChR2	Very slow photocycle kinetics and extended open-state
C128A/S/T; ChR2 D156A/C/N)	lifetimes, provides very high light sensitivity to cells, UV
	light- and green light-induced channel closure (bistable
	rhodopsins) ^{22,23}
ChIEF	Fast photocycle, low inactivation ¹³
C1V1 (E122T-E162T)	Green light-induced activation, high two-photon cross-
	section ^{15,24}
ReachR	Green light-induced activation, high expression level, and
	large photocurrents in mammalian cells ²⁵
PsChR	Violet light-induced activation, high unitary conductance and
	high Na+ conductivity, low inactivation, fast recovery
	kinetics ²⁰
Cs-Chrimson	Orange light-induced activation, improved membrane
	targeting by use of the CsChR (ChR87) N-terminus ²¹
Slow ChloC	Cl- selectivity, slow photocycle, improved photocurrents and
(ChR2 E90R-D156N T159C)	shifted reversal potential, used for voltage clamping to the Cl-
	reversal potential ²⁷
iC1C2	Cl- selectivity, used for voltage clamping to the Cl- reversal
	potential ²⁰

 Table 1: Overview of channelrhodopsin (ChR) variants relevant for optogenetic applications.

It should mentioned, that pumps, in contrast to ChRs, allow vectorial transport, almost irrespective of the actual gradient of the respective ion. However, apart from H+, Cl-, other ions even more relevant for the neuronal function have not been addressable by microbial rhodopsin pumps so far (Na+, K+, Ca2+). Particularly, an outward and inward K+ pumps would be highly desirable for neuroscience applications, as K+ is the main ion used for neuronal re- or hyper-polarization, and could, thus far, not be directly addressed using known rhodopsins. It was also demonstrated that the xenorhodopsin from Nanosalina (NsXeR) is a powerful pump which is able to elicit action potentials in rat hippocampal neuronal cells up to their maximal intrinsic firing frequency, proving that the inwardly directed proton pumps are suitable for light induced remote control of neurons¹¹.

3.2 Characterization and functional studies of optogenetic tools

To become a reliable candidate as an optogenetic instrument, each protein undergoes various characterization procedures, in which a detailed analysis of its properties is performed. First of all for rhodopsins the measurements of an absorption spectrum must be done. The wavelength of maximum absorption is crucial for optogentics, as it limits the depth, which can be reached by the laser, activating the protein.

The homogeneity and aggregation of protein is studied by means of electrophoresis under native conditions and also by analytical size-exclusion chromatography²⁹. The thermal stability is also determined by different techniques: determination of the melting temperature which can be derived from the first derivative of the denaturation curve obtained by a differential scanning fluorimetry or a thermal unfolding based essay methods^{30–32}.

To define the ion which is pumped throughput the biolipid membrane by the absorption of light-pulses complementary methods are used. The fastest is the measurement of light-induced change of pH in a suspension of E.coli with hight-yield heterologous expression of microbial rhodopsin in outer membrane with the electrodes^{11,33}.

When protein is purified and concentrated in can be reconstituted in lipidic liposomes. The similar to cell suspension pH changes upon illumination shows change of pH of the solution outside the liposome membrane. These pH changes is abolished, when protonophore CCCP (carbonylcyanide m-chlorophenylhydrazone) is added to the cells suspension. In experiments all the known outwardly directed proton pumps (like bR and pR) show the rise of pH. From that it can be concluded if studied rhodopsin inwardly is directed proton pump or not. Changing the composition of solution ions which are pumped can be defined. These experiments can be performed at a wide range of pH values (between pH 5 and 9) to test the possible changes of ion selectivity. Protein orientation in liposomes is important in this case and the most uniformity is achieved if the lipid vesicles of the smallest size is prepared. These lipid vesicles also can be used to prepare black-lipid membrane with orientated membrane proteins for direct photo-current measurements (BLM set up)^{11,34}.

For understanding of rate of pumping activity of the protein, its photocycle must be determined. The laser flash photolysis is common technique to study the microbial rhodopsins photocycle^{11,35–38}. The measurements performed in lipid vesicles with microbial rhodopsins incorporated are obviously more relevant to the native conditions than, for example, detergent micelles or nanodisks. Varying the lipid composition one can modulate the species-specific surrounding of the studied protein. Varying pH the photocycles corresponded to different ion selectivity could be studied.

Before a new potential optogenetic tool is used in neuronal cells, the protein undergoes several functional studies in model systems or another mammalian cells. Thus, the HEK293(T), neuroblastoma SH-SY5Y and NG108-15 cells are transfected with the plasmid with the protein of interest using existing protocols¹¹. After that electrophysiological characterization of the microbial rhodopsin in question by whole cell patch-clamp recordings are done for a rapid proof of its activity and applicability as an optogenetic tool. Photocurrents are measured in response to light pulses.

For in-cell visualization the protein is usually fused with a GFP variant (or Dendra2 fluorescent protein) and then high-resolution confocal laser scanning microscopy is used as in ³⁹.

3.3 Structure-based rational design of optogenetic tools

For rational design of the new optogenetic tools high-resolution structures are required. Thus, even having similar archtecture -7 transmembrane alpha-helices (named A to G or TM1 to TM2) bound with co-factor retianl covalently bound to the Lysine residue in the helix G or TM7 via Schiff base – each branch of rhodopsins has its own structural features which are vital for the function.

Only in 2017 crystal structure of the most important optogenetic tool ChR2 in the ground state was deciphered (Fig. 1)⁴⁰. It indicates the presence of unusial cavities inside the protein, which allows significant rearrengents inside the protein during ion translocation.



Fig. 1: Cavities and highly conserved amino acids of ChR2 and C1C2. a. Protomer A of ChR2 structure and the four main cavities: Extracellular Cavities 1 and 2 (EC1, EC2) and Inner Cavities 1 and 2 (IC1, IC2). DC pair is shown in red ellipse together with water molecule 5. b. C1C2 structure. DC pair is shown in black ellipse, lacks water molecule. The cavities were calculated using HOLLOW⁴¹. TM6 and TM7 helices are not shown. The hydrophobic membrane core boundaries are shown by the black lines.

Another unique feature of the wild-type ChR2 is the existance of the gates insite the pore, especially extracellular gate (EG), which was unknown before. Extremely dense hydrogen bonds network around it reveals the mechanism of its opening, creating wide field for rational design of new optogentic instruments.

As it was said previously, ChR2 lacks ion selectivity, which is the serious problem for optogenetics. To solve this problems, selective pumps can be used, such as light-driven sodium pump KR2, which was already shown to be incorporated into neuronal cells⁴². Crystal structures of different states of KR2 were determined in 2015 simultaniously by 2 groups of scientists. Interestingly, that the protein presents in diiferent oligomeric states in these structures: monomeric at pH around 4.3 and pentameric at pH above 4.9 (Fig. 2). The resolution varies from 1.45-2.3 Å for monomeric blue to 2.2-2.8 Å for pentameric red state. KR2 structures also reveal unique features of the protein, such as the existence of a big ion-uptake cavity, which plays a key role in the selectivity of the protein (especially residues Asn-61 and Gly-263 are involved) and the unusual location of proton acceptor Asp-116.



Fig. 2: Overall architecture of KR2. a, KR2 fold and its orientation in the membrane. b, KR2 pentamer. The novel N-terminal α -helix preceding the helix A is colored blue. The B-C loop is colored orange. The hydrophobic membrane core boundaries were calculated using the PPM server⁴³ and are shown by the black lines.

Using high-resolution structures of the protein several K+ pumping mutants were modelled, the most selective and efficient are G263F, G263W and N61P/G263W. Potassium pumps are of high interest for the optogenetics, as the concentraction of K+ ions is approximately 10 times more than that of Na+ in neurons.

Another perspective proteins to trigger action potential in neural cells as a substitution of ChR2 are inward proton pumps xenorhodopsins, and especially xenorhodopsin from *Nanosalina* (NsXeR). In 2017 functional and structural characterization of that protein was performed. The structure even at relatively low resolution reveals peculiar properties of the inward H+ pump. For instance, the set of 6 key amino acids (motif) differs a lot from those of known rhodopsins (Fig. 3). It also has a short N-terminal α -helix, which caps the inside of the protein.

Thus, having the structural information of the existing and potential optogenetic tools it is possible to optimize or modify them in the way to create new ones with novel functions, high selectivity and major efficiency.



Fig. 3: NsXeR structure. a. Comparison of NsXeR (yellow) and BR (magenta) motifs. Residues are shown as a NsXeR motif (WDSAPK) and BR motif (RDTDDK). Two residues, H48 and D220 in NsXeR are shown as an analogue of D96 residue in BR. b. Putative proton acceptor region in details. Distance between Schiff base and water molecule 2 is shown with red line with arrows (8.0Å). Distances between D76 (D85) and Schiff base in NsXeR (BR) are 4.9Å (3.8Å) respectively. Cavities inside the protein calculated by HOLLOW1.2 are shown transparent pink.

4 Perspectives and unsolved problems

Unfortunately, not all rhodopsins with excellent properties for optogenetic applications can be vectorially expressed in the membranes of interest. This is a serious limitation for the applications, and a repertoire of rhodopsins used in optogenetics must be considerably extended. Moreover, it also creates big difficulties in the study of the physiology of these proteins, for instance, by the patch clamp technique. For solution of this problem a rational screening of targeting leader peptides must be done.

Secondly, the poor choice of suitable proteins imposes a strict limitation to unleashing the huge potential of optogenetics. Selective light-driven cation pumps (Na+, K+, Ca2+) as well as selective ion channels are required for further progress of optogenetics.

Analyses of gene databanks and literature show that the known and characterized rhodopsins are only the tip of the iceberg and that the vast majority of important rhodopsin genes have yet to be discovered and rhodopsins with new properties should be identified and characterized. Moreover, the lack of structural and functional data on the diverse photoactivated proteins, also on those, which were recently identified as those with previously unknown functions (channel rhodopsins, xenorhodopsins, sodium pumps), does not yet allow rational protein engineering of useable tools. Deciphering of high-resolution structures of the target proteins at this point is of high importance.

All the characterized rhodopsins have 6 key functionally important amino acids (Table 2).

RDTDDK	Proton pumps (bR)
RDTEDK	Proton pumps (pR)
RTSADK	Chloride pumps (hR)
RNDQDK	Sodium Pumps (KR2)
RDTKDK	Proton pumps (ESR)
RDTFDK	Sensors (SRI, SRII)
RNTQDK	Chloride pumps (CIR)
RETHDK	Channel Rhodopsins
WDSAPK	Inward Proton Pumps (XeR)

Table 2: The known rhodopsins with different functions and/or properties differ from each other by a 6 letter motif (as a reference bR RDTDDK motif is usually used: R82, D85, T89, D96, D212, K216).

It is known that if the 6 letter motif differs then the functions and/or properties of the proteins differ. Bioinforamtic search of the new rhodopsins with unusual motifs is one of the ways to discover proteins with new functions, which will be applied as a novel optogenetic tools.

On the Fig. 4 the schematical pipeline of such search using artificial intelligence is shown. On this step the increase of determined structures number of rhodopsins from different classes is crucial. The examples of KR2 and NsXeR have already shown that in some cases the right motif can only be defined correctly after structure is solved.



Fig. 4: Left panel: given a set of target sequences and a set of structural templates, we compose structural models of the target sequence. The structural models (red points) are then

projected into the feature space and the artificial intelligence (AI) agent (black circle) is derived. Right panel: microbial rhodopsins are evaluated by different AI agents and the corresponding function is assigned for each microbial rhodopsin (red point s correspond to the proton pumps, blue - chloride pumps, green - sodium pumps, yellow - sensors). Large magenta circle is the most distant from the AI agents, and, thus, represents a promising microbial rhodopsin candidate.

Thus, correct analysis of gene databanks and careful functional and structural characterization of well-known and new rhodopsins will widen the limits of the optogenetics.

References

- Oesterhelt D, Stoeckenius W. Rhodopsin-like protein from the purple membrane of Halobacterium halobium. Nat New Biol. 1971;233(39):149-152. doi:10.1038/newbio233149a0.
- Stoeckenius W, Lozier RH, Bogomolni RA. Bacteriorhodopsin and the purple membrane of halobacteria. *BBA Rev Bioenerg*. 1979;505(3-4):215-278. doi:10.1016/0304-4173(79)90006-5.
- Bogomolni RA, Spudich JL. The photochemical reactions of bacterial sensory rhodopsin-I. Flash photolysis study in the one microsecond to eight second time window. *Biophys* J. 1987;52(6):1071-1075. doi:10.1016/S0006-3495(87)83301-5.
- Lanyi JK, Jurgen Weber H. Spectrophotometric identification of the pigment associated with light-driven primary sodium translocation in Halobacterium halobium. J Biol Chem. 1980;255(1):243-250.
- Wagner NL, Greco JA, Ranaghan MJ, Birge RR. Directed evolution of bacteriorhodopsin for applications in bioelectronics. *J R Soc Interface*. 2013;10(84):20130197-20130197. doi:10.1098/rsif.2013.0197.
- Beja O, Aravind L, Koonin E V., et al. Bacterial rhodopsin: Evidence for a new type of phototrophy in the sea. *Science (80-)*. 2000;289(5486):1902-1906. doi:10.1126/science.289.5486.1902.
- Nagel G, Ollig D, Fuhrmann M, et al. Channelrhodopsin-1: A light-gated proton channel in green algae. Science (80-). 2002;296(5577):2395-2398. doi:10.1126/science.1072068.
- Inoue K, Ono H, Abe-Yoshizumi R, et al. A light-driven sodium ion pump in marine bacteria. *Nat Commun.* 2013;4. doi:10.1038/ncomms2689.
- O'Malley M a. Exploratory experimentation and scientific practice: metagenomics and the proteorhodopsin case. *Hist Philos Life Sci.* 2007;29(3):337-360. doi:10.1007/s40656-014-0001-6.
- Inoue K, Nomura Y, Kandori H. Asymmetric functional conversion of eubacterial lightdriven ion pumps. J Biol Chem. 2016;291(19):9883-9893. doi:10.1074/jbc.M116.716498.
- Shevchenko V, Mager T, Kovalev K, et al. Inward H⁺ pump xenorhodopsin: Mechanism and alternative optogenetic approach. *Sci Adv.* 2017;3(9):e1603187. doi:10.1126/sciadv.1603187.
- 12. Kateriya S. "Vision" in Single-Celled Algae. News Physiol Sci. 2004;19(3):133-137. doi:10.1152/nips.01517.2004.
- 13. Lin JY, Lin MZ, Steinbach P, Tsien RY. Characterization of engineered

channelrhodopsin variants with improved properties and kinetics. *Biophys J.* 2009;96(5):1803-1814. doi:10.1016/j.bpj.2008.11.034.

- Kleinlogel S, Feldbauer K, Dempski RE, et al. Ultra light-sensitive and fast neuronal activation with the Ca 2+-permeable channelrhodopsin CatCh. *Nat Neurosci*. 2011;14(4):513-518. doi:10.1038/nn.2776.
- Prigge M, Schneider F, Tsunoda SP, et al. Color-tuned channelrhodopsins for multiwavelength optogenetics. J Biol Chem. 2012;287(38):31804-31812. doi:10.1074/jbc.M112.391185.
- Berndt A, Schoenenberger P, Mattis J, et al. High-efficiency channelrhodopsins for fast neuronal stimulation at low light levels. *Proc Natl Acad Sci.* 2011;108(18):7595-7600. doi:10.1073/pnas.1017210108.
- Ullrich S, Gueta R, Nagel G. Degradation of channelopsin-2 in the absence of retinal and degradation resistance incertain mutants. *Biol Chem.* 2013;394(2):11-20. doi:10.1515/bchm-2012-0256.
- Pan ZH, Ganjawala TH, Lu Q, Ivanova E, Zhang Z. ChR2 mutants at L132 and T159 with improved operational light sensitivity for vision restoration. *PLoS One*. 2014;9(6). doi:10.1371/journal.pone.0098924.
- 19. Schneider F, Gradmann D, Hegemann P. Ion selectivity and competition in channelrhodopsins. *Biophys J.* 2013;105(1):91-100. doi:10.1016/j.bpj.2013.05.042.
- Gunaydin LA, Yizhar O, Berndt A, Sohal VS, Deisseroth K, Hegemann P. Ultrafast optogenetic control. *Nat Neurosci.* 2010;13(3):387-392. doi:10.1038/nn.2495.
- 21. Klapoetke NC, Murata Y, Kim SS, et al. Independent optical excitation of distinct neural populations. *Nat Methods*. 2014;11(3):338-346. doi:10.1038/nmeth.2836.
- Bamann C, Gueta R, Kleinlogel S, Nagel G, Bamberg E. Structural guidance of the photocycle of channelrhodopsin-2 by an interhelical hydrogen bond. *Biochemistry*. 2010;49(2):267-278. doi:10.1021/bi901634p.
- Berndt A, Yizhar O, Gunaydin LA, Hegemann P, Deisseroth K. Bi-stable neural state switches. *Nat Neurosci*. 2009;12(2):229-234. doi:10.1038/nn.2247.
- Yizhar O, Fenno LE, Prigge M, et al. Neocortical excitation/inhibition balance in information processing and social dysfunction. *Nature*. 2011;477(7363):171-178. doi:10.1038/nature10360.
- Lin JY, Knutsen PM, Muller A, Kleinfeld D, Tsien RY. ReaChR: A red-shifted variant of channelrhodopsin enables deep transcranial optogenetic excitation. *Nat Neurosci*. 2013;16(10):1499-1508. doi:10.1038/nn.3502.
- Govorunova EG, Sineshchekov OA, Li H, Janz R, Spudich JL. Characterization of a highly efficient blue-shifted channelrhodopsin from the marine alga platymonas subcordiformis. J Biol Chem. 2013;288(41):29911-29922. doi:10.1074/jbc.M113.505495.
- Wietek J, Wiegert JS, Adeishvili N, et al. Conversion of channelrhodopsin into a lightgated chloride channel. *Science* (80-). 2014;344(6182):409-412. doi:10.1126/science.1249375.
- Berndt A, Lee SY, Ramakrishnan C, Deisseroth K. Structure-guided transformation of channelrhodopsin into a light-activated chloride channel. *Science (80-)*. 2014;344(6182):420-424. doi:10.1126/science.1252367.
- Cherezov V, Liu J, Griffith M, Hanson MA, Stevens RC. LCP-FRAP assay for prescreening membrane proteins for in meso crystallization. In: *Crystal Growth and Design.* Vol 8.; 2008:4307-4315. doi:10.1021/cg800778j.
- Zheng Y, Qin L, Zacarías NVO, et al. Structure of CC chemokine receptor 2 with orthosteric and allosteric antagonists. *Nature*. 2016;540(7633):458-461. doi:10.1038/nature20605.

- Biswas A, Shukla A, Vijayan RSK, Jeyakanthan J, Sekar K. Crystal structures of an archaeal thymidylate kinase from Sulfolobus tokodaii provide insights into the role of a conserved active site Arginine residue. J Struct Biol. 2017;197(3):236-249. doi:10.1016/j.jsb.2016.12.001.
- Javitt G, Ben-Barak-Zelas Z, Jerabek-Willemsen M, Fishman A. Constitutive expression of active microbial transglutaminase in Escherichia coli and comparative characterization to a known variant. *BMC Biotechnol*. 2017;17(1). doi:10.1186/s12896-017-0339-4.
- 33. Gushchin I, Shevchenko V, Polovinkin V, et al. Crystal structure of a light-driven sodium pump. *Nat Struct Mol Biol.* 2015;22(5). doi:10.1038/nsmb.3002.
- Bamberg E, Apell HJ, Dencher NA, Sperling W, Stieve H, Läuger P. Photocurrents generated by bacteriorhodopsin on planar bilayer membranes. *Biophys Struct Mech*. 1979;5(4):277-292. doi:10.1007/BF02426663.
- Chizhov I, Engelhard M. Temperature and halide dependence of the photocycle of halorhodopsin from Natronobacterium pharaonis. *Biophys J.* 2001;81(3):1600-1612. doi:10.1016/S0006-3495(01)75814-6.
- Chizhov I, Engelhard M, Chernavskii DS, Zubov B, Hess B. Temperature and pH sensitivity of the O(640) intermediate of the bacteriorhodopsin photocycle. *Biophys J*. 1992;61(4):1001-1006. doi:10.1016/S0006-3495(92)81907-0.
- Chizhov I, Schmies G, Seidel R, Sydor JR, Lüttenberg B, Engelhard M. The photophobic receptor from Natronobacterium pharaonis: Temperature and pH dependencies of the photocycle of sensory rhodopsin II. *Biophys J.* 1998;75(2):999-1009. doi:10.1016/S0006-3495(98)77588-5.
- Chizhov I, Chernavskii DS, Engelhard M, Mueller KH, Zubov B V., Hess B. Spectrally silent transitions in the bacteriorhodopsin photocycle. *Biophys J.* 1996;71(5):2329-2345. doi:10.1016/S0006-3495(96)79475-4.
- Bogorodskiy A, Frolov F, Mishin A, et al. Nucleation and growth of membrane protein crystals in meso - a fluorescence microscopy study. Cryst Growth Des. 2015:acs.cgd.5b01061. doi:10.1021/acs.cgd.5b01061.
- 40. Volkov O, Kovalev K, Polovinkin V, et al. Structural insights into ion conduction by channelrhodopsin 2. *Science (80-)*. 2017;358(6366). doi:10.1126/science.aan8862.
- Ho BK, Gruswitz F. HOLLOW: Generating accurate representations of channel and interior surfaces in molecular structures. *BMC Struct Biol.* 2008;8. doi:10.1186/1472-6807-8-49.
- Kato HE, Inoue K, Abe-Yoshizumi R, et al. Structural basis for Na + transport mechanism by a light-driven Na + pump. *Nature*. 2015;521(7550):48-53. doi:10.1038/nature14322.
- Lomize MA, Pogozheva ID, Joo H, Mosberg HI, Lomize AL. OPM database and PPM web server: Resources for positioning of proteins in membranes. *Nucleic Acids Res.* 2012;40(D1). doi:10.1093/nar/gkr703.

B 9 Synthetic Biology – Science and Engineering of Synthetic Biological Systems

Friedrich C. Simmel Physics of Synthetic Biological Systems Physics Department Technische Universität München

Contents

1	Intr	oduction	2			
2	Concepts for engineering of biological systems					
	2.1	From genes to circuits	3			
	2.2	Modularity	7			
	2.3	Programming	10			
3	Rese	earch directions in synthetic biology	10			
	3.1	Bionanoscience	10			
	3.2	Computing with biological circuits	11			
	3.3	Bioproduction and metabolic engineering	15			
	3.4	Cell-free synthetic biology & artificial cells	16			
4	Sum	mary & Conclusions	17			

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

A typical current 'definition' of synthetic biology reads: *Synthetic biology is the engineering of biology. It aims at the design and creation of novel biological parts, devices, and systems that perform useful functions.* The goal of synthetic biology is not so much the understanding of biological phenomena or the elucidation of basic biological principles, but rather the development of a technology based on biological components. Nevertheless, this does not mean that synthetic biology does not address any fundamental questions at all. On the contrary, looked at as a form of 'bottom-up' biology, it may actually contribute to our general understanding of life's processes.

Manipulation of living systems over the centuries

Mankind has manipulated biological systems for many thousands of years. With the emergence of agriculture and sedentism, human societies started to engage in the breeding of plants and animals, and biochemical processes such as fermentation were soon utilized for (bio)technological purposes¹. The way humans interfered with biological systems always reflected their contemporary technological abilities. With the establishment of 'modern' physics and chemistry during the enlightenment period (i.e., in the 17th century), new tools and capabilities for the investigation and manipulation of living systems emerged. We note in passing that the development of the light microscope led to the discovery of biological cells in 1665 [1], which then also sparked the cell hypothesis ('omnis cellula e cellula' [2], also implying that it is impossible to create living cells *de novo*). The beginning of organic chemistry - the chemistry of carbon compounds is mythically related to the synthesis of urea by Wöhler in 1828 (cf. [3]). With this and other achievements of the early organic chemists it became clear that 'naturally' occurring biological compounds were accessible through chemical synthesis. Already in 1915, Emil Fischer coined the term 'chemical-synthetic biology' [4], by which he suggested to modify the chemical building blocks of cells to alter biological processes. Up to a certain degree, the recent development of synthetic biology thus reflects the availability of powerful technologies for the manipulation of biological systems at the *molecular level* - this involves the availability of gene sequences of many organisms, the capability to synthesize gene length DNA relatively inexpensively, and novel microscopic as well as high-throughput methods for analysis.

Key experiments since 2000

The most recent wave of synthetic biology started with a series of key experiments published around the year 2000 - a synthetic bistable gene circuit [5], a genetic oscillator circuit [6], and a synthetic sender-receiver circuit [7]. Rather than analyzing a naturally occurring biological phenomenon, the authors combined existing biological components (in this case gene regulatory elements) to create synthetic dynamical systems ('gene circuits') inside of bacteria. The authors then went on to analyze these systems quantitatively, also involving mathematical modeling, to see whether their behavior met their 'design goal'. As indicated by these experiments, in contrast to standard genetic engineering synthetic biology typically aims at a more far-reaching modification of biological systems, involving multiple interacting components, which together generate novel behaviors.

¹Beer-brewing dates back around 10,000 years.



Fig. 1: Schematic representation of a genetic toggle switch, one of the first synthetic gene 'circuits' realized (image taken from Ref. [5]). In this circuit, two genes (each coding for repressor proteins) mutally repress each other. The system is stable in configurations, in which either one or the other promoter is in an active state. Switching from one state to the other is achieved by the addition of genetic 'inducers', which are small molecules that deactivate the repressors. For further explanations see also Figure 3.

2 Concepts for engineering of biological systems

As engineering biology involves the manipulation of highly complex systems, reliable methods are required to rationally build up large systems of interacting parts whose performance meets a certain design goal. Several engineering disciplines including computer and chemical engineering have developed methodologies that can be utilized in the context of synthetic biology. Typical steps in the engineering of complex systems involve *abstraction* & *modularization* followed by *standardization*, which facilitates the re-use of tested components. In some cases, a *scalable* technology can be developed - then defined procedures exist to rationally increase systems size while retaining its functionality. All of this requires that the procedures and modules are *reproducible* (within specifications) and *robust* with respect to disturbances.

These principles will be treated informally in the following paragraphs. As an example, we will first discuss the engineering of genetic circuits' in bacteria, which are closest in spirit to computer technology 2 . They will also exemplify some of the methods, concepts and challenges for synthetic biology.

2.1 From genes to circuits

Simple gene expression

A simple model of gene expression - production of a protein P from a gene G - is given by [8]:

$$\dot{r} = n_q \alpha_r - \delta_r r \tag{1}$$

$$\dot{p} = \alpha_p r - \delta_p p \tag{2}$$

Here, α_r and α_p represent the rates of RNA and protein production (i.e., transcription and translation), whereas $\delta_{r,p}$ denote the corresponding degradation rates (cf. Figure 2). Lower case letters denote concentrations of the species (i.e., r = [R], p = [P], etc.) and n_g is the copy number of the gene. Now one of the standard tasks is to generate a certain steady state concentration of protein P, which in this case is given by the ratio of the production and degradation rates:

$$p_0 = \frac{n_g \alpha_r \alpha_p}{\delta_r \delta_p} \tag{3}$$

²For simplicity, we will focus on cell-free and bacterial systems in this chapter.



Fig. 2: A simple model of gene expression. From left to right: RNAP binds to the promoter sequence of a gene and transcribes it into RNA with rate α_r . RNA molecules are then translated by ribosomes into proteins (rate α_p). RNA and proteins are degraded with rates δ_r , δ_p . In the case of stable proteins, δ_p corresponds to the dilution caused by bacterial growth.

How can we control these parameters at the molecular level? The copy number n_a represents the gene dosage of the gene coding for the protein. For instance, a gene may be present at several copies on the chromsome (or may be transcribed from several promoters, which complicates matters). When using plasmids - circular DNA molecules capable of replication within the host cell - as gene carriers, one typically has the choice between low (≤ 10), medium ($\approx 10-20$) and high copy numbers (> 20 up to hundreds). The choice of copy number not only has an influence on n_q , but also on noise properties [9] and on the metabolic load [10] put on the host cell. The transcription rate α_r depends, among others, on the choice of the promoter sequence (and the RNA polymerase), whereas α_p is influenced by the strength of the ribosome binding site (RBS). For the latter, a biophysical model has been developed, which allows the comparison of different RBS [11]. This model considers RNA secondary structure in the RBS region, which prevents access of the ribosome, as well as mechanical contributions (when the distance between RBS and start codon AUG does not match the optimum distance of 5 nt) and the free energy gained by ribosome binding. The degradation rates influencing the steady state protein concentrations can also be tuned up to a certain degree. In *E.coli*, RNA is degraded relatively quickly (with a typical lifetime of $\tau_r = 1/\delta_r = 2 \min [8]$), but its stability can be strongly affected by secondary structure or by the presence of RNA binding proteins. Proteins are much more stable than RNA and their lifetime $\tau_p = 1/\delta_p$ is well approximated by the generation time of the bacteria, i.e., typically 30 - 60 min. This means, that when production of a protein is stopped, its cellular concentration will fall over time with decay constant δ_n simply because of bacterial growth and division. If shorter lifetimes are required, proteins can be 'tagged' with degradation tags that direct the proteins to the degradation complex ClpXP, reducing τ_p to about 5 min [12]. How fast is the production of a protein? Equations 1 & 2 can be easily solved, resulting in

$$r(t) = \frac{n_g \alpha_r}{\delta_r} (1 - e^{-\delta_r t}) \tag{4}$$

$$p(t) = \frac{n_g \alpha_r \alpha_p}{\delta_r \delta_p} \left(1 - \frac{\delta_p e^{-\delta_r t} - \delta_r e^{-\delta_p t}}{\delta_p - \delta_r} \right) \stackrel{\delta_r \gg \delta_p}{\approx} \frac{n_g \alpha_r \alpha_p}{\delta_r \delta_p} \left(1 - e^{-\delta_p t} \right)$$
(5)

Thus the concentrations exponentially ramp up to the steady state values with the characteristic times τ_r, τ_p . In some cases, one has to consider that a protein may not be in its active form after production. In particular, fluorescent proteins such as GFP first have to 'mature' to reach a fluorescent state. This can introduce a considerable delay (up to hours) in the appearance of a fluorescence signal from such proteins, which makes their use as a readout for kinetic studies problematic. Fast variants of fluorescent proteins are available (maturation times on the order of 5 min), but in some cases even faster readout is desired, which may be achieved through the use of luciferase (photoluminescence) or other assays.



Fig. 3: Gene regulation and gene circuits. A, Repressor protein P_2 represses transcription of the gene for protein P_1 (e.g., by interference with transcription initiation or elongation). An inducer I can deactivate P_2 and thus activate expression. Transcription (TX) and translation (TL) are shown as a single step. B, Gene transfer function for repression, modeled by a Hill function. Shown is the steady state concentration of P_1 as a function of $[P_2]$. Indicated are design parameters important for the engineering of gene circuits. The baseline is influenced by the 'leakiness' of the promoter, the steepness of the transition by the Hill coefficient (cooperativity) of repression, threshold and fold-change are related to K_2 , copy number, promoter and RBS strength [13]. C, Gene autoregulation. The graph shows production and degradation rates of P_1 as a function of its concentration (Eq. 7). For low $[P_1]$ production exceeds degradation, for high $[P_1]$ degradation dominates, which 'tunes' $[P_1]$ into the transition region. D, Two proteins inhibiting each other's expression (cf. Fig. 1 [5]). The graph shows the 'nullclines' of the production rates of the proteins (cf. Eq. 6), with their three intersections corresponding to fixed points of the dynamical system. The intermediate fixed point is unstable, while the others are stable. The separatrix is indicated by the dashed line. E, Three repressor proteins arranged in a negative feedback loop, corresponding to the topology of the 'repressilator' [6, 14] showing oscillatory protein expression.

Gene regulation and circuits

Genes can influence each other's activity by a variety of regulatory mechanisms, and these can be used to link them together into artificial gene circuits. The canonical form of describing, e.g., *negative regulation* or *repression* of the expression of protein P_1 by another protein P_2 is (cf. Fig. 3A & B):

$$\dot{p}_1(t) = \frac{\alpha_1 K_2^{n_2}}{p_2^{n_2} + K_2^{n_2}} - \delta_p p_1, \tag{6}$$

i.e., repression is modeled as a Hill function with Hill exponent n_2 , whose sigmoidal shape approximates a step function with threshold K_2 for large n_2 . Maybe the simplest gene circuit is given by negative autoregulation - i.e., a protein inhibits its own production:

$$\dot{p}_1(t) = \frac{\alpha_1 K_1^{n_1}}{p_1^{n_1} + K_1^{n_1}} - \delta_p p_1 \tag{7}$$

For high Hill exponent n_1 the production of the protein essentially stops when $p_1 > K_1$, and thus the circuit produces a protein concentration of $\approx K_1$ (Fig. 3C). An important engineering



Fig. 4: Examples for RNA-based gene regulation mechanisms. A, The 'toehold riboregulator' [15] acts at the translational level and is based on a hairpin structure in the RBS region of an mRNA. The RBS is inaccessible to the ribosome in the inactive state. A trigger RNA molecule can invade the hairpin by toehold-mediated strand displacement. B, CRISPR interference [17] is based on sgRNA-directed binding of the dCas9 protein to a location on a gene. The complex interferes with transcription initiation or acts as a roadblock for elongation similar to a repressor protein. The binding site on the gene can be nearly freely chosen (via the sgRNA sequence), and only has to be next to a short 'protospacer adjacent motif' (PAM) sequence 5'-NGG-3'.

aspect of this result is that this concentration is not dependent on production and degradation rates (which may vary from cell to cell and also depends on factors such as the medium or growth state), but is only set by the binding properties of the protein to its regulatory region (and the topology of the circuit³). Negative autoregulatory loops typically have a stabilizing function, and they are frequently used also in synthetic gene circuits. Two repressor proteins can be arranged to mutually repress each other, forming a bistable gene circuit (Fig. 3D), which can be thought of as a 'toggle switch' or a 1-bit molecular memory. Three repressor proteins are used in the 'repressilator' circuit (Fig. 3E), mimicking the topology of a ring oscillator in electrical engineering.

RNA-based gene regulation

Gene expression can be regulated by other mechanisms than using activator or repressor proteins. In this context, RNA-based regulatory processes are particularly interesting for synthetic biology, as they are *sequence programmable*. In contrast to protein-protein and protein-DNA interactions RNA regulation is based on sequence complementarity and can thus be rationally designed. Due to the huge available sequence space, RNA regulatory sequences can be chosen uniquely and with good orthogonality properties. Two examples for synthetic RNA-based regulation are shown in Figure 4. The translational control mechanism of 'toehold switch' riboregulators [15] is related to that of naturally occurring riboswitches [16]. In such riboregulators, the RBS is sequestered within a stable RNA secondary structure, which makes it inaccessible for the ribosome. As shown in the figure, the secondary structure can be broken by short trigger RNA molecules in a process known as 'toehold-mediated strand invasion'. This reveals the RBS and thus activates translation of the mRNA into protein. Another RNA-based regulation mechanism is CRISPR interference (CRISPRi) [17], a derivative of the CRISPR gene editing technique. Similar to regulation by protein transcription factors, CRISPRi acts at the transcrip-

³One should note that it is good practice in electrical engineering to create circuits that are robust with respect to device variability - which essentially is similar to the situation here.

tional level, but the regulatory sequence can be freely chosen by means of so-called single guide RNA (sgRNA) molecules. CRISPR-associated proteins such as the Cas9 protein bind to these sgRNAs, which direct the Cas9-sgRNA complexes to sequence-complementary genetic target sites. CRISPRi utilizes a non-cleaving mutant of a Cas protein (termed dCas9), and thus the dCas9-sgRNA complex simply acts as a roadblock for RNA polymerase. Other RNA mechanisms are based on transcriptional termination or binding of antisense RNA to mRNA molecules [18].

The need for design rules

While traditional genetic engineering has been very successful, e.g., in the expression of certain proteins of interest in a host organism, the creation of synthetic gene circuits represents a considerably more complex engineering task. In the cellular context, the molecules will interact - more or less strongly - with a huge number of other components, which may lead to undesired feedback processes, a reduction of their effective concentration, and thus deterioration of the circuit behavior. Many systems parameters are unknown, influence each other, and in addition are dependent on the cell's growth state [10, 19]. Synthetic components always put a load on the cell's resources, which leads to reduced performance and resource-sharing problems. It is therefore still extremely challenging to effectively engineer large, synthetic circuits inside of living cells [20, 21]. In order to fully realize the grand vision of synthetic biology, effective design rules and strategies for engineering complex biological systems are thus desperately needed.

2.2 Modularity

One of the central engineering paradigms employed for the development of synthetic biology is the concept of *modularity*. Modularity is generally assumed to be a necessary requirement for the 'engineerability' of complex systems. Modules can be defined as discrete and relatively autonomous functional subsystems that can be combined to form a larger functional system according to specific rules. These rules govern the interactions of the modules at their interfaces. Functional modules as well as hierarchical organization can be identified also in biological systems, and there are general arguments why evolution of complex systems may actually lead to a modular structure [22]. Now the general assumption of synthetic biology is that one can create and add novel 'functional modules' to existing biological systems without a complete loss of function (of either the module or the pre-existing system). Furthermore, in order to create ever larger systems, multiple modules have to be combined in a rational manner and remain functional. In the following, we will briefly discuss some of the modularity challenges faced in synthetic biology, and proposed solutions.

Biobricks and genetic engineering standards

Already on a technical (i.e., chemical and biochemical) level, the engineering of biological systems with multiple components represents a considerable challenge. In contrast to traditional genetic engineering involving the introduction of 'only' a single protein-coding gene into a bacterial cell, methods have to be found to rationally synthesize systems with many interacting genetic 'modules'. Inspired by design rules in electrical engineering, researchers in synthetic biology created the so-called 'biobrick standard' (https://biobricks.org/) [23] for the concatenation of multiple gene modules (sometimes called 'biological parts' - see also the 'Registry



Fig. 5: An overview of the engineering challenges faced in synthetic biology and proposed solutions (image taken from Ref. [27]). A synthetic module has to be embedded into the complex regulatory and metabolic networks already present in the host organism. The scheme indicates modular and hierarchical design with a designed interface between module and 'chassis'.

of Standard Biological Parts' https://parts.igem.org/ connected to the iGEM competition) to a circuit. In recent years, novel cloning and gene assembly techniques (Gibson assembly, Golden Gate cloning, etc.) were developed that together with improved gene synthesis methods facilitate the creation of large DNA constructs containing multiple genes [24]. It is conceivable that - supported also by developments in lab automation - we will soon be able to 'write' arbitrarily long DNA sequences (cf. also the 'GP-Write' project - http://engineeringbiologycenter.org/). The main challenge then is not the physical creation of the DNA molecules coding for the synthetic biological systems - but getting them to work!

Retroactivity

One of the first modularity challenges arises simply when connecting two modules together e.g. one gene producing a transcription factor X that binds to the promoter of a downstream gene. Because of the binding interaction, the concentration of unbound $X(x_{free})$ will be different in the absence or presence of the downstream gene. Interpreting x_{free} as the output of the first gene, this value is obviously affected by the presence of the downstream component - an effect sometimes termed retroactivity [25]. Multiple downstream promoters will be occupied depending on their respective promoter strengths and the concentration of free transcription factors will be drained accordingly.

This situation is quite analogous to that found in electrical circuitry - the output voltage of a given device is different when measured in 'open circuit' configuration, or when put under load. As observed by Uri Alon [8, 26], 'modules in engineering, and presumably also in biology, have special features that make them easily embedded in almost any system. For example, output nodes should have 'low impedance', so that adding on additional downstream clients should not drain the output to existing clients.' As in electrical engineering, there are different strategies



Fig. 6: Automation of gene circuit synthesis using the program Cello [32]. A circuit is specified as a logic table that relates Boolean outputs with a set of inputs based on the hardware description language Verilog. This is automatically converted into a logic circuit, which can be rewritten in terms of combinations of NOR gates. These can be easily realized as a genetic circuit (cf. also Fig. 8.). In order to create a functioning circuit, the necessary genetic 'parts' are chosen from a gate library according to their transfer functions.

to avoid excessive retroactivity. One is to only use a small fraction of the output X as the input for an amplifying unit, whose output is then fed into the downstream device. The amplifier acts as an insulator that shields the upstream element from retroactivity effects.

Circuit-chassis interactions

Another important challenge is the interaction of a synthetically introduced component or circuit with the host organism - also known as the 'circuit-chassis' interaction (Fig. 5) [27]. Similar as in the retroactivity problem, the effective concentrations of circuit components of the synthetic circuit will be reduced because of undesired interactions with already present components. This will affect the input-output relations between the different stages of the circuit and potentially lead to unexpected feedback mechanisms. Furthermore, the presence of the circuit will draw on the host's resources and therefore affect its physiology and growth, in turn changing the circuit behavior. As a result, the same circuit may actually behave very differently in isolation (e.g., in a cell-free system), in different host organisms, and under different physiological conditions of the host. One approach to avoid at least the most obvious interactions is to use 'orthogonal' components - i.e., molecular components designed to only interact with the desired partners and not with any other component. In genetic circuits, this amounts to careful choice of transcription factors or design of RNA regulatory components. Chemical biologists approach the problem by developing 'orthogonal chemistries' that do not interfere with cellular biochemistry. There are also on-going attempts to develop artificial base-pairs (xeno-DNA [28]), introduce them into cells, and expand the genetic code [29] in order to establish chemical processes that do not interfere (up to a certain degree) with cellular processes.

Compartmentalization and spatial organization

One physical approach to define modules and reduce undesired interactions is the utilization of compartments for the separation of chemical processes from each other. The interface of the module with the chassis is then naturally defined by the boundary of the compartment (and thus potentially better controlled). In biology, this strategy is extensively used at various length scales: from the local (nanoscale) chemical environment in the active site of an enzyme, over the organization of multiple enzymes in clusters (e.g., carboxysomes, encapsulins, etc. [30]), to the organization of biochemical processes within membranous or also membrane-less organelles [31], or to cellular and even multicellular differentation and organization. In synthetic biology, there are thus ongoing efforts to create artificial enzyme clusters and organelles, and also artificial cell-scale compartments that interact with 'real' cells.

2.3 Programming

'Programming' and 'programmability' are frequently used terms in synthetic biology and related disciplines. In most cases, this simply means that the choice of a specific DNA (or RNA or protein) *sequence* controls a certain biological process or results in a defined molecular structure. Compared to conventional computer-programming this is similar to writing a program in machine code.

Researchers now hope to be able to drive the programming analogy much further. Rather than working at the hardware level, an abstract higher level description of synthetic biological processes would be desirable that can be compiled down to machine code (i.e., DNA sequence) in an automated manner. In order to realize this vision, modularity and well-defined interactions at the hardware level will be critically important. One of the most advanced efforts in this direction is the automated design of genetic circuits based on the hardware description language Verilog. Chris Voigt and coworkers developed the gene circuit design software Cello (http://www.cellocad.org/) that allows specifications of, e.g., logical truth tables that are then converted to a combination of suitable genetic logic gates (see also below) [32]. The circuit components are chosen from a database of gates with well-characterized transfer functions and matched in silico to optimally approximate the desired behavior. This approach appears to be rather successful for the implementation of logic gates for sensor/actuator applications and demonstrates that biological circuits can indeed be rationally engineered up to a certain degree. However, it is not clear whether the more general goal of abstract high-level programming can ever be achieved in the context of extant living systems, or whether simpler chassis such as minimal biological cells, artificial cellular or other cell-free systems will be required.

3 Research directions in synthetic biology

3.1 Bionanoscience

Arguably the simplest form of synthetic biology is the creation of synthetic biomolecular nanostructures. This is 'simple' in the sense that it 'only' amounts to the design of molecules that assume a specific desired conformation - i.e, it typically refers to the engineering of an equilibrium property and not the dynamics of a complex out-of-equilibrium system. Over the past years, maybe the most successful branch of bionanoscience has been DNA nanotechnology [33]. In this field base-pairing interactions between DNA molecules are used to create supramolecular


Fig. 7: Arbitrarily shaped biomolecular nanostructures can be created with the DNA origami technique [34]. A, In DNA origami, a long single-stranded DNA molecule (shown in gray) serves as a scaffold strand to which multiple shorter DNA molecules (called staples) can hybridize. The overall structure is created by multiple strand crossovers, at which the staples connect distant sequence domains on the scaffold. B, Schematic representation of a DNA origami structure with black scaffold strand and colored staples, omitting the helicity of DNA. Arrows indicate the $5' \rightarrow 3'$ direction of the DNA backbone.

structures by self-assembly. DNA has been used to create extended lattices and crystals or also discrete objects. The most well-known types of DNA assemblies are the DNA origami structures [34, 35]. In the DNA origami technique, on the order of 200 oligonucleotides of length 20 - 40 nucleotides (nt) are hybridized to a long single-stranded DNA molecule with a length of typically $7 - 8 \times 10^3$ nt. The short DNA 'staple' strands force the long 'scaffold' strand to fold into a three-dimensional shape, which is uniquely determined by the choice of the staple strands and their sequences. In DNA nanotechnology, there is a straightforward relation between a molecular code - the DNA sequence - and a three-dimensional molecular structure, which is due to the relatively well-behaved and well-understood, sequence-dependent interactions between DNA molecules. This property has facilitated the development of powerful computational tools that enable the design of DNA nanostructures at a relatively abstract level (i.e., without molecular details) using a CAD-like program interface (http://cadnano.org/) [36]. While DNA nanotechnology is very successful in the creation of complex molecular structures, it has its drawbacks in terms of chemical versatility, and also its potential use in a biological context. For this reason, efforts are underway to establish nanotechnologies based also on RNA [37, 38] or proteins [39, 40], or of complexes of these biomolecules. Similar to naturally occurring molecular complexes, biological nanostructures can be genetically encoded and therefore produced in vivo, which also allows their control by synthetic gene circuits.

3.2 Computing with biological circuits

A considerable fraction of work in synthetic biology is devoted to the implementation of computational functions into biological systems. There are several reasons for this: First, computer technology is the prime example for a technology based on modularity. Modern computer chips are among the most complex physical structures created by humankind, which has become possible through the application of modular and scalable design rules. It therefore seems that engineering principles developed in the context of computer science should be useful, up to a certain degree, also for the engineering of other complex systems and structures. Then, there is the superficial and popular analogy that genes are some kind of molecular program running on celullar hardware [41]. Accordingly, cells are just computers (or robots), and we have to learn how to program them. There is probably some truth in this statement, but one also has to recognize that biological hardware is very different from electronics hardware. Among others, information is not represented by electrons or voltages, but by conformational states, molecular modifications or concentrations - the components are (typically) not connected by wires (everything can interact with everything), and the computational processes are stochastic rather than clocked. Computing in such 'amorphous' media is possible, but follows other rules [42].

Hence, there are different and partly opposing strategies to realize computation in biological systems: One is to artificially impose computational models used in electronic computing on the biological system. This often amounts to the implementation of Boolean logic circuits (see below) in terms of biochemical reaction networks or gene circuits [8, 32, 43, 44], which requires threshold processes with a clear ON/OFF behavior. Once processes with sufficiently Boolean characteristics have been identified, they can be rationally connected into Boolean circuits (of course they have to be compatible, i.e., the output of one process has to serve as an input for another). The major advantage of this approach is that it enables universal, general purpose computing (in practice, however, this approach is limited by the imperfect Boolean nature and modularity of the components).

The opposite approach is to acknowledge that biochemical processes are, in reality, not Boolean and therefore treat cells as analog computing devices [45]. As with analog computers, such an approach does not allow the implementation of general purpose computing - essentially, a new device has to be constructed for each particular application. On the other hand, special purpose analog devices can be much more efficient in terms of space, time, or energy requirements.

There are a wide variety of computational models that are closer to the biological situation than Boolean electronic computers, and these are frequently discussed in the context of synthetic biological computation. For instance, analogies between gene or signaling networks and neural network computing can be made [8, 46, 47], agent-based computing more closely matches the distributed and amorphous nature of biological systems, and cellular automata capture aspects of biological growth processes [48]. In fact, in all these cases, biology served as the inspiration for the computational models themselves. One advantage of the availability of abstract computational models is the existence of rigorous proofs concerning their computational power, which gives an idea what could potentially be realized in a synthetic biological setting.

Logic gates and circuits

Logic functions can be implemented in a relatively straightforward way at the transcriptional level using combinations of transcription factors [49]. Transcription factors are typically activated or de-activated by small molecule inducers, which are regarded as the inputs of the gates. Promoter activity (with corresponding downstream gene expression) is taken as the output. Next to the trivial NOT (input negation) and YES (sensor/input repetition function) gates, AND and OR gates are most easily constructed in this context. From these, NAND and NOR gates (Fig. 8) can be produced, and combinations of these can be used - in principle - to generate any other Boolean function ({NAND}, {NOR}, and {AND; NOT} each are complete bases for



Fig. 8: Genetic implementation of a NOR gate. A, Symbolic representation of a NOR gate with inputs I_1 , I_2 and output, and corresponding truth table. B, Genetic realization of a NOR gate. Here two inducers are used as inputs, which switch ON transcription from their respective promoters. As a result, a repressor protein will be produced in the presence of either I_1 OR I_2 . The repressor protein represses production of the output, and therefore the output is only produced in the absence of the inducers. High and low concentrations of the participating molecules are interpreted as Boolean 1 and 0, respectively.

Boolean functions). Apart from transcription factors, also CRISPRi [43] or riboregulators [44] can utilized for the implementation of logic functions.

In contrast to the realization of single logic gates, it is considerably more challenging to connect multiple gates into logic circuits. An obvious problem is the wiring of the gates, for which 'signal homogeneity' is required as the output of each gate must be usable as an input for another gate. Another challenge for circuitry are mismatched dynamic ranges of inputs and outputs, which prevents proper signal propagation through a series of gates. As mentioned above, this problem has been addressed using a computational approach that matches the transfer functions of the components of gene circuits *in silico* [32].

Memory

A fundamental memory unit is a bistable component that can be controllably switched from one state to the other. The 'toggle switch' [5] based on two mutually repressing genes (cf. Figs. 1, 3D) was inspired by the RS flip flop, one of the traditional memory circuits in electronics. Other types of memory have been based on recombinase enzymes that modify the orientation of gene coding regions with respect to promoter regions by inversion [50, 51]. Such components have also been used, in a very limited way, to implement arithmetic functions such as counting [52]. From a computer science view, systems with memory - i.e., internal states - are much more powerful than simple logic circuits and allow the solution of wider class of problems. It will therefore be of great interest to combine synthetic gene circuitry with molecular memory in the future.

Communication

While gene circuits and memory can endow single cells with computational functions, cell-tocell communication can further enhance their capabilities in the context of distributed, multicellular systems. Among others, signaling and communication is required for the generation of spatial responses such as pattern formation [53] and differentiation [54]. In particular, it may be advantageous to implement different functions into specialized cells (or compartments) and let them cooperatively interact via defined communication channels (cf. also Fig. 5).

The most popular signaling mechanism employed in (bacterial) synthetic biology is based on the naturally occurring 'quorum sensing' phenomenon. In quorum sensing bacteria can sense their own population density by sending out and receiving diffusible inducer signals (called





Fig. 9: Synthetic 'communication channels' for bacteria. A, Bacteria may communicate via small diffusible inducers similar as in naturally occurring quorum sensing systems. In the scheme, the 'sender' bacteria contain a gene for the AHL synthase LuxI (AHL stands for acyl homserine lactone), which catalyzes the production of the AHL signal. AHL can diffuse out of the bacterial cell and enter other cells. 'Receiver' bacteria contain a gene for the transcriptional activator LuxR, which activates transcription from the P_{lux} promoter in the presence of AHL (image taken from Ref. [58]). B, Communication via phages. A DNA 'message' written into a phagemid can be packaged into bacteriophage M13 particles, which upon release can infect other bacteria with the message (image taken from Ref. [57]).

auto-inducers, AIs) that stimulate the production of even more signals through a positive feedback mechanism. As a result, the bacteria can switch on specific genes only when a critical population density has been reached [55]. At the molecular level, quorum sensing involves an autoinducer synthase that produces the signals, and a transcriptional activator that responds to this signal by gene activation. In one of the original works in synthetic biology, Ron Weiss and Thomas Knight dissected the natural quorum sensing systems of *Vibrio fischeri* consisting of the AI synthase LuxI and the activator LuxR into two subsystems, which only contained the sender and receiver parts, respectively [7]. This system together with a few other QS systems is widely in use in synthetic biology.

Other bacterial communication mechanism such as horizontal gene transfer (conjugation [56]) or communication via phages [57] have been investigated for application in synthetic biology. While being more complex, they have the advantage of signaling information-rich genetic messages rather than small molecules.

Applications

Apart from engineering aspects and addressing fundamental conceptual issues, implementation of computational functions into biological systems has a wide range of potential applications [13]. Genetic circuits can be used to evaluate sensory input patterns received by the bacteria (e.g. the presence of externally supplied or intracellularly generated molecules, or other stimuli such as light), and calculate appropriate responses. This is of interest for medical applications (sensing of disease-related molecules potentially coupled to the delivery of drugs [59]), or sense & destroy applications [60] in medicine or in an environmental context. Biocomputing is also of interest for control systems for bioproduction, where regulatory circuits could be used to autonomously tune parameters towards optimum production conditions [13].

3.3 Bioproduction and metabolic engineering

One of the major potential application areas of synthetic biology is the engineering and improvement of biotechnological production processes. Metabolic engineering - the creation of industrially relevant chemicals through re-engineering of existing metabolic processes - has been developed independently of synthetic biology already decades ago. Starting already in the 1960s, mathematical tools have been developed to describe the 'metabolic flux' through networks such as 'Biochemical Systems Theory' by Savageau [61] or 'Metabolic Control Analysis' by Kacser & Burns and Heinrich & Rapoport [62]. Many of the concepts developed in metabolic engineering anticipate later work in systems biology.

Synthetic biology now promises to allow a more far-reaching engineering of biological systems for the production of useful metabolites [63]. This also involves technological advances such as large-scale gene synthesis, the creation of cells with reduced or strongly altered genomes, and the expansion of the genetic code - topics we cannot cover here due to limited space. Control circuits such as discussed in the previous section could help to operate cells under optimum conditions for bioproduction, e.g., in terms of usage of resources, energy consumption, cell viability, or reactant stoichiometry.

Spatial organization of metabolic flux

Interesting *physical* approaches to improve the performance of metabolic pathways involve the spatial organization of cooperating enzymes into clusters, along supramolecular scaffolds or inside of compartments. Scaffolding or compartmentalization may be achieved using in vivo-produced RNA or protein nanostructures as briefly discussed above. Spatial organization of enzymes is supposed to be beneficial for a variety of reasons. In a metabolic pathway involving multiple enzymatic steps, intermediate reactants may be lost by diffusion or degradation, or to enzymes not participating in the pathway. Metabolic flux through a pathway may be optimized by controlling the stoichiometry of the participating enzymes [64]. In some cases the order of chemical steps may be controlled by a corresponding spatial ordering of the enzymes (a popular analogy is an 'assembly line').

A quantitative treatment shows that the effectiveness of enzyme clustering or co-localization strongly depends on the enzymes' catalytic rates and the diffusion properties of the intermediates. For instance, in a simple two-enzyme pathway

$$S \xrightarrow{E_1} I \xrightarrow{E_2} P,$$
 (8)

in the regime $[S] \ll [E_1]$ a single enzyme E_1 will produce $k_{cat,1}[S]/K_{M,1}$ molecules of the intermediate I per second (here $k_{cat,1}, K_{M,1}$ are the enzyme's catalytic rate and Michaelis constant, respectively). Within a region of characteristic size L, the intermediates will either be converted to product P with a rate $k_{cat,2}[E_2]/K_{M,2}$ (where $[E_2]$ is the mean concentration of E_2 in the region), or leave the region without reaction by diffusion with rate $\sim D/L^2$. It turns out that for slow E_2 enzymes, i.e., when $k_{cat,2}[E_2]/K_{M,2} \times L^2/D < 1$, clustering of the enzymes around the source E_1 is most effective in terms of overall metabolic flux. For faster reactions, a distribution of the enzymes in space is more effective [65, 66].

3.4 Cell-free synthetic biology & artificial cells

One possible solution to the problems arising from the implementation of synthetic modules into extant living systems - the 'circuit-chassis' interactions mentioned above - is to get rid of the 'chassis' and create synthetic biological systems from the bottom up. Working with a limited (and known) number of molecular components also promises to be more amenable to quantitative and predictive modeling, and therefore rational design. Furthermore, in the absence of the biological context, components of very different origin - biological or chemical - can be combined more freely. Accordingly, in the past few years there has been an increasing interest in the realization of 'cell-free' synthetic biological systems [67, 68]. These were either based on cell-free gene expression systems, or using systems assembled from purified biochemical components.

Cell-free gene expression

The use of cell extracts marks the beginning of modern biochemistry - in 1907 Eduard Buchner won the Nobel Prize in Chemistry for his demonstration of the fermentation of sugar in yeast extract (rather than using yeast cells). Later cell extracts were used for the elucidation of the molecular mechanisms of protein expression as well as for small-scale protein synthesis for structural biology and other applications [69]. Until recently, cell-free gene expression systems were optimized for protein expression, but not used for the implementation and study of synthetic biological circuits. In recent years, however, improved protocols have been developed for the generation of bacterial cell extracts, in which the operation of dynamical gene circuits as well as the assembly of relatively complex multi-protein structures has become possible [70]. One remarkable achievement in this context was the complete assembly of a bacteriophage in a cell extract to which only the genome of the phage had been added [71].

As cell extracts contain all the protein content of a cell, many biological processes can be faithfully 'simulated' in them. For instance, they not only include RNA polymerase and ribosomes, but also the protein and RNA degradation machinery, which are important for the realization of dynamical systems. Due to the production process, components of a cell extract are diluted and less active compared to the living cell, and thus reaction kinetics are typically slowed down [72]. Still, cell extracts are black boxes with unknown components and idiosyncrasies. Cell-free transcripton/translation systems containing only purified (and known) components such as the PURExpress system therefore are a popular (but more costly) alternative [73].

A general problem of cell-free systems is the absence of a metabolism - during the operation of the systems, energy rich compounds such as ATP, NTPs and amino acids are used up, while waste products accumulate. In order to extend the lifetime of the reactions, ATP regeneration systems have been developed [74]. Alternatively the systems may be operated in open reaction chambers with permanent exchange of nutrients and waste products. In one of the most advanced works in this context, genetic oscillators were operated inside of a microfluidic chemostat - here continuous supply of nutrients and genes combined with dilution was used to 'simulate' the dynamical effect of cell growth and division [75]. In this case, cell-free gene expression was successfully used for 'rapid prototyping' of gene circuits, that also worked when implemented in bacteria.

can be motivated in various ways. On the one hand, the synthesis of an artificial cellular system with a metabolism and the ability to grow, divide, and potentially evolve will throw light on the nature of the transition of matter from an abiotic to a living state. When also restricting the chemistry of artificial cells to pre-biotically plausible molecules, the resulting 'protocells' could give us an idea about the emergence of the earliest cellular life forms [76, 77]. Also from a more biotechnological viewpoint, compartmentalized systems acting as cell-scale bioreactors are of considerable interest [78] - as already discussed above, for bioproduction applications compartmentalization could help to avoid loss of intermediate compounds by diffusion or sidereactions and separate competing biochemical processes from each other (cf. the related 'modularity' discussion above). In most cases the term 'artificial cell' has been applied rather loosely to compartmentalized biochemical systems in general, as no-one has succeeded in creating a genuine artificial living cell-scale system so far.

Technically, cell-scale compartmentalization can be achieved in severalfold ways. Using microfluidics, biochemical mixtures can be easily encapsulated into water-in-oil-emulsion droplets with diameters in the range of $10 - 100 \,\mu m$ [79]. Such emulsion droplets have the disadvantage that they are essentially closed systems and do not allow materials exchange with their environment. On the other hand, they can be generated in large numbers and with good monodispersity, which is useful for screening and molecular evolution applications. Vesicles with lipid bilayer boundaries more closely resemble cellular compartments, but they are more difficult to prepare [80]. Peptide vesicles or vesicles defined by membranes made from amphiphilic polymers here provide interesting alternatives.

As indicated, among the major challenges for the creation of a genuine artificial cell is the realization of a metabolism that allows execution of cellular processes over extended periods of time. This should also enable cells to grow, and growth would have to be coupled to an appropriate cell division mechanism. Artificial cells may be operated in a 'heterotrophic' mode - in this case, energy-rich compounds such as ATP are supplied externally [82] and transported through the cell boundaries via appropriate membrane channels [80]. When autonomous ATP generation is desired, however, the reconstitution of an ATP synthase will be necessary, which further requires the establishment of a proton-motive force across the membrane. Growth has been demonstrated in cases where polymerization of RNA inside of a vesicle led to an increase in osmotic pressure and thus vesicle swelling, where the additional membrane area was provided by the uptake of lipids from the surroundings [81]. Thus several of the defining processes for an artificial cell have already been realized, but so far there is no example of a working system, in which all of these have been coupled and integrated into an autonomous, self-perpetuating synthetic cell-scale system.

4 Summary & Conclusions

Synthetic biology is a fast-moving field that aims at the (re-)engineering of biological devices and systems. If successful, the research program of synthetic biology will have a wide-ranging impact on future biotechnology, medicine and sustainability. In order to realize a genuine synthetic biology, also fundamental questions will have to be addressed, which have to do with the complexity of biological systems and their 'engineerability'.

We may finally ask what physics can do for synthetic biology - and vice versa? On the one

hand, a thorough quantitative understanding of biological processes will be the basis for the engineering of synthetic biological systems. Quantitation enables computational modeling and supports abstraction and modularization, which are required for the engineering of complex biological systems. On the other hand, being able to create synthetic biological systems will allow us to test physicochemical hypotheses about the living state in a well-defined experimental setting, and in this sense synthetic biology simply follows the successful tradition of the physical sciences.

References

- R. Hooke, Micrographia: Or, Some Physiological Descriptions of Minute Bodies Made by Magnifying Glasses, with Observations and Inquiries Thereupon, J. Martyn and J. Allestry, London (1665).
- [2] R. Virchow (1855), see https://en.wikipedia.org/wiki/Rudolf_Virchow
- [3] D. McKie, Nature 153, 608 (1944).
- [4] E. Fischer, Die Kaiser-Wilhelm-Institute und der Zusammenhang von organischer Chemie und Biologie, in: Untersuchungen aus Verschiedenen Gebieten: Vorträge und Abhandlungen Allgemeinen Inhalts, M. Bergmann (ed.), Springer Berlin Heidelberg (1924).
- [5] T. S. Gardner, C. R. Cantor and J. J. Collins, Nature 403, 339-342 (2000).
- [6] M. B. Elowitz and S. Leibler, Nature 403, 335-338 (2000).
- [7] R. Weiss, T. F. Knight, A. E. Condon and G. Rozenberg, DNA Computing, 6th International Workshop on DNA-Based Computers, DNA6 2054, 1-16 (2000).
- [8] U. Alon, An introduction to systems biology: design principles of biological circuits, CRC Taylor & Francis, Boca Raton, FL (2007).
- [9] J. Paulsson and M. Ehrenberg, Q. Rev. Biophys. 34, 1-59 (2001).
- [10] S. Klumpp, Z. Zhang and T. Hwa, Cell 139, 1366-1375 (2009).
- [11] H. M. Salis, E. A. Mirsky and C. A. Voigt, Nature Biotechnol. 27, 946-U112 (2009).
- [12] J. Andersen, C. Sternberg, L. Poulsen, S. Bjorn, M. Givskov and S. Molin, Appl. Environ. Microbiol. 64, 2240-2246 (1998).
- [13] J. A. N. Brophy and C. A. Voigt, Nat. Meth. 11, 508-520 (2014).
- [14] L. Potvin-Trottier, N. D. Lord, G. Vinnicombe and J. Paulsson, Nature 538, 514-517 (2016).
- [15] A. A. Green, P. A. Silver, J. J. Collins and P. Yin, Cell 159, 925-939 (2014).
- [16] M. Mandal and R. R. Breaker, Nat. Rev. Mol. Cell Biol. 5, 451-463 (2004).
- [17] Lei S. Qi, M.H. Larson, L.A. Gilbert, J.A. Doudna, J.S. Weissman, A.P. Arkin and W. A. Lim, Cell 152, 1173-1183 (2013).
- [18] F. J. Isaacs, D. J. Dwyer and J. J. Collins, Nat. Biotechnol. 24, 545 (2006).
- [19] M. Scott, C. W. Gunderson, E. M. Mateescu, Z. Zhang and T. Hwa, Science 330, 1099-102 (2010).
- [20] P. E. M. Purnick and R. Weiss, Nat. Rev. Mol. Cell Biol. 10, 410-422 (2009).
- [21] R. Kwok, Nature 463, 288-290 (2010).
- [22] L. H. Hartwell, J. J. Hopfield, S. Leibler and A. W. Murray, Nature 402, C47-C52 (1999).
- [23] B. Canton, A. Labno and D. Endy, Nat. Biotechnol. 26, 787-793 (2008).
- [24] P. A. Carr and G. M. Church, Genome engineering, Nat. Biotechnol.27, 1151-1162 (2009).
- [25] D. Del Vecchio, A. J. Ninfa and E. D. Sontag, Mol. Syst. Biol. 4, 161 (2008).
- [26] U. Alon, Science **301**, 1866-1867 (2003).

- [27] S. Rolli, M. Mangold and K. Sundmacher, Chemical Engineering Science 69, 1-29 (2012).
- [28] V. B. Pinheiro, A. I. Taylor, C. Cozens, M. Abramov, M. Renders, S. Zhang, J. C. Chaput, J. Wengel, S.-Y. Peak-Chew, S. H. McLaughlin, P. Herdewijn and P. Holliger, Science 336, 341-344 (2012).
- [29] D. A. Malyshev, K. Dhami, T. Lavergne, T. Chen, N. Dai, J. M. Foster, I. R. Corrła and F. E. Romesberg, Nature 509, 385-388 (2014).
- [30] T. O. Yeates, C. A. Kerfeld, S. Heinhorst, G. C. Cannon and J. M. Shively, Nat. Rev. Microbiol. 6, 681-691 (2008).
- [31] A. A. Hyman, C. A. Weber and F. Jülicher, Annu. Rev. Cell Dev. Biol. 30, 39-58 (2014).
- [32] A. A. K. Nielsen, B. S. Der, J. Shin, P. Vaidyanathan, V. Paralanov, E. A. Strychalski, D. Ross, D. Densmore and C. A. Voigt, Science 352, aac7341-aac7341 (2016).
- [33] M. R. Jones, N. C. Seeman and C. A. Mirkin, Science 347, 1260901-1260901 (2015).
- [34] P. W. K. Rothemund, Nature **440**, 297-302 (2006).
- [35] S. M. Douglas, H. Dietz, T. Liedl, B. Högberg, F. Graf and W. M. Shih, Nature 459, 414-8 (2009).
- [36] S. M. Douglas, A. H. Marblestone, S. Teerapittayanon, A. Vazquez, G. M. Church and W. M. Shih, Nucleic Acids Res. 37, 5001-5006 (2009).
- [37] P. Guo, Nat. Nanotechnol.5, 833-842 (2010).
- [38] C. Geary, P. W. K. Rothemund and E. S. Andersen, Science 345, 799-804 (2014).
- [39] H. Gradisar, S. Bozic, T. Doles, D. Vengust, I. Hafner-Bratkovic, A. Mertelj, B. Webb, A. Sali, S. Klavzar and R. Jerala, Nat. Chem. Biol. 9, 362-366 (2013).
- [40] P.-S. Huang, S. E. Boyken and D. Baker, Nature 537, 320-327 (2016).
- [41] D. Bray, Wetware: A Computer in Every Living Cell, Yale University Press (2009).
- [42] H. Abelson, D. Allen, D. Coore, C. Hanson, G. Homsy, T. F. Knight, R. Nagpal, E. Rauch, G. J. Sussman, R. Weiss and G. Homsy, Commun. ACM 43, 74-82 (2000).
- [43] A. A. Nielsen and C. A. Voigt, Mol. Syst. Biol. 10, 763-763 (2014).
- [44] A. A. Green, J. Kim, D. Ma, P. A. Silver, J. J. Collins and P. Yin, Nature 548, 117-121 (2017).
- [45] R. Daniel, J. R. Rubens, R. Sarpeshkar and T. K. Lu, Nature 497, 619-623 (2013).
- [46] J. Vohradsky, FASEB J. 15, 846-854 (2001).
- [47] J. Kim, J. J. Hopfield and E. Winfree, Advances in Neural Information Processing Systems 17, 681-688 (2004).
- [48] S. Wolfram, Cellular automata as models of complexity, Nature **311**, 419-424 (1984).
- [49] N. E. Buchler, U. Gerland and T. Hwa, On schemes of combinatorial transcription logic, Proc. Natl. Acad. Sci. USA 100, 5136 (2003).
- [50] J. Bonnet, P. Subsoontorn and D. Endy, Proc. Natl. Acad. Sci. U.S.A. 109, 8884-8889 (2012).
- [51] P. Siuti, J. Yazbek and T. K. Lu, Nat. Biotechnol. **31**, 448-452 (2013).
- [52] A. E. Friedland, T. K. Lu, X. Wang, D. Shi, G. Church and J. J. Collins, Science 324, 1199-1202 (2009).
- [53] S. Basu, Y. Gerchman, C. H. Collins, F. H. Arnold and R. Weiss, Nature 434, 1130-4 (2005).
- [54] M. Isalan, C. Lemerle and L. Serrano, Plos Biol. 3, 488-496 (2005).
- [55] B. A. Hense, C. Kuttler, J. Mller, M. Rothballer, A. Hartmann and J. U. Kreft, Nat. Rev. Microbiol. 5, 230-9 (2007).
- [56] C. Smillie, M. P. Garcillan-Barcia, M. V. Francia, E. P. C. Rocha and F. de la Cruz, Microbiol. Mol. Biol. Rev. 74, 434-452 (2010).

- [57] M. E. Ortiz and D. Endy, J. Biol. Engin. 6, 16 (2012).
- [58] T. Ramalho, A. Meyer, A. Muckl, K. Kapsner, U. Gerland and F. C. Simmel, PLoS One 11, e0145829 (2016).
- [59] M. Mimee, A. C. Tucker, C. A. Voigt and T. K. Lu, Cell Systems 1, 62-71 (2015).
- [60] S. Gupta, E. E. Bram and R. Weiss, ACS Synth. Biol. 2, 715-723 (2013).
- [61] M. A. Savageau, Biochemical systems analysis. A study of function and design in molecular biology, Addison-Wesley, Reading, MA (1976).
- [62] D. A. Fell, Biochem. J. 286, 313-330 (1992).
- [63] J. Nielsen and JayD. Keasling, Cell 164, 1185-1197 (2016).
- [64] J. E. Dueber, G. C. Wu, G. R. Malmirchegini, T. S. Moon, C. J. Petzold, A. V. Ullal, K. L. J. Prather and J. D. Keasling, Nat. Biotechnol. 27, 753-759 (2009).
- [65] A. Buchner, F. Tostevin and U. Gerland, Phys. Rev. Lett. 110, 208104 (2013).
- [66] M. Castellana, M. Z. Wilson, Y. Xu, P. Joshi, I. M. Cristea, J. D. Rabinowitz, Z. Gitai and N. S. Wingreen, Nat. Biotechnol. 32, 1011-8 (2014).
- [67] D. C. Harris and M. C. Jewett, Curr. Op. Biotechnol. 23, 672-678 (2012).
- [68] M. T. Smith, K. M. Wilding, J. M. Hunt, A. M. Bennett and B. C. Bundy, FEBS Lett, 588, 2755-2761 (2014).
- [69] F. Katzen, G. Chang and W. Kudlicki, Trends Biotechnol. 23, 150-156 (2005).
- [70] Z. Z. Sun, C. A. Hayes, J. Shin, F. Caschera, R. M. Murray and V. Noireaux, Journal of visualized experiments: JoVE (79), e50762 (2013).
- [71] J. Shin, P. Jardine and V. Noireaux, ACS Synth. Biol. 1, 408-413 (2012).
- [72] E. Karzbrun, J. Shin, R. Bar-Ziv and V. Noireaux, Phys. Rev. Lett. 106, 048104 (2011).
- [73] Y. Shimizu, T. Kanamori and T. Ueda, Methods 36, 299-304 (2005).
- [74] F. Caschera and V. Noireaux, Biochimie 99, 162-168 (2014).
- [75] H. Niederholtmeyer, Z. Z. Sun, Y. Hori, E. Yeung, A. Verpoorte, R. M. Murray and S. J. Maerkl, eLife 4, e09771 (2015).
- [76] S. Rasmussen, Protocells: Bridging Nonliving and Living Matter, MIT Press, Cambridge, MA (2009).
- [77] I. A. Chen and P. Walde, CSH Persp. Biol. 2, a002170-a002170 (2010).
- [78] A. Pohorille and D. Deamer, Trends Biotechnol. 20, 123 (2002).
- [79] R. K. Shah, H. C. Shum, A. C. Rowat, D. Lee, J. J. Agresti, A. S. Utada, L. Y. Chu, J. W. Kim, A. Fernandez-Nieves, C. J. Martinez and D. A. Weitz, Materials Today 11, 18-27 (2008).
- [80] V. Noireaux, Y. T. Maeda and A. Libchaber, Proc. Natl. Acad. Sci. U. S. A. 108, 3473-3480 (2011).
- [81] M. M. Hanczyc, S. M. Fujikawa and J. W. Szostak, Science 302, 618-22 (2003).
- [82] S. S. Mansy, J. P. Schrum, M. Krishnamurthy, S. Tob, D. A. Treco and J. W. Szostak, Nature 454, 122-125 (2008).

C 1 Membranes and vesicles

Timon Idema Department of Bionanoscience Kavli Institute of Nanoscience Delft University of Technology, The Netherlands

Contents

1	Intr	oduction	2			
2	Mer	ane physics				
	2.1	Membrane free energy	3			
	2.2	Membrane fluctuations	5			
3	Membrane shapes					
	3.1	Spheres	6			
	3.2	Tubes	6			
	3.3	Constrained volume shapes	7			
	3.4	Axisymmetric vesicles	8			
4	Multi-component membranes 9					
	4.1	Line tension and membrane shape	9			
	4.2	Measuring line tensions and Gaussian moduli	11			
5	Membrane-mediated interactions					
	5.1	Interactions on an asymptotically flat surface	12			
	5.2	Interactions on closed surfaces	13			
	5.3	Interactions in living systems	14			
A	Som	ne useful differential geometry	15			

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Membranes are ubiquitous in the cell. Google 'cell' and you will find a large collection of images like figure 1, which shows a schematic of an animal cell. Almost everything you see in the image is a membrane. Even though the basic physical structure of all those membranes is the same, what is immediately clear from figure 1 is that they exhibit many different shapes, from simple spheres to complex networks. Unsurprisingly, the shapes of the various parts of the cell are strongly coupled to their specific function [1]. For instance, small spherical vesicles are used for transport and signaling, whereas the extended endoplasmic reticulum network is home to a large number of ribosomes, whose function is to synthesize membrane proteins. The shape of the plasma membrane surrounding the entire cell depends strongly on the cell type. It can create folds used for absorbing or releasing chemicals, and both long finger-shaped and broad sheet-like protrusions for cell locomotion.



Fig. 1: Anatomy of a typical animal cell. Figure from [2].

The basic building blocks of all biological membrane are lipid molecules. Lipids consist of a polar 'head' group, and one or more hydrocarbon chains, also known as 'tails'. Because the heads are hydrophilic and the tails are hydrophobic, the whole molecule becomes ampiphilic. When dissolved in water, it will therefore spontaneously self-organize into structures that allow for both minimum contact between the water and the tails, and maximum contact between the water and the heads. For lipids with a fairly large head group and only one tail, this structure is a micelle: a spherical construct with tails on the inside and heads on the surface. For lipids with two tails, the optimal structure is a bilayer of lipids, consisting of two leaflets with again tails on the inside and heads on both surfaces, which is the structure of a membrane.

A typical lipid bilayer membrane is only 4-5 nm thick, but can easily span many microns in both lateral directions. When left to their own in a water-rich environment, membranes will spontaneously from closed structures to prevent any exposure of the lipid tails to the water, therefore also defining an inside and an outside. Artificial membrane vesicles can easily be created in the lab, using either bulk [3–5] or microfluidic [6–9] methods.

Although water can (slowly) diffuse across a lipid bilayer membrane, larger molecules cannot. The volume enclosed by a membrane can therefore be controlled by applying an osmotic pressure: changing the concentration of some solute in the outside environment will cause a flow of water, resulting in a change of the enclosed volume. Initially spherical membranes can thus be deformed to more interesting shapes. Likewise, membranes can act as an electrical capacitor, allowing for control of some of their properties by the application of electric fields. With large enough fields, the energy penalty for the exposure of lipid tails to water can even be overcome, resulting in the creation of holes in the membrane. This electroporation is the primary method by which pieces of DNA are introduced into cells in genetic engineering. In living cells, membrane composition, osmotic pressure, and electrical potential is actively controlled by proteins. Specialized proteins can cause lipids to move from one leaflet to the other, split off small membrane vesicles, or pump ions and macromolecules across the membrane, allowing the cell to communicate, eat, and move [1].

2 Membrane physics

2.1 Membrane free energy

Lipid bilayer membranes constitute a fascinating physical system in their own right. Within the plane of the membranes, the lipids can move about freely, and the membrane thus forms a 2D-fluid. However, when we try to bend the membrane in the direction normal to its plane, it behaves like an elastic solid: any deformation away from the equilibrium shape will carry an energy penalty that is quadratic in that deformation (which is simply Hooke's law). Any physical model of the membrane has to account for at least these two facts. The fluid (or more precisely liquid) nature of the membrane has two consequences. First, like any fluid, the membrane cannot support shear forces (or in other words, a shear will cause it to flow). Second, like all liquids, the membrane is effectively incompressible, which in our current context means that its total surface area is conserved. The membrane shares these two properties with soap films, which are also (effectively) 2D liquids but without the elastic bending property. To account for the incompressibility, we impose conservation of membrane area, which we can do by a Lagrange multiplier term in the free energy:

$$E_{\rm area} = \int_{S} \sigma dA, \tag{1}$$

where S is the entire surface (usually closed), dA the area element, and σ the 2D equivalent of the pressure, which has the dimensions of surface tension (and is therefore usually also called the surface tension, even though the membrane is not the surface of anything). Equation (1) is the 2D version of the well-known pdV term in the first law of thermodynamics: it simply represents the work necessary to compress the fluid. Because the fluid does not support any shear stress, it is also the only contribution the fluid gives to the free energy.

Describing the bending energy of the membrane requires a little more work. To get it, we make the assumption that the membrane is a thin elastic sheet. Calculating the curvature at any point



Fig. 2: Curved surfaces. (a) A saddle point on a two-dimensional surface embedded in \mathbb{R}^3 . The thick red lines indicate the principal directions. If the positive and negative curvatures are equal, the mean curvature at the saddle point is zero. If the surface extends to infinity, its Gaussian curvature is negative. (b) Coordinate system on an axisymmetric vesicle. The z-axis coincides with the axis of symmetry. The vesicle is parametrized using the arc length s along the contour. The radial coordinate r gives the distance from the symmetry axis and the coordinate z the distance along that axis. The shape of the vesicle be given as r(z), r(s), or in terms of the contact angle ψ of the contour as a function of either s or r. The geometric relations between r, z and ψ are given in equation (15).

of a thin rod is easy: simply find the tangent circle to the rod at that point, and the curvature is the inverse of that circle's radius. In practice, the easiest way to do this is by calculating the second derivative of the function describing the shape of the rod in space. On a surface, we can define a local 2D coordinate system, and describe the surface as a function of these two coordinates. Like for a line, we can calculate the second derivatives of this function, which now gives us a symmetric matrix. The eigenvalues of this matrix are the curvatures corresponding to the tangent circles in the direction of the two basis vectors, giving us two curvatures c_1 and c_2 (known as the principal curvatures, see figure 2a). Of course, the values of these curvatures depend on the coordinates chosen. We can get a coordinate-independent version by taking the trace and the determinant of the curvature matrix, which are simply the sum and product of the two eigenvalues, and give us two curvatures of the surface:

$$H = \frac{1}{2}(c_1 + c_2) \qquad \text{mean curvature} \tag{2}$$

$$K = c_1 c_2$$
 Gaussian curvature (3)

Technical details on how to calculate H and K are given in appendix A. The bending energy of the surface is now given by an expansion in the curvature up to quadratic terms [10, 11]:

$$E_{\text{bend}} = \int_{S} \left[\frac{\kappa}{2} (2H - C_0)^2 + \kappa_{\text{G}} K \right] \mathrm{d}A,\tag{4}$$

where κ and κ_{G} are known as the bending and Gaussian modulus respectively. C_0 is any spontaneous (i.e. pre-bending) curvature the membrane may have due to the shape of the lipids; for

symmetric bilayers we typically have $C_0 = 0$ (membranes with nonzero spontaneous curvature are discussed in chapter C2). The factors 2 are historical.

Although the mean and Gaussian curvature seem similar in nature, they are mathematically very different. The mean curvature is extrinsic, which means that it can only be calculated by embedding the surface in a higher-dimensional space. The Gaussian curvature on the other hand is intrinsic (a fact that Gauss himself found so intriguing that he called this observation his Theorema Egregium, or remarkable theorem). To illustrate the difference, let's take the surface of the Earth, which you (hopefully) know is curved. We surface-dwellers can determine the planet's Gaussian curvature without ever going to space, simply by drawing a triangle on the ground and measuring its internal angles, which won't add up to 180° . Thanks to another theorem (Gauss-Bonnet), we can even use the triangle (or any other shape we draw on the ground) to calculate the Gaussian curvature of the patch of Earth inside the shape by simply measuring the curvature of the boundary of our shape. In particular, if we don't draw any boundary at all, but take the whole surface as our shape, the theorem tells us that the total Gaussian curvature is a constant, which depends on the topology of the shape alone. For a closed surface, the second term in equation (4) thus integrates to a constant, which is why it is often left out of the description of the membrane. Since we'll be discussing shapes with a boundary as well (for membranes with multiple domains), we'll retain it where applicable.

The total free energy of the membrane is given by the Canham-Helfrich model, which is simply the sum of the area and bending terms:

$$E_{\rm CH} = \int_{S} \left[\frac{\kappa}{2} \left(2H - C_0 \right)^2 + \kappa_{\rm G} K + \sigma \right] \, \mathrm{d}A. \tag{5}$$

2.2 Membrane fluctuations

The energy given by equation (5) is a functional of the membrane shape: we can calculate its value for any shape and thus find the shape that minimizes the energy. However, that will not be the exact shape the membrane assumes at finite temperature, which is reflected by the fact that $E_{\rm CH}$ is a free energy: the observed shape will only minimize $E_{\rm CH}$ on average. Thermal fluctuations will cause the actual shapes to deviate from this average, and rather than one fixed shape, we'll get an ensemble of shapes, which can all be observed with a probability equal to their Boltzmann factor, $P(E) = \exp(-E_{\rm CH}/k_{\rm B}T)$.

For biological and biomimetic membranes, the bending modulus κ is of the order of 20 $k_{\rm B}T$. Large fluctuations away from the equilibrium membrane shape will therefore be suppressed, but the fluctuations that are present are large enough to be observed. This is actually fortunate, as from the fluctuation spectrum, we can determine the value of κ and σ . To see how this works, we consider a piece of membrane that in equilibrium is flat. We'll use regular Cartesian coordinates x and y to parametrize this piece of membrane, and describe any deviation from the flat surface in the z direction by a function u(x, y) (this parametrization, known as the Monge gauge, is used frequently in the membrane literature). A straightforward exercise in calculus allows us to express $E_{\rm CH}$ as an integral over u and its derivatives:

$$E_{\rm CH} = \frac{1}{2} \int_{L \times L} \left[\kappa (\nabla_{\perp}^2 u)^2 + \sigma (\nabla_{\perp} u)^2 \right] \mathrm{d}x \mathrm{d}y, \tag{6}$$

plus constant terms with $\kappa_{\rm G}$ and σ that do not influence our results. In equation (6), ∇_{\perp} refers to the gradient with respect to the 'flat' coordinates x and y, and we have taken our piece of

membrane to be a square of size $L \times L$. We now expand u(x, y) in Fourier modes

$$u(\mathbf{x}) = \sum_{\mathbf{q}} u_{\mathbf{q}} e^{i\mathbf{q}\cdot\mathbf{x}}, \qquad u_{\mathbf{q}} = \frac{1}{L^2} \int_{-L/2}^{L/2} \mathrm{d}x \int_{-L/2}^{L/2} \mathrm{d}y \, u(\mathbf{x}) e^{-i\mathbf{q}\cdot\mathbf{x}}, \tag{7}$$

where $\mathbf{q} = (q_x, q_y) = \frac{2\pi}{L}(l_x, l_y)$ with $l_x, l_y \in \mathbb{Z}$. Substituting this expansion in equation (6), we obtain an expression that is quadratic in each Fourier mode

$$E = \frac{L^2}{2} \sum_{\mathbf{q}} \left[\kappa (\mathbf{q} \cdot \mathbf{q})^2 + \sigma (\mathbf{q} \cdot \mathbf{q}) \right] u_{\mathbf{q}} u_{\mathbf{q}}^*.$$
(8)

The equipartition theorem tells us that each quadratic mode in the energy contributes $\frac{1}{2}k_{\rm B}T$ to the thermal energy of a system. If we equate the energy of each of the modes with that thermal energy, we find for the thermal average of a mode

$$\left\langle |u_q|^2 \right\rangle = \left\langle u_{\mathbf{q}} u_{\mathbf{q}}^* \right\rangle = \frac{1}{L^2} \frac{k_{\mathrm{B}} T}{\kappa q^4 + \sigma q^2},\tag{9}$$

where q is the length of q. Equation (9) can be fitted to the fluctuation spectrum of an actual membrane to obtain the values of the material parameters σ and κ [14]. A similar expression can be derived for the fluctuations around a spherical reference shape [15] and even used for vesicles containing multiple membrane domains [16].

3 Membrane shapes

3.1 Spheres

It probably won't surprise you that the shape for which the bending energy of a membrane is minimized is a sphere. In fact, for a sphere, the bending energy is a constant, independent of the radius: $E_{\text{bend}} = 8\pi\kappa + 4\pi\kappa_{\text{G}}$. We can use this observation to derive bounds on the value of the (rather elusive) Gaussian modulus. Suppose that through some process we manage to split an originally spherical membrane in two smaller spheres, with the same total area as the original sphere. The difference in energy between the new and the original configuration would then be $8\pi\kappa + 4\pi\kappa_{\text{G}}$. If our original sphere is to be stable, this energy difference had better be a positive number, so we impose that $\kappa_{\text{G}} > -2\kappa$. On the other hand, the Gauss-Bonnet theorem tells us that for a closed surface, the integral over the Gaussian curvature is given by $4\pi(1-g)$, where g is the genus of the surface, i.e., the number of holes in it. If κ_{G} were positive, a spherical membrane would therefore be unstable to the formation of holes. We thus get stable spheres as long as $-2\kappa < \kappa_{\text{G}} < 0$.

3.2 Tubes

To get more interesting shapes, we should apply additional conditions. One option is to actively perturb the membrane by exerting a point force on a large spherical membrane vesicle. Experimental results show that applying such a force on a 'giant' unilamellar vesicle (or GUV, with a radius of $10 - 50 \ \mu\text{m}$) results in the extraction of a cylindrical membrane tube with uniform cross section [17]. In this case the total energy of the system is given by

$$E_{\text{tube}} = \int_{S} \left(\frac{\kappa}{2} (2H)^2 + \sigma\right) dS - fL, \qquad (10)$$

where f is the applied force and L is the displacement of the point where the force is attached in the direction of that force. Specifically, for a cylindrical tube of radius R and length Lequation (10) reads

$$E_{\text{tube}} = \left(\frac{\kappa}{2}\frac{1}{R^2} + \sigma\right)2\pi RL - fL.$$
(11)

Equation (11) shows a competition between two effects: the bending rigidity term tends to increase the tube radius, whereas the surface tension term tends to reduce it. A stable solution for an applied force f_0 can be obtained by choosing the proper radius R_0 such that the two effects exactly cancel. The values of f_0 and R_0 for given κ and σ are found from the stability condition that the derivatives of E_{tube} with respect to R and L should vanish. They give [18–20]:

$$R_0 = \sqrt{\frac{\kappa}{2\sigma}}$$
 and $f_0 = 2\pi\sqrt{2\kappa\sigma}$. (12)

For typical values of $\kappa \approx 10 \ k_{\rm B}T \approx 40 \ {\rm pN} \cdot {\rm nm}$ and $\sigma = 0.05 \ {\rm pN/nm}$ we get $R_0 \approx 20 \ {\rm nm}$ and $f_0 \approx 13 \ {\rm pN}$. The tube radius is thus several orders of magnitude smaller than that of an experimental vesicle, which means that the implicit assumptions that any surface and volume constraints on the tube could be ignored, were justified. Although the typical tube radius is too small to resolve in an optical setup, the force can be readily measured in an experiment where a tube is extracted from a vesicle using optical tweezers [21]. Because the force is comparable to the force that can be exerted by molecular motors (a typical kinesin motor for example can exert a maximum force of about 5 pN), several motors acting together can also extract a tube, as has been demonstrated experimentally [22–24].

3.3 Constrained volume shapes



Fig. 3: Shape of red blood cells. (a) Micrograph of human red blood cells, showing their distinct biconcave shape. Image courtesy of the National Institutes of Health (U.S.A.), scalebar 5 μ m. (b) Numerically obtained shape of a red blood cell, from the minimization of the bending energy (5), for a fixed enclosed volume and membrane area. The calculations were performed using the Surface Evolver software package by Brakke [25].

An alternative additional condition is to fix the volume enclosed by the membrane. The sphere is the shape that encloses the maximal volume given its area; by forcing the volume to be less than that of a sphere we therefore create some 'excess area'. One particular such shape is the biconcave one of the red blood cell, where the enclosed volume is about half that of the sphere with the same area. Analytical expressions for such shapes are not easy to obtain, but numerically minimizing the curvature energy of a uniform closed membrane given an enclosed volume and total membrane area is a tractable task. The software package Surface Evolver by Brakke [25] does just that. Figure 3 shows an example numerical result, where we begin with an arbitrary shape with the set amount of enclosed volume and surface area, and allow the curvature energy to relax. Independent of the original shape, we invariably find the biconcave shape of the red blood cell.

In general, a differential equation for the mean curvature of a closed uniform vesicle with specified area and enclosed volume can be obtained through variational analysis. The energy is in this case given by the (mean) curvature energy with two Lagrange multiplier terms, one for the area (where the multiplier is the surface tension) and one for the volume (where the multiplier is the pressure difference across the membrane):

$$E = \int_{S} \left(\frac{\kappa}{2} (2H)^{2} + \sigma \right) dA + p \int dV$$
(13)

The calculation of the first variation of this energy is lengthy but straightforward and was first performed by Ou-Yang and Helfrich [26]. An alternative approach that uses Lagrange multipliers to enforce all relevant geometrical relations can be found in [27]. The condition that this variation should vanish for an equilibrium shape results in the shape equation

$$p - 2\sigma H + 4\kappa H(H^2 - K) + 2\kappa \Delta H = 0, \tag{14}$$

where Δ is is the Laplace-Beltrami differential operator on the membrane surface (equation 30).

3.4 Axisymmetric vesicles

Equation (14) becomes a lot more tractable if we apply it to axisymmetric vesicles. Such vesicles are completely specified by giving the contour shape in a plane which contains the axis of rotation. Typically the axes of this plane are labeled r (horizontal) and z (vertical), where the z-axis is the axis of rotation. Because the contour is a curve in \mathbb{R}^2 we can parametrize it using the arc length along the contour from an arbitrary starting point, typically the topmost point of the contour. The coordinates r(s) and z(s) of any point on the contour are then related via the contact angle $\psi(s)$ on any point of the contour (see figure 2b):

$$\dot{r} = \frac{\mathrm{d}r}{\mathrm{d}s} = \cos\psi(s)$$
 and $\dot{z} = \frac{\mathrm{d}z}{\mathrm{d}s} = -\sin\psi(s).$ (15)

Substituting the axisymmetric expressions in the shape equation (14) gives third-order differential equation for $\psi(s)$ [28]. This equation can be integrated once to give a second order equation in ψ [29]:

$$\ddot{\psi}\cos\psi = -\frac{1}{2}\sin\psi\dot{\psi}^2 - \frac{\cos^2\psi}{r}\dot{\psi} + \frac{\cos^2\psi + 1}{2r^2}\sin\psi + \frac{\sigma}{\kappa}\sin\psi + \frac{p}{2\kappa}r - \frac{1}{2\pi\kappa}\frac{f}{r}.$$
 (16)

In equation (16), f is an integration constant, which we readily identify as the applied force in the z-direction (the same force we had in equation 10). Because equation (16) is nonlinear, finding a general analytical solution has proved challenging (nobody has succeeded thus far). Numerical solutions are easy to find however, though surprisingly even easier for the third-order equation, which turns out to be more stable [19].

There is an alternative way of deriving the differential equation (16), by writing the energy (13) as an action, or an integral over a Lagrangian $\mathcal{L} = \mathcal{L}(\psi, \dot{\psi}, r, \dot{r}, z, \dot{z})$. This approach has the advantage that it gives us the proper differential equation for each axisymmetric patch of the vesicle surface, and also the conditions at their boundaries [30, 31]. For a patch that runs from $s = s_1$ to $s = s_2$ we have

$$E = 2\pi\kappa \int_{s_1}^{s_2} \mathcal{L} \mathrm{d}s,\tag{17}$$

with

$$\mathcal{L} = \frac{r}{2} \left(\dot{\psi} + \frac{\sin\psi}{r} \right)^2 + \frac{\sigma}{\kappa} r + \frac{p}{2\kappa} r^2 \sin\psi + \gamma (\dot{r} - \cos\psi) + \eta (\dot{z} + \sin\psi).$$
(18)

In equation (18) we added two additional Lagrange multipliers γ and η to enforce the geometrical relations (15). Variation of the functional E with respect to the variables ψ , r, z, γ and η gives their respective Euler-Lagrange equations, which for any variable x read

$$\frac{\mathrm{d}}{\mathrm{d}s}\frac{\partial\mathcal{L}}{\partial\dot{x}} - \frac{\partial\mathcal{L}}{\partial x} = 0.$$
(19)

Unsurprisingly, from this process we recover equations (15) and (16). As an added bonus however, we also get conditions on the boundary of our patch, which allow us to study membranes with multiple domains.

4 Multi-component membranes

Biological membranes contain many different kinds of lipids, which a wide range of characteristics. Some lipids are charged while others are not, some have very large head groups, some have longer tails than others, and some have fully saturated tails (meaning that there are no double bonds between the carbon atoms) while others have one or more unsaturated link. Membranes also contain cholestorol, which looks a bit like a lipid in the sense that it has a polar head and a hydrophobic tail, though both are much smaller than those of a typical lipid, and cholesterol therefore can fill up gaps between lipids. In order to mimic the behavior of biological lipids more closely, one can make vesicles that contain multiple lipid types in the lab. Mixing different components of course may not work, as they may phase-separate into differently composed phases, just like water and oil tend to do. A well-studied example of such a phase-separating mixture consists of a lipid with saturated tails, one with tails containing an unsaturated bond, and cholesterol. At physiological conditions such a system will typically mix into a liquid-ordered phase rich in saturated lipids (and cholesterol), and a liquid-disordered phase rich in unsaturated lipids. The names of the phases refer to the presence or absence of long-range order between the direction of the lipid tails, see figure 4.

4.1 Line tension and membrane shape

Phase separation into two (or more) domains leads to the formation of a line tension at the domain boundary. In the case of lipid bilayers, there are several factors that contribute to the emergence of this line tension. First, like for any two coexisting phases, the boundary is not perfectly sharp: the concentrations of the various molecules present do not make a jump at the



Fig. 4: *Phase separation of a ternary mixture (saturated and unsaturated-tail lipids and cholesterol) into a liquid-ordered and a liquid-disordered phase. (a) Schematic cross-section of a phase-separated lipid bilayer. Saturated lipids in blue, unsaturated in red, cholesterol in green. (b) Gibbs phase triangle showing the composition of the lipid bilayer, with each of the vertices corresponding to a membrane that consists of a single component (unsaturated lipid, saturated lipid, and cholesterol), and points in the interior representing mixtures. The ternary mixture exhibits phase separation into two coexisting phases, connected by tie lines (black lines). The blue dots indicate critical points. (c) At lower temperatures, there may even be three coexisting phases, liquid-ordered, liquid-disordered, and gel (pink region), bordered by regions with two coexisting phases. (d) Gibbs prism showing Gibbs triangles at different temperatures.*

domain boundary but rather have a smooth transition when we go from one domain to the other. The resulting concentration gradient carries a free energy penalty, which for a two-dimensional surface corresponds to an effective energy per unit length, i.e. a line tension. Second, because there are differences in size between the lipids, the membrane may have a different thickness in different domains. Again, there will not be a sharp transition, but rather a smooth change from one thickness to the other, carrying an energy penalty that contributes to the line tension. The total energy of a membrane with two domains is given by

$$E = \sum_{i=1}^{2} \int_{S_i} \left(\frac{\kappa_i}{2} (2H)^2 + \bar{\kappa}_i K + \sigma_i \right) \mathrm{d}A + p \int \mathrm{d}V + \tau \oint_{\partial S} \mathrm{d}l,$$
(20)

where τ is the magnitude of the line tension, and ∂S indicates the boundary of the domains. The combination of coexisting domains with a line tension in a two-dimensional membrane and the possibility of bending in the direction perpendicular to the membrane give rise to an energy trade-off. On the one hand, the bending energy term tends to keep the membrane as flat as possible, while on the other hand the line tension term tries to get the domain boundary to be as short as possible. If the line tension is high enough, it can therefore force the membrane to create 'bulges': the piece of membrane inside a domain can bulge out into the third dimension, reducing its boundary length at the price of increased curvature, see figure 5.

If a vesicle is prepared from a uniformly mixed lipid phase (typically at high temperature) and allowed to demix over time (for instance after cooling down), domains will initially grow within the plane of the membrane. These domains freely diffuse around, and will therefore frequently run into each other and merge. Once the domains have reached a critical size, they will bulge out, as shown in figure 5a. We can estimate that critical size from a dimensional argument: the line tension has dimension of energy per unit length, and the bending modulus of energy, so the ratio $\xi = \kappa/\tau$ defines a length, known as the invagination length [32]. For domains larger than this length, budding is advantageous, and can be further stimulated by applying an osmotic pressure difference, lifting the strict constraint on the enclosed volume.



Fig. 5: *Membrane shapes governed by the interplay between curvature and line tension. (a) The presence of a line tension between domains may cause the domains to bulge in the third dimension. (b) On a vesicle, many domain bulges may coexist for a long time, even though merging them would lead to a lower total energy. (c) Fully phase-separated 'snowman' or 'barbapapa' vesicle with two domains. Fitting the shape of this vesicle allows us to extract values for the line tension and the difference in Gaussian modulus between the domains. Figures by S. Semrau [16, 34].*

Once domains have budded, they can coexist for a long time, as they repel each other via the deformations they impose on the membrane [33–35]. These membrane-mediated interactions also occur for proteins, which are much smaller (see chapter C2); domains are one possible way of quantifying these interactions and studying their (non-linear) effects [34–36]. Domains have also been hypothesized to occur in the membranes of living cells [37], but never been directly observed. Because extracts of cellular plasma membranes do exhibit domain formation, the most likely reason why large domains do not form in living cells is that their membranes are actively recycled and mixed at high rate, preventing domains from growing large enough to be seen. Nonetheless, small domains may be important for grouping proteins, and membrane-mediated interactions between proteins certainly do occur, for example in viral coat assembly and clathrin-mediated endocytosis.

4.2 Measuring line tensions and Gaussian moduli

Even though bulged domains may coexist for a long time, the ground state of a multi-domain vesicle remains full separation. The domains will therefore eventually merge, creating a two-domain vesicle with a typical snowman or barbapapa shape, as shown in figure 5c. While these shapes have no biological relevance, they are useful experimentally as they allow us to determine the line tension directly. Moreover, in addition to a line tension due to differences in composition or thickness, there is a third energetic cost associated with a domain boundary in a lipid bilayer membrane that can be estimated from the snowman shape. Although the Gauss-Bonnet theorem tells us that the integral over the Gaussian curvature over a closed surface is a constant, the Gaussian modulus may be different for different domains. We can still use the Gauss-Bonnet theorem to relate the integral over the surface to a contribution on the boundary, but when the moduli of the two domains are not equal, the boundary terms will not cancel. At the boundary we thus pick up an additional term that depends on the difference of the Gaussian moduli, $\Delta \kappa_{\rm G}$. This term is qualitatively different from the line tension term in equation (20), which is most easily seen from its effect on the boundary conditions we get for

the shape equation (16) [31]:

$$\lim_{\varepsilon \downarrow 0} (\kappa_2 \dot{\psi}(\varepsilon) - \kappa_1 \dot{\psi}(-\varepsilon)) = -(\Delta \kappa + \Delta \bar{\kappa}) \frac{\sin \psi_0}{r_0}, \tag{21}$$

$$\lim_{\varepsilon \downarrow 0} \left(\kappa_2 \ddot{\psi}(\varepsilon) - \kappa_1 \ddot{\psi}(-\varepsilon) \right) = \left(2\Delta \kappa + \Delta \bar{\kappa} \right) \frac{\cos \psi_0 \sin \psi_0}{r_0^2} + \frac{\sin \psi_0}{r_0} \tau, \tag{22}$$

where the boundary is at s = 0, $r_0 = r(0)$, $\psi_0 = \psi(0)$, $\Delta \kappa = \kappa_2 - \kappa_1$, and $\Delta \bar{\kappa} = \bar{\kappa}_2 - \bar{\kappa}_1$. As the snowman shape is axisymmetric, it can be found by solving equation (16) for each domain, with boundary conditions (21) and (22), plus continuity of r and ψ at the boundary. The bending moduli and surface tensions can be determined from the fluctuation spectrum of each domain independently, as discussed in section 2.2, and the pressure difference across the membrane can be estimated from the Laplace pressure. We are then left with two unknowns, the line tension and difference in Gaussian modulus, which can be determined from fluctuations of the domain boundary, which is done most easily in the regime where the domains do not bulge [38].

5 Membrane-mediated interactions

If you put two bowling balls (or two bodies) on a mattress, they'll locally deform the mattress, and if placed close enough together, attract each other. The attractive force originates in the elastic deformation of the mattress due to the presence of the balls, which is minimized for the case that the two balls are at one spot. Since membranes exhibit elastic deformations as well, objects included or adhered to the membrane experience similar interactions. These objects may be lipid domains, colloids or polymers bound to some lipids, or proteins included in the lipid bilayer; the basic physical principles underlying their interactions are all the same. A major difference between the membrane and the mattress however is that the membrane is also fluid, so the adhered or included objects can move around in the membrane as well, leading to qualitatively different interactions as compared to the purely elastic mattress. Unsurprisingly, it also matters if the membrane is curved at the scale of the imposed deformation or not, especially if it is also closed on this scale (e.g. if it forms a cylinder or sphere with a radius not that different from that of the imposed deformation).

5.1 Interactions on an asymptotically flat surface

The simplest case, for an asymptotically flat membrane with point-like inclusions, was already studied in the 1990s [39, 40]. The presence of the inclusions locally puts constraints on the membrane shape. If the inclusions are linked to an underlying substrate (such as the cytoskeleton in a cell, or a coverslip in an experiment with a supported lipid bilayer), the constraints simply fix the position and possibly the slope of the membrane at the location of the inclusion. For free membranes, the position obviously isn't fixed, as the membrane can adapt its shape to minimize the total energy. Likewise, an inclusion can be tilted if that results in a lower total energy. For point-like inclusions in free membranes, what is imposed is therefore the curvature, or the second derivative of the shape. For a rotationally symmetric conical inclusion, the imposed curvature c equals the ratio of the inclusion's opening angle α to its radius $a, c = \alpha/a$. Other

possible inclusions are banana-shaped (imposing a curvature in one direction only, like BARdomain proteins do) or saddle-shaped (imposing opposite curvatures in the x and y directions). To determine the equilibrium shape of a flat membrane with point-like inclusions, we can use equation (6) for the Canham-Helfrich energy in the Monge gauge. To this energy we can add a number of Lagrange multipliers imposing the constraints, which are all of the form $(\partial_{ij}^2 u)\delta(\mathbf{x} - \mathbf{x}_p) = c$, where $\mathbf{x} = (x, y)$, \mathbf{x}_p is the position of the inclusion, and *i* and *j* can be either *x* or *y*. We can then calculate the change in the energy δE due to a small change in the shape δu , which needs to vanish for an equilibrium shape. This procedure gives us the shape equation for *u*:

$$\kappa \nabla_{\perp}^4 u - \sigma \nabla_{\perp}^2 u = \sum_a \Lambda_a D_a(\mathbf{x}), \tag{23}$$

where the elements of D_a are of the form $\partial_{ij}\delta(\mathbf{x} - \mathbf{x}_p)$ and the Λ_a are the Lagrange multipliers. As equation (23) is linear, it can be solved if we have its Green's function, which is given by [41]

$$G(\mathbf{x}) = \frac{1}{2\pi\sigma} \left[K_0(\lambda r) + \log(\lambda r) \right], \tag{24}$$

where $\lambda = \sqrt{\sigma/\kappa}$, $r = |\mathbf{x}|$, and K_0 is the zeroth order modified Bessel functions of the second kind. The shape is now given by $u(\mathbf{x}) = \sum_a \Lambda_a G_a(\mathbf{x})$, where the elements of G_a are of the form $\partial_{ij}G(\mathbf{x} - \mathbf{x}_p)$ [42]. Given the shape, the total energy of the membrane can be calculated easily. For two conical inclusions a distance R apart, it is given by:

$$E(R) = 2\pi\kappa\alpha_1\alpha_2(\lambda a)^2 K_0(\lambda R) + \pi\kappa(\alpha_1^2 + \alpha_2^2)(\lambda a)^2 K_2^2(\lambda R) + \dots$$
$$\approx -2\pi\kappa\alpha_1\alpha_2(ka)^2 \log(kR) + 4\pi\kappa(\alpha_1^2 + \alpha_2^2) \left(\frac{a}{R}\right)^4 + \mathcal{O}\left(\frac{1}{R^5}\right), \tag{25}$$

where α_1 and α_2 are the opening angles of the cones and *a* their radius. In the case that there is no tension, the interaction energy scales with the particle distance as $1/R^4$. In that case membrane-mediated interactions are weaker than electrostatic but stronger than van der Waals interactions. This term, which originates in the bending of the membrane, is always repulsive. In the presence of tension, we get an additional term which is repulsive if the conical inclusions are on the same side of the membrane, and attractive if they are on opposite sides, see figure 6a. For multiple inclusions, the same formalism can be used, though the expression for the minimum energy configuration grows quickly. Numerically minimizing the energy as a function of the inclusion positions is straightforward however. Two examples of equilibrium configurations of saddle-shaped inclusions are shown in figure 6b.

5.2 Interactions on closed surfaces

On intrinsically curved surfaces, such as tubes and spheres, membrane inclusions still experience membrane-mediated interactions. Compared to the flat case, there are two complicating factors that make the calculation of the interaction energies more difficult for closed systems. First, we can no longer use the Monge gauge description of the membrane in Cartesian coordinates. Fortunately, we can use similar expansions around cylindrical and spherical (or any well-defined equilibrium) shapes. However, for these cases the Green's functions are significantly more complicated than for the flat membrane. The second complication, which also feeds back into the Green's function, is that the surface is closed at a finite distance. Apart from





Fig. 6: Interactions on an asymptotically flat membrane. (a) Interaction energy of two identical conical inclusions as function of their distance R (scaled by the inclusion radius a). The inclusions touch at R/a = 2. In both plots the value of the interaction length is $1/\lambda = 2a$. Blue line: inclusions on the same side of the membrane ($\alpha_1 = \alpha_2 = \alpha$) always repel. Orange line: inclusions on opposite sides of the membrane ($\alpha_1 = -\alpha_2 = \alpha$) exhibit long-range attraction due to the tension term in the energy, and short-range repulsion due to the bending term in the energy. (b) Equilibrium configurations of multiple saddle-shaped inclusions. Small numbers of saddles spontaneously aggregate into rings, larger numbers build more complex structures.

appearing in the boundary conditions, the finite size of the surface also means that the inclusions can 'feel' each other in multiple directions, not just along the shortest line connecting them. For the simplest case, point-like inclusions on a tube, an analytical treatment is still possible. It turns out that the curved and closed nature of the shape fundamentally changes the interactions. Consequently, identical inclusions, that would repel on a flat surface, now attract each other in the angular direction. For multiple inclusions, the total energy is minimized when they form a ring around the tube. If there are more inclusions than fit into a ring, the remaining ones will build a second ring, and those rings themselves also interact, resulting in a stable separation, much like the rings observed in neuronal axons (figure 7a) [43].

For vesicles, which are curved in both directions, the analytical calculations become very difficult. Numerical results show that two identical inclusions on a sphere will attract, as also measured explicitly in an experiment with colloids adhered to a vesicle (figure 7b) [44]. For multiple inclusions, rings similar to those predicted in tubular membranes also form on both spherical [45] and elliptical [46] vesicles, where for prolate ellipsoidal shapes the ring spontaneously locates at the midplane, and induces a vesicle shape very similar to the snowman shape of a two-component vesicle (figure 7c).

5.3 Interactions in living systems

Ultimately, the goal of these is to not just describe biomimetic systems created in the lab, but configurations in living cells as well. There have been two recent successes in this direction. In 2013, EM imaging revealed that the endoplasmic reticulum has a 'parking garage' structure with multiple layers coupled by spiral ramps [47]. An adaptation of the flat-membrane approximation was used successfully to explain these shapes as a tradeoff between bending and surface tension contributions to the configuration of pairs of spiral proteins [48]. Very recently, constrictions of the tubular networks formed by mitochondria were observed to be closely coupled to the presence of a given protein, which can act as both a curvature sensor (in small concentrations) and curvature inducer (at large concentrations), an effect corroborated by numerical results [49].



Fig. 7: Interactions of inclusions on curved and closed membranes. (a) Point-like inclusions spontaneously form rings around tubular membranes. Identical rings experience long-ranged attraction and short-ranged repulsion (blue dashed line), resulting in a stable equilibrium distance. The red line shows the interaction energy for rings of opposite inclusions. Figure from [43]. (b) Interaction energy between two colloids adhered to a vesicle measured experimentally and determined numerically. Both results show long-range attraction with an energy minimum of about $3k_{\rm B}T$. Figure from [44]. (c) Numerical minimization of the energy of a prolate ellipsoidal vesicle containing multiple colloidal inclusions shows that the inclusions spontaneously form a ring around the vesicle's midplane. Figure from [46].

Appendices

A Some useful differential geometry

In order to calculate the Canham-Helfrich free energy for arbitrary membrane shapes, we need a way to calculate the area element and curvatures from a description of the membrane. In general, a two-dimensional surface embedded in three-dimensional space can be described by a function $\mathbf{r}(x_1, x_2)$, where x_1 and x_2 parametrize the surface. For example, for a sphere we could set $x_1 = \theta$, $x_2 = \phi$, and we have $\mathbf{r}(\theta, \phi) = R(\sin \theta \cos \phi, \sin \theta \sin \phi, \cos \theta)$. Given any point (x_1, x_2) on the surface, we can define a local coordinate system consisting of two tangents and one normal to the surface:

$$\mathbf{e}_{i} = \frac{\partial \mathbf{r}}{\partial x_{i}} \quad (i = 1, 2), \qquad \hat{\mathbf{n}} = \frac{\mathbf{e}_{1} \times \mathbf{e}_{2}}{|\mathbf{e}_{1} \times \mathbf{e}_{2}|}.$$
(26)

We can also define the surface's metric tensor (or 'first fundamental form') and curvature tensor (or 'second fundamental form'), the components of which are given by

$$g_{ij} = \mathbf{e}_i \cdot \mathbf{e}_j, \qquad L_{ij} = \frac{\partial^2 \mathbf{r}}{\partial x_i \partial x_j} \cdot \hat{\mathbf{n}} = \frac{\partial \mathbf{e}_i}{\partial x_j} \cdot \hat{\mathbf{n}} = -\mathbf{e}_i \cdot \frac{\partial \hat{\mathbf{n}}}{\partial x_j}.$$
 (27)

The last equality in the expression for L_{ij} is known as the Weingarten equation. The area element of our surface is related to the area element in parameter space through the Jacobian, which is simply the square root of the determinant of the metric:

$$\mathrm{d}A = \sqrt{\mathrm{det}(g_{ij})}\mathrm{d}x_1\mathrm{d}x_2. \tag{28}$$

The mean H and Gaussian K curvatures are the invariants of the curvature tensor, i.e., its trace and determinant. Because our surface is curved, when calculating these quantities we also need to take the metric into account; they are given by

$$H = -\frac{1}{2}g^{ij}L_{ij} = \frac{1}{2}\nabla \cdot \hat{\mathbf{n}}, \qquad K = \frac{\det(L_{ij})}{\det(g_{ij})},$$
(29)

where g^{ij} is the inverse of g_{ij} . Finally, the Laplacian (or Laplace-Beltrami operator) on the surface is given by

$$\Delta = \frac{1}{\sqrt{\det(g_{ij})}} \partial_i \left(g^{ij} \sqrt{\det(g_{ij})} \partial_j \right).$$
(30)

References

- [1] B. Alberts et al., Molecular biology of the cell (Garland Science, 2008).
- [2] Blausen.com staff, WikiJournal of Medicine 1, 10 (2014).
- [3] J. P. Reeves and R. M. Dowben, J. Cell Physiol. 73, 49 (1969).
- [4] F. Olson, C. A. Hunt, F. C. Szoka, W. J. Vail, and D. Papahadjopoulos, Biochim. Biophys. Acta 557, 9 (1979).
- [5] M. I. Angelova and D. S. Dimitrov, Faraday Discuss. Chem. Soc. 81, 303 (1986).
- [6] J. C. Stachowiak et al., Proc. Natl Acad. Sci. USA 105, 4697 (2008).
- [7] S. Ota, S. Yoshizawa, and S. Takeuchi, Angew. Chem. Int. Ed. 48, 6533 (2009).
- [8] S. Matosevic and B. M. Paegel, J. Am. Chem. Soc. 133, 2798 (2011).
- [9] S. Deshpande, Y. Caspi, A. E. C. Meijering and C. Dekker, Nat. Commun. 7, 10447 (2016).
- [10] P. B. Canham, J. Theoret. Biol. 26, 61 (1970).
- [11] W. Helfrich, Z. Naturforsch. C 28, 693 (1973).
- [12] M. Deserno, Chem. Phys. Lipids 185, 11 (2015).
- [13] G. R. Lázaro, I. Pagonabarraga, and A. Hernández-Machado, Chem. Phys. Lipids 185, 45 (2015).
- [14] M. Mutz and W. Helfrich, J. Phys. (France) 51, 991 (1990).
- [15] J. J. Pécréaux, H.-G. Döbereiner, J. Prost, J.-F. Joanny, and P. Bassereau, Eur. Phys. J. E 13, 277 (2004).
- [16] S. Semrau, T. Idema, L. Holtzer, T. Schmidt, and C. Storm, Phys. Rev. Lett. 100, 088101 (2008).
- [17] E. Evans and A. Yeung, Chem. Phys. Lipids 73, 39 (1994).
- [18] E. Evans, H. Bowman, A. Leung, D. Needham, and D. Tirrel, Science 273, 933 (1996).
- [19] I. Derényi, F. Jülicher, and J. Prost, Phys. Rev. Lett. 88, 238101 (2002).
- [20] T. Idema and C. Storm, Eur. Phys. J. E 34, 67 (2011).
- [21] G. Koster, A. Cacciuto, I. Derényi, D. Frenkel, and M. Dogterom, Phys. Rev. Lett. 94, 068101 (2005).
- [22] G. Koster, M. Van Duijn, B. Hofs, and M. Dogterom, Proc. Natl. Acad. Sci. USA 100, 15583 (2003).
- [23] C. Leduc et al., Proc. Natl. Acad. Sci. USA 101, 17096 (2004).
- [24] P. M. Shaklee, T. Idema, G. Koster, C. Storm, T. Schmidt, and M. Dogterom, Proc. Natl. Acad. Sci. USA 105, 7993 (2008).
- [25] K. A. Brakke, Exper. Math. 1, 141 (1992).

- [26] Ou-Yang Z.-C. and W. Helfrich, Phys. Rev. A 39, 5280 (1989).
- [27] M. M. Müller, M. Deserno, and J. Guven, Phys. Rev. E 72, 061407 (2005).
- [28] Hu J.-G. and Ou-Yang Z.-C., Phys. Rev. E 47, 461 (1993).
- [29] W.-M. Zheng and J. Liu, Phys. Rev. E 48, 2856 (1993).
- [30] F. Jülicher and U. Seifert, Phys. Rev. E 49, 4728 (1994).
- [31] F. Jülicher and R. Lipowsky, Phys. Rev. E 53, 2670 (1996).
- [32] R. Lipowsky, J. Phys. II (France) 2, 1825 (1992).
- [33] T. Baumgart, S. T. Hess, and W.W. Webb, Nature 425, 821 (2003).
- [34] S. Semrau, T. Idema, T. Schmidt, and C. Storm, Biophys. J. 96, 4906 (2009).
- [35] T. Ursell, W. S. Klug, and R. Phillips, Proc. Natl. Acad. Sci. USA 106, 13301 (2009).
- [36] T. Idema, S. Semrau, C. Storm, and T. Schmidt, Phys. Rev. Lett. 104, 198102 (2010).
- [37] K. Simons and E. Ikonen, Nature 387, 569 (1997).
- [38] A. R. Honerkamp-Smith, P. Cicuta, M. D. Collins, S. L. Veatch, M. den Nijs, M. Schick, and S. L. Keller, Biophys. J. 95, 236 (2008).
- [39] M. Goulian, R. Bruinsma, and P. Pincus, Europhys. Lett. 22, 145 (1993).
- [40] T. R. Weikl, M. M. Kozlov, and W. Helfrich, Phys. Rev. E 57, 6988 (1998).
- [41] A. R. Evans, M. S. Turner, and P. Sens, Phys. Rev. E 67, 041907 (2003).
- [42] P. G. Dommersnes and J.-B. Fournier, Biophys. J. 83, 2898 (2002).
- [43] A. Vahid and T. Idema, Phys. Rev. Lett. 117, 138102 (2016).
- [44] C. van der Wel, A. Vahid, A. Šarić, T. Idema, D. Heinrich, and D. J. Kraft, Sci. Rep. 6, 32825 (2016).
- [45] A. Šarić and A. Cacciuto, Phys. Rev. Lett. 108, 118101 (2012).
- [46] A. Vahid, A. Šarić, and T. Idema, Soft Matter 13, 4924 (2017).
- [47] M. Terasaki et al., Cell 154, 285 (2013).
- [48] J. Guven, G. Huber, and D. M. Valencia, Phys. Rev. Lett. 113, 188101 (2014).
- [49] S. C. J. Helle et al., eLife 6, e30292 (2017).

C 3 Cytoskeletal Filaments - Actinflaments, Microtubules and Intermediate Filaments

Sarah Köster Institut für Röntgenphysik Georg-August-Universität Göttingen

Contents

1	Introduction					
2	The	role of the different filament types	types 3			
	2.1	Actin filaments	3			
	2.2	Microtubules	4			
	2.3	Intermediate filaments	4			
	2.4	Interactions between the filament types	4			
3	Dynamic assembly of filaments					
	3.1	Nucleotide-based polymerization	5			
	3.2	Self-assembly into filaments	8			
	3.3	Polyelectrolyte nature of protein filaments	9			
4	Mechanics of cytoskeletal filaments					
	4.1	Filament bending	10			
	4.2	Filament stretching	13			
5	Summary					

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Biological cells and tissues are extremely well adapted to the mechanical challenges they face in the context of the respective tissue or organ. For example, muscle cells are resistant to tension, epithelial cells bear strong shear forces and neurons are exposed to hydrostatic pressure. It is well accepted that the so called cytoskeleton, a composite intracellular biopolymer network plays a decisive role for cell mechanics (see chapter D7). Thus, the cytoskeleton is of high importance for maintaining cell shape and enabling migration, contraction, and division. Indeed, much of the diversity among the approximately 200 different cell types in the human body is due to the cytoskeleton. The main components of the cytoskeleton are three distinct protein filament systems, namely actin filaments, microtubules and intermediate filaments, which will be introduced and discussed in this chapter. The biological basics are found in more detail in cell biology textbooks, see, *e.g.* Refs. [1, 2]. It is a key feature of the cytoskeleton that the three types of filaments, which have distinctly different physical properties, can be combined in many different ways and thus provide a "construction kit" to meet the cell's mechanical needs. At the same time, however, all filament types share common traits, in particular from a polymer physics point of view.



Fig. 1: Fluorescence micrographs and sketches of the three filaments systems in the cytoskeleton of 3T3 fibroblasts: a) Actin filaments (red), b) microtubules (blue), c) intermediate filaments, in this case vimentin, (green) and d) composite image of all three filament types. The distinct structure of each individual network hints at differing mechanical properties of each filament type. All scale bars are 20 µm. Images recorded by Ulrike Rölleke.

Fig. 1 shows fluorescence micrographs of each of the three filaments types in a fibroblast cell. The actin filaments (shown in red) were labeled by a phalloidin conjugate (see below), the microtubules (shown in blue) and vimentin intermediate filaments (shown in green) by specific antibodies conjugated to fluorescent labels. The nucleus is not stained in this case but can easily be assigned to the round void space in the cell center (indicated in orange in the sketches). The images show the distinctly different structures formed by the filaments in cells. For actin, we observe parallel bundles (so-called stress-fibers), a "cortex" close to the cell membrane (see chapter C1 and C2) and cellular protrusions like filopodia. Tubulin forms a cell-spanning system of "tracks", originating from a joint "microtubule organizing center" (MTOC). Vimentin

intermediate filaments, finally, organize into networks with a small mesh size. The mere appearance of these distinct biopolymer networks gives important hints at their specific role and function in the cell. At the same time, however, it is clear that the three networks are mechanically integrated, interact with each other and even overlap in certain regions within the cytoplasm. Cross-linked by associated proteins and driven by dynamic molecular motors (see chapter C5 and C6), which are not shown in Fig. 1, the cytoskeleton is a prime example for active biological matter. Importantly, it is the *emergent* properties of the *combined* networks and the adaptability that makes the cytoskeleton such an intriguing viscoelastic biological material.



2 The role of the different filament types

Fig. 2: Schematic showing the interactions of the three filament types in eukaryotes, mediated by cross-linkers and motor proteins. From Ref. [3].

2.1 Actin filaments

Actin filaments, or microfilaments, (see Fig. 1a) are responsible for defining and maintaining cell shape, for cell locomotion and contraction. They are involved in different structures in the cell, as shown in Fig. 2, lower right: (i) In the cortex, the filaments are cross-linked at large angles and form a gel-like structure that is located just underneath the fragile plasma membrane and preservers it. (ii) In stress fibers, filaments are arranged in an anti-parallel manner and cross-linked by α -actinin, which allows for directed myosin motors to enter the fibers. (iii) In lamellipodia and filopodia, the filaments are tightly packed in a parallel manner, crosslinked by fimbrin, and push the plasma membrane outwards at the leading edge of the cell, thus enabling migration and shape changes. Actin is a globular, 42 kDa protein which is highly conserved, with about 90% similarity in the amino acid sequence between species. There are three isoforms, α -actin, which is found in muscle cells, and β - and γ -actin, which are found in non-muscle cells. However, for our considerations in the context of this chapter, we will not distinguish between the isoforms, but regard them together.

2.2 Microtubules

Microtubules (see Fig. 1b) are stiff, hollow cylinders, as suggested by their name, and play a pivotal role for transport of organelles and vesicles within the cell, where they serve as "tracks", and during cell division, where they pull the chromosomes apart. These mirotubule tracks start from the MTOC and explore the whole cell. Transport along microtubules is enabled by directed dynein and kinesin motors, to which vesicles bind as cargo. In the mitotic spindle, microtubules and motors act together during cell division to pull the chromosomes into the daughter cells. Microtubules are also found inside cilia, flagella, and microvili – cellular appendices, which cells use to move in fluids (see chapter E4). Like actin, tubulin is highly conserved with 75 % similarity even between yeast and human. Two closely related, 55 kDa proteins, α - and β -tubulin, form heterodimers, which then further assemble into the microtubule (see below).

2.3 Intermediate filaments

Intermediate filaments have a diameter of around 10 nm, which lies between actin filaments (7 nm) and microtubules (25 nm), therefore their name. Unlike actin filaments and microtubules, however, intermediate filaments are not conserved, but expressed in a cell-type specific manner with further differences between species. In humans, about 70 different genes encode for the family of intermediate filament proteins [4]. All members of this family share two important features: (i) The rod-shaped monomers, which, (ii) self-assemble in a hierarchical manner into extended filaments. Prominent examples are vimentin, which is expressed in cells of mesenchymal origin, desmin, found in muscle cells, and keratin as the main intermediate filament protein in epithelial cells¹. The molecular weights differ between the specific proteins and are roughly 50 kDa. Intermediate filaments provide cells with mechanical strength and resistance to shear. Interestingly, they are the only cytoskeletal filaments which do not display polarity and thus posses no directly associated motor proteins.

2.4 Interactions between the filament types

As shown in Fig. 2 the three filaments types are closely interlinked by motor proteins (myosin, kinesin, dynein) and by associated proteins. The list of possible factors and interactions is long and complex (see Ref. [3] and references therein): some proteins link the same type of

¹This section focuses on cytoskeletal intermediate filaments. It should be noted, however, that intermediate filaments also occur in the nuclear lamin (lamin A and B).

filaments to one another (*e.g.* α -actinin, filamin and fascin for actin, or some distinct microtubule associated proteins (MAPs) for tubulin), others link different filaments types together (*e.g.* plectin which binds to all three cytoskeletal filaments). Additionally, as explained further below, the charged filaments may also be cross-linked by the addition of counter-ions. Together with unavoidable steric interactions within the cytoplasm, including entanglement of the filaments, these active and passive interactions give rise to the mechanics and dynamics of cells and tissues.

3 Dynamic assembly of filaments

3.1 Nucleotide-based polymerization

Actin monomers are globular (G-actin) with a diameter of about 5 nm and a molecular weight of 42 kDa. During assembly, the monomers bind head-to-tail, thus resulting in polar filaments (F-actin, microfilaments or actin filaments) with two distinct ends that are called the plus- and the minus end (or, historically, the "barbed" and the "pointed" end, because the filaments look like arrows, when decorated by myosin heads).² The fully assembled filament resembles a double helix of two protofilaments that are twisted around each other (see Fig. 3c) in a double helix with 37 nm pitch. These filaments have a diameter of merely 7 nm and can be many μ m long. This enormous aspect ratio is shared among all cytoskeletal filaments. From a physics point of view, they can thus be "corse grained" and regarded as biopolymers, including the whole wealth of polymer physics descriptions available. Note, however, that the bonds between the monomers are not covalent, but rather hydrophobic interactions and other non-covalent interactions. Thus, the filaments may assemble and disassemble in a highly dynamic manner, thus adapting to the cellular environment.

Actin assembly may be divided into three distinct phases, as depicted schematically in Fig. 3a. During nucleation, or lag phase, a certain number of monomers (typically three for actin) bind and form a stable nucleus. Subsequently, during elongation, or growth phase, long filaments extend from these nuclei, before finally the steady state, or equilibrium phase, is reached, where the length stays mostly constant. Owing to the three-phase assembly process, the length of the filaments depends on the number of nuclei at the beginning (the more nuclei, the more, but shorter filaments) and thereby on protein concentration and assembly time.

Actin filaments are not static structures. Even during equilibrium phase, monomers are constantly added and released from the filaments, however in a steady-state manner, such that the filament length remains constant. Free actin monomers, G-actin, bind one ATP (adenosinetriphosphate) molecule, which is hydrolyzed to ADP (adenosinediphosphate) upon binding. Thus, to be precise, we have to consider separate binding and unbinding constants for ATPand ADP-actin and for the plus and the minus end, as summarized in Tab. 1. In general, rates for the plus end are higher, thus both assembly and disassembly occur faster than for the minus end. ATP-actin plays a more important role for assembly, while ADP-actin, which results from hydrolysis of the ATP bound to the actin monomers upon integration into the filament, is prone to disassembly. In equilibrium, *i.e.* when assembly and disassembly occur at exactly the same

²It should be noted that the terms "plus" and "minus" do not denote charges associated with the molecules, but merely a faster and a more slowly assembling filament end.



Fig. 3: *a)* Schematic representation of the assembly of actin monomers into filaments, including lag phase, growth phase and equilibrium phase. b) Phenomenon of treadmilling. There are as many subunits added as are released per unit time, thus the filament length remains constant. c) *Fully assembled, helical actin filament with a diameter of 7 nm and a pitch of 37 nm. Graphics from Ref. [1].*

rate and the filament length thus stays constant, we obtain

$$\frac{\mathrm{d}n^{+,\mathrm{on}}}{\mathrm{d}t} = k_{\mathrm{on}}^+[C_c] \tag{1}$$

for the assembly at the plus end and

$$\frac{\mathrm{d}n^{+,\mathrm{off}}}{\mathrm{d}t} = -k_{\mathrm{off}}^{+} \tag{2}$$

for the dissassembly at the plus end, and thus in total:

$$\frac{\mathrm{d}n^+}{\mathrm{d}t} = k_{on}^+ [C_c]^+ - k_{\rm off}^+ = 0.$$
(3)

Thus the critical concentration of free monomers (G-actin) in equilibrium is

$$[C_c]^+ = \frac{k_{\text{off}}^+}{k_{\text{on}}^+}.$$
(4)

 $[C_c]^+$ is what remains as a soluble pool in steady state phase and what is required as a minimum concentration for assembly to start.

Tab. 1 lists critical concentrations for the plus and minus end. Obviously, in general $[C_c]^+ \neq [C_c]^-$. Thus, for G-actin concentrations $[C_s]$, such that $[C_c]^+ < [C_s] < [C_c]^-$ we observe that
monomer	$k_{\rm on}^+$	$k_{\rm off}^+$	$k_{\rm on}^-$	$k_{\rm off}^-$	$[C_c]^+$	$[C_c]^{-}$
ATP-actin	$11.6 (\mu M s)^{-1}$	1.4 s^{-1}	$1.3 \ (\mu M \ s)^{-1}$	$0.8 \ {\rm s}^{-1}$	0.12 μM	$0.6 \ \mu M$
ADP-actin	$3.8 \ (\mu M \ s)^{-1}$	7.2 s^{-1}	$0.16 \ (\mu M \ s)^{-1}$	0.27 s^{-1}	1.9 μM	$1.7 \ \mu M$
GTP-tubulin	8.9 $(\mu M s)^{-1}$	44 s^{-1}	$4.3 \ (\mu M \ s)^{-1}$	23 s^{-1}	4.9 μM	$5.3 \ \mu M$
GDP-tubulin	$0 (\mu M s)^{-1}$	733 s^{-1}	$0 (\mu M s)^{-1}$	915 s^{-1}	n.a.	n.a.

Table 1: Binding and unbinding constants for the plus and the minus end of actin filaments and microtubules as well as for ATP-/ADP-actin and GTP-/GDP-tubulin. Note the differing units for the assembly $(M \ s^{-1})$ and disassembly (s^{-1}) rates since the polymerization depends on the free monomer concentration, whereas the depolymerization does not. Numbers taken from Ref. [5].

the plus end growths and the minus end shrinks. If the growth and the shrinking occur at the same rate, *i.e.*

$$\frac{\mathrm{d}n^+}{\mathrm{d}t} = -\frac{\mathrm{d}n^-}{\mathrm{d}t} \tag{5}$$

we obtain the steady state concentration $[C_s]$

$$[C_s] = \frac{k_{\text{off}}^+ + k_{\text{off}}^-}{k_{\text{on}}^+ + k_{\text{on}}^-}.$$
(6)

This phenomenon (see Fig. 3b) is called "treadmilling" since, similar to a treadmill in the gym, as many monomers are added to the plus and of the actin filament as are removed from the minus end.

In many *in vitro* experiments, the dynamic assembly and disassembly of actin filaments is undesired, for example when determining mechanical properties, as described in the following section. An easy way to avoid these dynamics is the addition of phalloidin, the toxin of *Amanita phalloides*, the death cap mushroom, which binds and stabilizes actin filaments. Phalloidin can also be conjugated to fluorescent labels and thus be used to stain F-actin specifically (see Fig. 1a). For cell experiments, latrunculin, the toxin of the sea sponge *Latrunculia magnifica* is frequently employed, which binds and stabilizes monomers, thus decreasing the number of filaments, as well as cytochalasin, a fungal metabolite, which caps the filament plus end. Both these latter drugs can thus be used to interfere with the actin cytoskeleton in cell experiments.

The assembly of tubulin into microtubules shares many characteristics with actin assembly. Here as well, small subunit (dimers of α and β -tubulin with a 55 kDa each) bind head-to-tail, thus assembling into polar filaments. By contrast to the helical actin filaments, microtubules are hollow cylinders (see Fig. 4a), made of 13 (in rare cases 11 or 15) longitudinal units, arranged in parallel as "protofilaments". In these protofilaments, each dimeric repeat is 8 nm long. Like actin filaments, microtubules posses a plus- and a minus-end and the plus end is favored for polymerization (see Tab. 1). Free tubulin dimers bind a guanosinetriphosphate (GTP) molecule in each monomer and the GTP molecule bound to the β -tubulin hydrolyzes into guanosinediphosphate (GDP) upon integration into the microtubule.

Our discussion on critical concentrations and treadmilling is valid for tubulin polymerization as well. Additionally, however, for microtubules, phases of very rapid disassembly (called "catastrophe") are observed between phases of growth. This phenomenon is termed "dynamic instability" and is related to GTP hydrolysis. Tab. 1 shows that GTP-tubulin favors growth,



Fig. 4: *a)* Fully assembled, cylindrical microtubule with a diameter of 25 nm and a dimer size of 8 nm. b) Phenomenon of catastrophe (rapid disassembly). Graphics from Ref. [1].

whereas GDP-tubulin tends to disassemble³. Thus, as long as GTP-tubulin is bound to the end of the microtubule, but not hydrolyzed to GDP yet, it forms a so-called GTP cap, that stabilizes the microtubule. If this cap is lost, however, the microtubule rapidly shrinks (see Fig. 4b), until rescued by a new GTP cap. Just like for actin, there are drugs, which interfere with tubulin assembly and disassembly. Paclitaxel, the poison of the pacific yew, stabilizes microtubules and thus interferes with chromosome segregation and cell division. For this reason, it is also used in cancer therapy. For *in vitro* microscopy experiments, paclitaxel conjugated to fluorescent dyes is widely used. By contrast, nocodazole interferes with microtubule polymerization and is employed for testing the effect of the absence of microtubules.

3.2 Self-assembly into filaments

By contrast to actin filaments and microtubules, intermediate filaments do not polymerize (*i.e.* addition of monomers one by one) in a strict way, but self-assemble in a hierarchical manner [6]. The intermediate filament monomer is rod shaped with a structured, α -helical rod domain in the center flanked by intrinsically disordered head and tail domains. The monomers form parallel dimers and antiparallel tetramers at physiological conditions. Upon the addition of monovalent ions, the tetramers start to assemble first laterally into unit length filaments (ULFs) and then elongate into extended filaments with a diameter of about 10 nm and a length of many μ m (see Fig. 5a). The hierarchical structure of the filaments, where many monomers are arranged per cross-section (typically 32 for vimentin and 16 for keratin) gives rise to intriguing mechanical properties, such as high flexibility and enormous stretchability [7, 6].

Furthermore, the distinct architecture of intermediate filaments gives rise to phenomena such as subunit exchanged from mature filaments, where a filament subunit like a tetramer or an octamer leaves the filaments or is integrated in the filaments (Fig. 5c). Subunit exchange has been observed in cells [8] and *in vitro* [9]. The elongation reaction also allows for the annealing of ULFs, or of filaments of any length [10]. Thus, end-to-end annealing of long filaments does occur (Fig. 5b), albeit at a small rate that is associated to the rotational diffusion coefficient of long filaments. Notably, there are no nucleotides involved in intermediate filament assembly and the resulting filaments are not polar (*c.f.* the antiparallel organization of the dimers in

³We assume that rapid disassembly is associated with GDP-tubulin.



Fig. 5: *a)* Assembly path way of intermediate filaments. The rod shaped monomers first form dimers, tetramers and ULFs in a lateral manner. Subsequently, the ULFs associate longitudinally and form extended, 10 nm diameter filaments. In vitro, the assembly reaction can be initiated by the addition of physiological concentrations of monovalent salts, such as KCl or NaCl. b) End-to-end annealing of two fully assembled filaments. c) Subunit exchange in a fully assembled filament. Graphic from Ref. [6]

a tetramer). Thus, there are no motors known, that move along intermadiate filaments like myosin, kinesin and myosin do on actin filaments and microtubules.

3.3 Polyelectrolyte nature of protein filaments

Like all proteins, cytoskeletal filaments are polypeptide chains built up from amino acids. These amino acids may either be uncharged, positively charged, or negatively charged. Consequently, the polypeptide chain carries a net charge and a distinct charge pattern that is transferred to the assembled filament (see Fig. 6a) and complemented by hydrophobic regions arranged along the filament. Thus, protein filaments may be regarded as polyelectrolytes [11]. Interestingly, as shown in Fig. 6b, the intracellular biological filaments (cytoskeletal filaments and DNA) carry a high negative net charge, whereas extracellular filaments (fibrin and collagen) are not very highly charged. Fig. 6c shows the direct observation of two vimentin intermediate filaments binding to each other in the presence of MgCl⁺ ions. For intermediate filaments, in particular, charge interactions play an important role in assembly and disassembly: during cell division, the highly stable filaments become phosphorylated, which adds extra negative charge to specific amino acids, thus disassemble, and reassemble after the formation of two daughter cells.



Fig. 6: *a)* Surface potential of actin and vimentin filaments (as well as DNA and Pf1 virus). *b)* Negative surface charge for different biological filaments. c) Micrographs (fluorescence, inverted gray scale) of vimentin intermediate filaments in the precence of Mg^{2+} ions. The time lapse shows direct binding of two filaments due to charge interactions. a,b from Ref. [11]; c from Ref. [12].

4 Mechanics of cytoskeletal filaments

4.1 Filament bending

As discussed in the introduction, cytoskeletal filaments are mechanical elements of the cell. Thus, determination of their mechanical properties is thus of great value. Here, we discuss individual filaments, whereas in chapter E1 rheological methods are introduced, which are frequently used to study the mechanics of filament networks. In principle, we can distinguish between two modes of mechanic impact on the filament: (i) bending and (ii) stretching⁴. The bending rigidity κ of a filament is directly proportional to the so-called persistence length L_P , scaled by k_BT :

$$L_P = \kappa / (k_B T). \tag{7}$$

Biological filaments may be grouped into (i) stiff rods for $L_P >> L$, (ii) semi-flexible polymers, where $L_P \simeq L$, and (iii) flexible polymers, where $L_P << L$. Examples are microtubules with $L_P \simeq 1$ mm, actin filaments ($L_P \simeq 15 \ \mu m$)⁵ as well as intermediate filaments ($L_P \simeq 1 \ \mu m$) and DNA with $L_P \simeq 50$ nm, respectively, assuming a length of 10 to 20 μm for each of them. Fig. 7a-c shows examples of all four biological filaments as sketches and in fluorescence micrographs. Flexible polymers can be described by a (self-avoiding) random walk.

⁴We could also regard filament twisting, which is, however, beyond the scope of this chapter.

⁵Note that this value was determined for phalloidin-stabilized actin filaments; without the stabilization, the persistence length is a factor of about 2 smaller [13].



Fig. 7: *a-c)* From left to right, biopolymers with increasing persistence length, i.e. flexible, semiflexible and stiff. Below: fluorescence micrographs of DNA, vimentin intermediate filaments, actin filaments and microtubules, respectively. The DNA molecules coil up so strongly that in the fluorescene image only small dots are visible. d) Schematic, determination of the persistence length of a filament or polymer according to the worm-like chain model.

For cytoskeletal filaments, however, this model is not valid and instead of discrete, independent polymer segments, a continuous model, the so-called worm-like chain model is typically applied [14]. For a filament in equilibrium, *i.e.* no length changes and no external forces, which is thermally fluctuating, the persistence length is determined as sketched in Fig. 7d. The contour is parametrized by the arc length s and the Hamiltonian, which contains merely the bending energy (that is, we assume non-extendable filaments and ignore twisting or rotation), is

$$H = \int_0^L \mathrm{d}x \left[\frac{\kappa}{2} \frac{\mathrm{d}^2 \vec{r}(s)}{\mathrm{d}s^2}\right]^2. \tag{8}$$

From this simple Hamiltonian, the tangent correlation can be calculated analytically:

$$\langle \vec{t}(s) \cdot \vec{t}(s-l) \rangle = \langle \cos \theta(l) \rangle = \exp\left(-\frac{l}{L_P}\right),$$
(9)

hence the term "persistence length": fluctuations decay with L_P , the length over which they are persistent. A large L_P corresponds to stiff filaments, and a small one to flexible polymers. The mean squared end-to-end distance $\langle \vec{R}^2 \rangle$ amounts to

$$\langle \vec{R}^2 \rangle = \left\langle \int_0^L \int_0^L \mathrm{d}l_1 \mathrm{d}l_2 \vec{t}(l_1) \cdot \vec{t}(l_1) \right\rangle = 8L_P^2 \left(\frac{L}{2L_P} - 1 + \exp\left(-\frac{L}{2L_P}\right) \right). \tag{10}$$

For the case of very flexible ($L_P << L$) or very stiff ($L_P >> L$) filaments, Eq. 10 simplifies to $\langle \vec{R}^2 \rangle \approx 2LL_P$ and $\langle \vec{R}^2 \rangle \approx L^2$, respectively.

Filaments may be fluorescently labeled and directly observed by fluorescence microscopy [15, 16, 17, 18] (see chapter A1). This approach offers time and ensemble averaging of the filament contours and thus statistically valid values. Moreover, if the filament is equilibrated in solution at constant temperature, external influences, such as filament-surface interactions are avoided. Further methods employed for determination of the persistence length are atomic force microscopy (AFM), electron microscopy (EM, see chapter A2), rheology (see chapter E1), mechanical stretching (see below), and light scattering (see chapter A4).

The persistence length of intermediate filaments has been measured in analogy to actin filaments. Interestingly, the filaments assembled from different intermediate filament proteins have different persistence lengths, varying by about one order of magnitude.

Actin filaments ($L_P \simeq 15 \ \mu$ m) and intermediate filaments $L_P \simeq 1 \ \mu$ m) differ in persistence length by one order of magnitude, despite similar diameter. This difference is most likely due to the strongly differing molecular structure: as described above, the actin filament is a double helix of globular monomers, whereas the intermediate filament consists of rod-shaped monomers, which are arranged in parallel, resulting in a much more "open" structure. As detailed below, however, it is important to note that the values for L_P , which are measured in *in vitro* settings are relevant in the cell in a more indirect way: by cross linking and bundling much stronger structures are created than the individual filaments could ever provide.

The bending rigidity κ is directly related to the Young's modulus E of a filament as

$$\kappa = EJ,\tag{11}$$

where J is the geometrical moment of inertia. For cylindrical rods,

$$J = \frac{\pi R^4}{4},\tag{12}$$

whereas for hollow cylinders

$$J = \frac{\pi (R_a^4 - R_i^4)}{4}.$$
 (13)

Thus,

$$L_P = \frac{EJ}{k_B T} = \frac{\pi E R^4}{4k_B T} \tag{14}$$

scales with the radius to the fourth power and the determination of the persistence length via this relation depends on a very accurate measurement of the radius.

The persistence length of microtubules (\simeq mm) is much larger than the size of a typical cell (10 - 100 μ m) and thus, these filaments are typically mostly straight in the cell and buckle under mechanical load. As microtubules in cells are embedded in the viscoelastic cytoplasm, their buckling behavior differs from classical Euler buckling (see Fig. 8a, top). The competition between bending energy (of the rod) and elastic deformation energy (of the matrix) leads to a buckling wavelength λ ,

$$\lambda = 2\pi \left(\frac{\kappa}{\alpha}\right)^{1/4},\tag{15}$$

where κ is the bending rigidity of the microtubule and α is proportional to the shear modulus of the matrix, see Fig. 8a, bottom. Fig. 8b shows such buckling events in microtubules that occur spontaneously in cells [19]. Bundles of keratin filaments show the same phenomenon (Fig. 8c), however, in this case, the persistence length, and thereby bending rigidity of the bundle depends on how strongly the individual filaments (see electron micrograph in Fig. 8d) are coupled to each other. If N ins the number of filaments in the bundle,

$$L_{P_{\text{bundle}}} = N^{\gamma} L_{P_{\text{filament}}},\tag{16}$$

with $\gamma = 1$ for uncoupled filaments and $\gamma = 2$ for fully coupled filaments. Using this analysis method, the persistence length of keratin bundles amounts to about 1 mm with strong coupling between the individual filaments [20].



Fig. 8: a) Compression of a rod in viscous (top) or elastic (bottom) environment. b) Buckling events in microtubules (scale bar 5 μ m) and c) in keratin bundles. d) Electron micrographs of transversal and longitudinal cross-sections of keratin bundles. a,b from Ref. [19]; c,d from Ref. [20].

4.2 Filament stretching

Apart from bending, filaments may be mechanically manipulated by stretching. The non-linear and highly complex stress-strain behavior of cytoskeletal filaments reveals intriguing properties that hints at their important physiological roles. Fig. 9a shows an early example [21] that compares the three filament types as well as fibrin, an important blood protein that is instrumental for forming blood clots. Interestingly, both microtubules and actin filaments (F-actin) display strictly linear stress-strain behaviors and rupture at comparatively low strains of 80% and 20%, respectively. By contrast, fibrin and vimentin intermediate filaments do not rupture in the setting of this particular experiment. These experiments were performed using rheology (see chapter E1) and thus on entangled, non-cross-linked networks of the respective filament types.

A more direct way of determining the stress-strain behavior of filament is by pulling individual filaments by using optical traps of atomic force microscopy (AFM). An examples of such an experiment is shown in Fig. 9b. Here, individual vimentin intermediate filaments were stretched at different loading rates (pulling velocities) and the data reveal more detail than the rheology data [21], which probe a full polymer network. Due to the molecular architecture (see above, Fig. 5) of the intermediate filament with mostly α -helices in the monomer and multiple monomers arranged per cross-section, different regimes are observed: (i) At low strains and stresses, we observe a linear regime, where the slope of the curve gives rise to the elasticity and thereby persistence length of the filament (see Eq. 14). (ii) A plateau regime follows, where with comparatively low stress, the filaments may be pulled much further. In this regime presumably the α -helices transition into β -sheets, thereby elongating [22, 23, 24, 25, 26]. (iii) Finally, the filaments stiffen and enter a linear regime again, supposedly due to pulling on the β -sheets. AFM



Fig. 9: *a)* Stress-strain behavior of the three types of cytoskeletal filaments (and fibrin in addition), measured by rheological methods, from Ref. [21]. b) Stress-strain behavior of individual intermediate filaments measured using optical traps, from Ref. [22].

experiments, where the applied force reaches one order of magnitude more than in the optical trap experiments shown in Fig. 9b show additional regimes and strains up to at least least 3.5, *i.e.* to 4.5 times their length, without breaking, which is quite remarkable. Interestingly, the plateau region is more extended for smaller loading rates. This complex stress-strain behavior of intermediate filaments is special among the three types of cytoskeletal filaments and hints at an important mechanical role of the filaments in the cell: The filaments act much like a "safety belt" that can easily be pulled at small rates and small forces, but stiffens rapidly upon high rates and forces [26]. Thereby, intermediate filaments are mechanically "invisible" in motile cells, but are still able to protect cells against severe impact.

5 Summary

To summarize, the cytoskeleton of eukaryotes comprises three protein filament types with distinctly different mechanical properties, both in terms of bending stiffness and concerning stressstrain behavior. By combining bundles and networks of these filaments in a tailored way and employing passive cross-linkers and active motor proteins, the cell adapts to the specific mechanical challenges in its tissue context. Actin filaments and microtubules are highly dynamic: they possess motor systems, with which they interact directly, and they constantly assemble and disassemble in the cell. Thus, actin filaments are a key player in cell motility and contraction, whereas microtubules play a leading role in intracellular transport and in chromosome segregation during cell division. Microtubules are hollow tubes with a diameter of 25 nm and are thus stiff on the length scale of a cell. Actin filaments are double helices build up from globular monomers and possess a persistence length on the order of the size of a cell. As a consequence of their molecular architecture, both filaments types are barely extendable and rupture at comparatively low strains. Intermediate filaments differ from the other two filament types: they are non-polar and consequently no motors are associated with them. The are open structures, leading to a small persistence length but enormous extensibility up to at least 4.5 times their original length.

References

- B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts, and P. Walter, *Molecular Biology* of the Cell (Garland Science, New York, 2014).
- [2] H. Lodish, A. Berk, C. A. Kaiser, M. Krieger, M. Scott, A. Bretscher, H. Ploegh, and P. Matsudaira, *Molecular Cell Biology* (W. H. Freeman, New York, 2007).
- [3] F. Huber, A. Boire, M. P. Lopez, and G. Koenderink, Current Opinion in Cell Biology 32, 39 (2015).
- [4] I. Szeverenyi, A. J. Cassidy, C. W. Chung, B. T. K. Lee, J. E. Common, S. C. Ogg, H. Chen, S. Y. Sim, W. L. P. Goh, K. W. Ng, J. A. Simpson, L. L. Chee, *et al.*, Human Mutation 29, 351360 (2007).
- [5] D. Boal, Mechanics of the Cell (Cambridge University Press, Cambridge, 2012).
- [6] J. Block, V. Schroeder, P. Pawelzyk, N. Willenbacher, and S. Köster, Biochimica Biophysica Acta: Molecular Cell Biology 1853, 3053 (2015).
- [7] S. Köster, D. Weitz, R. Goldman, U. Aebi, and H. Herrmann, Current Opinion in Cell Biology 32, 82 (2015).
- [8] G. Çolakoğlu and A. Brown, Journal of Cell Biology 185, 769 (2009).
- [9] B. Nöding, H. Herrmann, and S. Köster, Biophysical Journal 107, 2914 (2014).
- [10] S. Winheim, A. R. Hieb, M. Silbermann, E.-M. Surmann, T. Wedig, H. Herrmann, J. Langowski, and N. Mücke, PloS One 6, e19202 (2011).
- [11] P. A. Janmey, D. R. Slochower, Y.-H. Wang, and A. Ceber, Soft Matter 10, 1439 (2014).
- [12] C. Dammann, H. Herrmann, and S. Köster, Israel Journal of Chemistry 56, 614 (2016).
- [13] H. Isambert, P. Venier, A. C. Maggs, A. Fattoum, K. Ridha, D. Pantaloni, and M.-F. Carlier, Journal of Biological Chemistry 270, 11437 (1995).
- [14] O. Kratky and G. Porod, Receuil des Traveaux Chimiques des Pays-Bas 68, 1106 (1949).
- [15] L. Le Goff, O. Hallatschek, E. Frey, and F. Amblard, Physical Review Letters 89, 258101 (2002).
- [16] G. Frederick, B. Mickey, J. Nettleton, and J. Howards, Journal of Cell Biology 120, 923 (1993).
- [17] J. Käs, H. Strey, J. X. Tang, D. Finger, R. Ezzell, E. Sackmann, and P. A. Janmey, Biophysical Journal 70, 609 (1996).
- [18] A. Ott, M. Magnasco, A. Simon, and A. Libchaber, Physical Review E 48, R1642 (1993).
- [19] C. P. Brangwynne, F. C. MacKintosh, S. Kumar, N. A. Geisse, J. Talbot, L. Mahadevan, K. K. Parker, D. E. Ingber, and D. A. Weitz, Journal of Cell Biology **173**, 733 (2006).
- [20] J.-F. Nolting, W. Möbius, and S. Köster, Biophysical Journal 107, 2693 (2014).
- [21] P. A. Janmey, U. Euteneuer, P. Traub, and M. Schliwa, Journal of Cell Biology 113, 155 (1991).
- [22] J. Block, H. Witt, A. Candelli, E. J. G. Peterman, G. J. L. Wuite, A. Janshoff, and S. Köster, Physical Review Letters 118, 048101 (2017).
- [23] L. Kreplak, H. Bär, J. F. Leterrier, H. Herrmann, and U. Aebi, Journal of Molecular Biology 354, 569 (2005).
- [24] L. Kreplak, H. Herrmann, and U. Aebi, Biophysical Journal 94, 2790 (2008).
- [25] C. Guzman, S. Jeney, L. Kreplak, S. Kasas, A. J. Kulik, U. Aebi, and L. Forro, Journal of Molecular Biology 360, 623 (2006).
- [26] Z. Qin, L. Kreplak, and M. J. Buehler, PLoS One 4, e7294 (2009).

C 4 Theory of Semiflexible Network Materials

Cornelis Storm Theory of Polymers and Soft Matter group Department of Applied Physics Eindhoven University of Technology

Contents

1	Introduction	2
2	Worm-Like and Semiflexible Chains: Model Definitions	2
3	Nonlinear Force-Extension of a Single Semiflexible Chain	3
4	Static Response of an Affinely Deforming Semiflexible Network	6
5	Dynamic Response of Semiflexible Chains and Networks	10
A	Fourier transform convention	14

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

In this lecture, I summarize the theory of semiflexible filaments and their networks. I will begin by deriving the static elastic response of a single semiflexible chain, and networks composed of such chains, to full nonlinear order. Then, I address the linear dynamical response of single polymers and networks to compute the frequency-dependent rheological response. In preparing these notes, I have followed derivations and lines of reasoning in the excellent textbook by Doi [1] and a seminal paper by Gittes and MacKintosh [2]. The discussion on network energies and stress tensors provides background to [3].

2 Worm-Like and Semiflexible Chains: Model Definitions

Consider a polymeric (line-like) object, with a total contour length ℓ_c [L]. At any point t [T] in time, it traces out a space curve parameterized by the arc length s [L] as $\mathbf{r}(s,t) = \{x(s,t), y(s,t), z(s,t)\}$. Inextensibility of the backbone is ensured by enforcing that the tangent vector t [dimensionless] has unit length, everywhere:

$$|\mathbf{t}| = \left|\frac{\partial \mathbf{r}}{\partial s}\right| = 1.$$
(1)

The Worm-Like Chain model assigns to such a conformation an elastic energy, given by the square local curvature, integrated along the entire chain:

$$\mathcal{H}_{WLC} = \frac{\kappa}{2} \int_0^{\ell_c} \mathrm{d}s \, \left| \frac{\partial^2 \mathbf{r}}{\partial s^2} \right|^2 \,. \tag{2}$$

Where **r** is subject to the inextensibility condition Eq. (1). Here, κ [ML³T⁻²] is the flexural rigidity of the polymer. κ is related to the persistence length ℓ_p [L] of the polymer as

$$\kappa = k_{\rm B} T \ell_{\rm p} \,. \tag{3}$$

The persistence length is the correlation lenth for tangent autocorrelation along the chain:

$$\langle \mathbf{t}(s) \cdot \mathbf{t}(s + \Delta s) \rangle = e^{-\Delta s/\ell_{\rm p}}.$$
 (4)

The WLC model describes all polymers that possess some degree of directional persistence. A sub-class of these are the so-called *semiflexible* polymers: chains whose contour length ℓ_c is of the order of their persistence length ℓ_p . By Eq. (4), $\ell_p \sim \ell_c$ implies that a semiflexible chain has a fair amount of orientational correlation, even between its two end points. In other words, the polymer is *weakly bent*, and mostly oriented along a single direction in space. We will choose our *z*-axis to lie along this direction, and will split up our tangent vector in a parallel, *z*-component $t_{\parallel} = t_z$ and a 2-dimensional perpendicular component $t_{\perp} = \{t_x, t_y\}$. Because t_z is much larger than both t_x and t_y , for a semiflexible chain we have

$$\mathbf{t}_{\parallel} \gg |\mathbf{t}_{\perp}| \,. \tag{5}$$

We may use this to expand the inextensibility condition Eq. (1) to second order in $|\mathbf{t}_{\perp}|$:

$$\mathbf{t}_{\parallel} = \left(1 - |\mathbf{t}_{\perp}|^2\right)^{1/2} \approx 1 - \frac{1}{2} |\mathbf{t}_{\perp}|^2.$$
 (6)

The semiflexible approximation consists of terminating this expansion at second order, which in effect—approximately implements the inextensibility condition Eq. (1) by eliminating one degree of freedom; the parallel component \mathbf{t}_{\parallel} . Using this expansion to eliminate, likewise, the parallel component \mathbf{r}_{\parallel} from the Hamiltonian, Eq. (2), and with the obvious notation $\mathbf{r}_{\perp} = {\mathbf{r}_x, \mathbf{r}_y}$, we find that the semiflexible energy functional \mathcal{H}_{SF0} [ML²T⁻²] may be written as

$$\mathcal{H}_{SF0} = \frac{\kappa}{2} \int_0^{\ell_c} \mathrm{d}s \, \left| \frac{\partial^2 \mathbf{r}_\perp}{\partial s^2} \right|^2 \,. \tag{7}$$

In the presence of an external force $\mathbf{f} [\mathsf{MLT}^{-2}]$, this energy is augmented with the addition of a term $\mathcal{H}_f = -\mathbf{f} \cdot (\mathbf{r}(\ell_c) - \mathbf{r}(0))$. Without loss of generality, we may choose \mathbf{f} to lie along the *z*-axis ($\mathbf{f} = f\hat{z}$), and pin one end of the polymer to the origin by setting $\mathbf{r}(0) = \vec{0}$. The potential energy term then reduces to $\mathcal{H}_f = -f\mathbf{r}_{\parallel}(\ell_c)$, with \mathbf{r}_{\parallel} the *z*-component of \mathbf{r} . Note, that $\mathbf{r}_{\parallel}(\ell_c)$ is the projected length of the polymer in the *z*-direction, which we may write as an integral over *s* and approximate to second perpendicular order as before in Eq. (6):

$$\mathbf{r}_{\parallel}(\ell_{\rm c}) = \int_{0}^{\ell_{\rm c}} \mathrm{d}s \left(\frac{\partial \mathbf{r}_{\parallel}}{\partial s}\right) \approx \ell_{\rm c} - \frac{1}{2} \int_{0}^{\ell_{\rm c}} \mathrm{d}s \left|\frac{\partial \mathbf{r}_{\perp}}{\partial s}\right|^{2} \,. \tag{8}$$

In this last equation, the second term represents the reduction of the projected length due to the transverse fluctuations. Only if these are zero everywhere is the projected length equal to the full contour length. Dropping a (constant) term $-f\ell_c$ from the total energy functional $\mathcal{H}_{SF} = \mathcal{H}_{SF0} + \mathcal{H}_f$, we arrive at the full Hamiltonian for a semiflexible polymer, subject to an external force f:

$$\mathcal{H}_{SF} = \frac{1}{2} \int_0^{\ell_c} \mathrm{d}s \left[\kappa \left| \frac{\partial^2 \mathbf{r}_\perp}{\partial s^2} \right|^2 + f \left| \frac{\partial \mathbf{r}_\perp}{\partial s} \right|^2 \right] \,. \tag{9}$$

For what follows, it will be instructive to gather the two independent components of \mathbf{r}_{\perp} into a single complex-valued function r(s,t) as

$$\mathbf{r} \equiv \mathbf{r}_x + i\mathbf{r}_y \,. \tag{10}$$

In terms of this complex variable,

$$\mathcal{H}_{SF} = \frac{1}{2} \int_0^{\ell_c} \mathrm{d}s \left[\kappa \left| \frac{\partial^2 \mathbf{r}}{\partial s^2} \right|^2 + f \left| \frac{\partial \mathbf{r}}{\partial s} \right|^2 \right] \,, \tag{11}$$

where norms are now to be interpreted as $|\sigma|^2 = \sigma \sigma^*$.

3 Nonlinear Force-Extension of a Single Semiflexible Chain

As a first step towards computing the mechanical response of networks, we derive the forceextension characteristics of a single chain. That is, we will be computing the expectation value of the projected length, $\ell = \langle \mathbf{r}_{\parallel}(\ell_c) \rangle$ [L] as a function of the applied force f. First, we Fourier transform Eq. (11) (for Fourier transform conventions, see Appendix A). For didactic purposes, let us spell out the transformation of one of the terms in the Hamiltonian.

$$\begin{split} \int_{0}^{\ell_{c}} \mathrm{d}s \left| \frac{\partial^{2} \mathbf{r}}{\partial s^{2}} \right|^{2} &= \int_{0}^{\ell_{c}} \mathrm{d}s \left(\frac{1}{\ell_{c}} \sum_{q} (-q^{2}) \, \mathbf{r}_{q} \, e^{iqs} \right) \left(\frac{1}{\ell_{c}} \sum_{q'} (-q'^{2}) \, \mathbf{r}_{q'}^{\star} \, e^{-iq's} \right) \,, \\ &= \frac{1}{\ell_{c}^{2}} \sum_{qq'} \int_{0}^{\ell_{c}} \mathrm{d}s \, (q^{2}q'^{2}) \, e^{i(q-q')s} \, \mathbf{r}_{q} \mathbf{r}_{q'}^{\star} \,, \\ &= \frac{1}{\ell_{c}} \sum_{qq'} (q^{2}q'^{2}) \delta_{qq'} \mathbf{r}_{q} \mathbf{r}_{q'}^{\star} \,, \\ &= \frac{1}{\ell_{c}} \sum_{q} q^{4} \, |\mathbf{r}_{q}|^{2} \,, \end{split}$$
(12)

where we have used the resolution of the Kronecker δ , Eq. (88) in Appendix A. Note the dimensions of $q [L^{-1}]$ and $r_q [L^2]$. Transforming the other term in completely analogous fashion, we arrive at the following expression for the Hamiltonian

$$\mathcal{H}_{SF} = \frac{1}{2\ell_c} \sum_q \left(\kappa q^4 + f q^2 \right) |\mathbf{r}_q|^2 \,. \tag{13}$$

This energy functional is a sum of independent, quadratic terms and thus we may apply the equipartition theorem to it: the expectation value for each of its terms is $\frac{1}{2}k_{\rm B}T$ [ML²T⁻²], per degree of freedom. Taking into account the fact that r_q is a complex variable and thus represents *two* degrees of freedom, we obtain

$$\frac{1}{2\ell_{\rm c}} \left(\kappa q^4 + fq^2\right) \left\langle |\mathbf{r}_q|^2 \right\rangle = k_{\rm B}T \,, \tag{14}$$

or, equivalently, that the power spectrum is given by

$$\langle |\mathbf{r}_q|^2 \rangle = \frac{2k_{\rm B}T\ell_{\rm c}}{\kappa q^4 + fq^2} \,. \tag{15}$$

This, coincidentally, provides one way to measure the persistence length (or κ) using the fluctuation spectrum: FT'ing the transverse fluctuations (at zero force) and fitting these to a q^4 power law. We will use the expression to compute the expectation value of the end-to-end length which, as we had derived earlier in Eq. (8), may be written in terms of the transverse displacement field as

$$\mathbf{r}_{\parallel}(\ell_{\rm c}) \approx \ell_{\rm c} - \frac{1}{2} \int_0^{\ell_{\rm c}} {\rm d}s \left| \frac{\partial \mathbf{r}_{\perp}}{\partial s} \right|^2 = \ell_{\rm c} - \frac{1}{2\ell_{\rm c}} \sum_q q^2 |\mathbf{r}_q|^2 \,. \tag{16}$$

From this, we can compute the equilibrium expectation value of the remaining 'slack' $\Delta(f)$ [L] in the polymer (the length still stored in transverse fluctuations, *i.e.* the contour length minus the projected length) by the following expression, where we write $\ell = \langle \mathbf{r}_{\parallel}(\ell_c) \rangle$,

$$\ell_{\rm c} - \ell(f) = \frac{1}{2\ell_{\rm c}} \sum_{q} q^2 \langle |\mathbf{r}_q|^2 \rangle,$$

$$\equiv \langle \Delta(f) \rangle.$$
(17)

Substituting the equipartition value for $\langle |\mathbf{r}_q|^2 \rangle$ yields

$$\ell_{\rm c} - \ell(f) = \frac{1}{2\ell_{\rm c}} \sum_{q} q^2 \langle |\mathbf{r}_q|^2 \rangle ,$$

= $k_{\rm B} T \sum_{q} \frac{1}{\kappa q^2 + f} \qquad \left(q = \frac{n\pi}{\ell_{\rm c}}; n = 1, 2, \cdots \right) .$ (18)

This sum can be computed analytically, but first we inspect the high-force regime. Passing to the continuum limit (see Appendix A) yields

$$\ell_{\rm c} - \ell(f) \approx \frac{k_{\rm B} T \ell_{\rm c}}{\pi} \int_0^\infty dq \, \left(\frac{1}{\kappa q^2 + f}\right) = \frac{\ell_{\rm c}}{2} \left(\frac{k_{\rm B} T}{f \ell_{\rm p}}\right)^{1/2} \,. \tag{19}$$

That is, the force *diverges*, as full extension (zero slack) is approached, with a characteristic power of -2 in the slack:

$$f \approx \left(\frac{k_{\rm B}T\ell_{\rm c}^2}{4\ell_{\rm p}}\right) \left(\ell_{\rm c} - \ell(f)\right)^{-2}.$$
(20)

Returning to the full Eq. (18), it is instructive to first introduce the dimensionless force φ

$$\varphi \equiv \frac{f\ell_{\rm c}^2}{\kappa}\,,\tag{21}$$

in terms of which the expectation value of the slack as a function of force relation reads

$$\langle \Delta(f) \rangle = \ell_{\rm c} - \ell(f) = \left(\frac{\ell_{\rm c}^2}{2\ell_{\rm p}}\right) \frac{1}{\varphi} \left[\sqrt{\varphi} \coth\sqrt{\varphi} - 1\right].$$
⁽²²⁾

This expression contains the full, nonlinear force-extension of the semiflexible chain, but not in its most obvious form. The slack is not equal to the extension, of course—a more natural extension variable to consider is the extension away from the equilibrium length that is prompted by the force f. To compute this, we first note that taking the limit $\varphi \to 0$ gives the equilibrium slack:

$$\langle \Delta(0) \rangle = \ell_{\rm c} - \ell(f=0) = \frac{\ell_{\rm c}^2}{6\ell_{\rm p}}.$$
 (23)

This is the amount of contour length that a semiflexible polymer, in equilibrium, typically stores in transverse fluctuations. As a result, the equilibrium length ℓ_0 is the full contour length minus the equilibrium slack

$$\ell_0 = \ell_c \left(1 - \frac{\ell_c}{6\ell_p} \right) \,. \tag{24}$$

We now have all we need to define the actual extension away from equilibrium. Introducing the nondimensionalized extension $\delta \ell$ as the difference between the total length at force f minus the equilibrium length ℓ_0 , normalized by ℓ_c^2/ℓ_p ;

$$\tilde{\delta}\ell = \left(\frac{\ell_{\rm c}^2}{\ell_{\rm p}}\right)^{-1} \left(\ell(f) - \left[1 - \frac{\ell_{\rm c}^2}{6\ell_{\rm p}}\right]\right),\tag{25}$$

allows us to express the full, nonlinear force-extension of a semiflexible chain as

$$\tilde{\delta}\ell(\varphi) = \frac{1}{6\varphi} \left(\varphi - 3\sqrt{\varphi} \coth\sqrt{\varphi} + 3\right).$$
(26)

This expression highlights the universality of the force-extension response; upon proper scaling with the relevant length- and energy scales, a single master curve describes all semiflexible polymers. At small forces, we find by Taylor expansion that

$$\tilde{\delta}\ell(\varphi) \approx \frac{\varphi}{90}$$
. (27)

Transforming back to dimensional variables, this means that at small forces the semiflexible chain (with contour length ℓ_c and persistence length ℓ_c) behaves like a linear (Hookean) spring, with a spring constant $k_{\rm sp}$ [MT⁻²] and a rest length ℓ_0 [L] that may be computed from the chain's characteristic length scales and the temperature as

$$\mathcal{E}_{\rm lin} = \frac{1}{2} k_{\rm sp} \left(\ell - \ell_0\right)^2 , \quad \left(\text{with } k_{\rm sp} = 90 k_{\rm B} T \left(\frac{\ell_{\rm p}}{\ell_{\rm c}^2}\right)^2, \, \ell_0 = \ell_{\rm c} \left(1 - \frac{\ell_{\rm c}}{6\ell_{\rm p}}\right) \right). \tag{28}$$

For higher forces, terms of higher order in the extension $\delta \ell = \ell - \ell_0$ come into play (straighforwardly computed by higher order Taylor expansion of Eq. (26)), until the energy diverges in the limit $\ell \rightarrow \ell_c$. It is, however, still possible to define a nonlinear energy function which, upon derivation with respect to $\delta \ell$, gives the force by integrating (the dimensionalized version of) the inverse of Eq. (26)

$$\mathcal{E}(\delta\ell) \equiv \int_0^{\delta\ell} \mathrm{d}(\delta\ell') f(\delta\ell') \,. \tag{29}$$

The upshot of all of this is, that a semiflexible polymer behaves as a nonlinear spring, whose elastic response may be systematically and analytically computed over the entire range of allowed extensions.

4 Static Response of an Affinely Deforming Semiflexible Network

As we focus now on *networks* of such nonlinear polymers, we are going to be dealing with more complicated deformations than simple longitudinal extensions. It will be useful to first consider the geometry of stresses and strains in three dimensions.

A general deformation R maps points \mathbf{x} [L] in a reference space onto images \mathbf{x}' [L] in a target space as

$$\mathbf{x} \mapsto \mathbf{x}'(\mathbf{x}) = \mathbf{R}(\mathbf{x}) \,. \tag{30}$$

There is no deformation when $\mathbf{R}(\mathbf{x}) = \mathbf{x}$. To distinguish non-trivial deformations, we split of the *displacement vector* \mathbf{u} [L]:

$$\mathbf{R}(\mathbf{x}) = \mathbf{x} + \mathbf{u}(\mathbf{x}) \,. \tag{31}$$

Assuming that distortions vary slowly in space, we may linearize this around the origin O to find

$$R_i(\mathbf{x}) \approx R_i(\mathbf{O}) + \left(\frac{\partial R_i(\mathbf{x})}{\partial x_j}\right) x_j.$$
 (32)

In what follows, we will subtract any uniform translations (*i.e.*, $\mathbf{R}(\mathbf{O}) = \mathbf{O}$), and summation over repeated indices will be implied. The first derivatives appearing in this equation define the deformation tensor $\Lambda(\mathbf{x})$ [dimensionless]

$$\Lambda_{ij}(\mathbf{x}) = \frac{\partial R_i(\mathbf{x})}{\partial x_j} = \delta_{ij} + \frac{\partial u_i}{\partial x_j} \equiv \delta_{ij} + \eta_{ij}.$$
(33)

The tensor η [dimensionless] is called the displacement gradient tensor. A transformation is called *affine* when the deformation tensor is constant throughout the the entire body or volume that is deformed. In other words, the deformation is affine if and only if $\Lambda \neq \Lambda(\mathbf{x})$. If this is the case, the constant deformation tensor Λ effects a true linear mapping

$$\mathbf{x} \mapsto \mathbf{x}' = \mathbf{\Lambda} \cdot \mathbf{x} \,. \tag{34}$$

While affinely deforming systems are rarely, if ever, encountered in disordered networks in *crystals*, affinity is the norm: the local symmetries, in fact, dictate that it should apply. The variable that is energy-conjugate to the displacement gradient tensor (and therefore to Λ , too) is the first Piola-Kirchhoff stress tensor σ^I . The elastic energy \mathcal{F}_{el} [ML²T⁻²] required to deform an infinitesimal reference volume Ω of the system, which may be computed as the integral of the elastic energy density w [ML⁻¹T⁻²],

$$\mathcal{F}_{\rm el} = \int_{\Omega} w(\Lambda) \, \mathrm{d}^3 \mathbf{x} \,, \tag{35}$$

may be used to compute the first Piola-Kirchhoff stress tensor σ^{I} [ML⁻¹T⁻²] by simple derivation of the energy density

$$\sigma_{ij}^{I} = \frac{\partial w(\Lambda)}{\partial \Lambda_{ij}} \,. \tag{36}$$

Because Λ is defined by derivatives of \mathbf{R} with respect to positions in the reference volume, σ^I is a so-called mixed stress tensor - it measures the force in target space per unit area in reference space. Likewise, $w(\Lambda)$ is the elastic energy density defined relative to the reference volume Ω [L⁻³]. Experiments, however, typically record the force in deformed space per unit area in the deformed space. The corresponding stress tensor is the true or Cauchy stress tensor σ^C . By using the facts that upon transforming to the reference space \mathbf{x}' the total elastic energy \mathcal{F}_{el} should not change, that $d^3\mathbf{x} = (\det \Lambda)^{-1}d^3\mathbf{x}'$ and that $\partial_{\mathbf{x}'} = \Lambda^T \partial_{\mathbf{x}}$, we can compute the Cauchy stress σ^C [ML⁻¹T⁻²] from the first Piola-Kirchhoff stress directly as

$$\sigma^C = \frac{1}{\det \Lambda} \sigma^I \cdot \Lambda^T \,. \tag{37}$$

The strain measure conjugate to the Cauchy stress tensor is $\frac{\partial u_i}{\partial x_{\alpha}'}$, which is generally not simple to compute (it leads to the so-called Almansi strain tensor), and computing the Cauchy stress in general settings is typically most straightforwardly achieved by first computing $w(\Lambda)$, taking derivatives w.r.t. the components of Λ , and transforming the resultant σ^I to σ^C using Eq. (37). This is the procedure we will adopt in the following.

Note, that this formalism applies equally to linearized elasticity theories and finite-strain systems. The constitutive behavior underlying the relation between stress and strain need not be linear (Hookean)—in fact Eq. (36) derives from general thermodynamic considerations. The assumption is made that the strain varies slowly over Ω but neither it, nor the stresses, need be small.

Consider now a network of polymers subjected to a deformation Λ . The elastic energy density $w(\Lambda)$ for a volume Ω is then simply given by the sum of all contributions from individual springs:

$$w(\Lambda) = \frac{1}{\Omega} \sum_{\alpha} \mathcal{E}^{\alpha}(\Lambda) \,. \tag{38}$$

The index α labels the springs, and $\mathcal{E}^{\alpha}(\Lambda)$ is the elastic energy of each single spring subject to Λ —the quantity defined in Eq. (29) with Λ generating the single chain extensions $\delta \ell^{\alpha}$ [L]. For central force networks, this energy is a function of the *length* of the deformed chain only: denoting by \mathbf{L}_{α} [L] the undeformed end-to-end vector of polymer α we therefore have

$$\mathcal{E}^{\alpha}(\Lambda) = \mathcal{E}^{\alpha}(|\Lambda \cdot \mathbf{L}^{\alpha}|).$$
(39)

The first Piola-Kirchhoff stress tensor for such a network can be computed using Eq. (36):

$$\sigma_{ij}^{I} = \frac{\partial w(\Lambda)}{\partial \Lambda_{ij}} = \frac{1}{\Omega} \sum_{\alpha} \frac{\partial}{\partial \Lambda_{ij}} \mathcal{E}^{\alpha}(|\Lambda \cdot \mathbf{L}^{\alpha}|) \,. \tag{40}$$

By the chain rule, the differential appearing in the summation may be rewritten

$$\frac{\partial}{\partial \Lambda_{ij}} \mathcal{E}^{\alpha}(|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|) = \left(\frac{\partial \mathcal{E}^{\alpha}(|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|)}{\partial |\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|}\right) \left(\frac{\partial |\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|}{\partial \Lambda_{ij}}\right) \,.$$

The first term between brackets is simple the derivative of the spring energy w.r.t. the spring length—*i.e.*, the *force* as a function of the length $f^{\alpha}(|\Lambda \cdot \mathbf{L}^{\alpha}|)$. The second term may be rewritten as

$$\frac{\partial |\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|}{\partial \Lambda_{ij}} = \frac{\partial}{\partial \Lambda_{ij}} \left(\Lambda_{kl} \mathbf{L}_{l}^{\alpha} \Lambda_{km} \mathbf{L}_{m}^{\alpha} \right)^{1/2},$$

$$= \frac{1}{2|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|} \left(\delta_{ki} \delta_{lj} \mathbf{L}_{l}^{\alpha} \Lambda_{km} \mathbf{L}_{m}^{\alpha} + \Lambda_{kl} \mathbf{L}_{l}^{\alpha} \delta_{ki} \delta_{mj} \mathbf{L}_{m}^{\alpha} \right),$$

$$= \frac{(\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha})_{i} (\mathbf{L}^{\alpha})_{j}}{|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|},$$
(41)

so that the components of the first Piola-Kirchhoff stress tensor for the network are given by

$$\sigma_{ij}^{I} = \frac{1}{\Omega} \sum_{\alpha} f^{\alpha}(|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|) \left(\frac{(\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha})_{i}(\mathbf{L}^{\alpha})_{j}}{|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|} \right) \,. \tag{42}$$

Recalling that we may transform σ^{I} to the Cauchy stress σ^{C} using Eq. (37), we find

$$\sigma_{ij}^{C} = \frac{1}{\Omega \det \Lambda} \sum_{\alpha} f^{\alpha}(|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|) \left(\frac{(\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha})_{i}(\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha})_{j}}{|\mathbf{\Lambda} \cdot \mathbf{L}^{\alpha}|}\right) \,. \tag{43}$$

We now pass to a continuum limit, replacing the sum by an integral over the distribution $\mathcal{P}(\mathbf{L})$ [L⁻³] of undeformed end-to-end lengths, which for an arbitrary function g happens as:

$$\sum_{\alpha=1}^{N} g(\mathbf{L}^{\alpha}) \longrightarrow N\langle g(\mathbf{L}) \rangle = N \int \mathcal{P}(\mathbf{L}) g(\mathbf{L}) \,\mathrm{d}^{3}\mathbf{L} \,. \tag{44}$$

Letting ρ [L⁻³] denote the number density of polymers in the reference configuration, *i.e.* $\rho = N/\Omega$ this yields the expression for the Cauchy stress appropriate for large polymer networks:

$$\sigma_{ij}^{C} = \frac{\rho}{\det \Lambda} \int \mathcal{P}(\mathbf{L}) f(|\mathbf{\Lambda} \cdot \mathbf{L}|) \left(\frac{(\mathbf{\Lambda} \cdot \mathbf{L})_{i}(\mathbf{\Lambda} \cdot \mathbf{L})_{j}}{|\mathbf{\Lambda} \cdot \mathbf{L}|} \right) d^{3}\mathbf{L}$$
$$= \rho_{\mathbf{\Lambda}} \left\langle f(|\mathbf{\Lambda} \cdot \mathbf{L}|) \left(\frac{(\mathbf{\Lambda} \cdot \mathbf{L})_{i}(\mathbf{\Lambda} \cdot \mathbf{L})_{j}}{|\mathbf{\Lambda} \cdot \mathbf{L}|} \right) \right\rangle_{\mathcal{P}(\mathbf{L})}.$$
(45)

The integral is over all of 3-space. A nonlinear dependence of the stress tensor component σ_{ij}^C on the strain Λ may thus originate from two sources: a nonlinear force-extension curve *or* a geometrically incurred nonlinearity arising from the second, bracheted term in the average—this term does not depend on the force-extension behavior. The prefactor ρ_{Λ} [L⁻³] is the number density in the deformed configuration:

$$\rho_{\Lambda} = \left(\frac{N}{\Omega}\right) \left(\frac{\Omega}{\Omega'}\right) = \frac{\rho}{\det\Lambda}.$$
(46)

For incompressible systems, det $\Lambda = 1$ and the distinction between reference and deformed densities is irrelevant. Eq. (45) does not assume isotropy, in fact $\mathcal{P}(\mathbf{L})$ may take on any form including the collection of δ -peaks appropriate for crystalline structures. Provided one knows the force-extension curve, the initial radial distribution of spring end-to-end vectors, and, obviously, the strain Λ this expression allows one to directly compute, to all desired nonlinear orders, the affine Cauchy stress.

This equation is the basis for the analysis in [3], which amounts to substituting the nonlinear, semiflexible force-extension Eq. (26) into Eq. (45), assuming an isotropic distribution of filaments which all have the same rest length ℓ_0 ;

$$\mathcal{P}(\mathbf{L}) = \frac{1}{4\pi\ell_0^2} \delta(|\mathbf{L}| - \ell_0).$$
(47)

Extracting all components of the stress tensor is, in principle, possible but cumbersome due to the required inversion of the extension-force relation. As a simple demonstration, however, we can compute the linear shear stress of a network of semiflexible polymers. For the case of simple shear in the xz-direction, the deformation tensor is given by

$$\Lambda(\gamma) = \begin{pmatrix} 1 & 0 & \gamma \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} ,$$
 (48)

with γ [dimensionless] the shear strain. Expanding Eq. (45) to lowest order in γ , and using the linear single-chain energy \mathcal{E}_{lin} given by Eq. (28), we find

$$\sigma_{xz}^{C} \approx \left(\frac{\rho_{\Lambda}k_{\rm sp}\ell_{0}^{2}}{4\pi}\right)\gamma \int_{0}^{\pi} \mathrm{d}\theta \int_{0}^{2\pi} \mathrm{d}\varphi \left[\cos^{2}\theta \sin^{3}\theta \cos^{2}\varphi\right]$$
$$= \left(\frac{1}{15}\rho_{\Lambda}k_{\rm sp}\ell_{0}^{2}\right)\gamma. \tag{49}$$

By definition, the linear shear modulus G_0 [ML⁻¹T⁻²] is the coefficient of linear response for this loading protocol:

$$\sigma_{xz}^C \equiv G_0 \,\gamma \,, \tag{50}$$

from which we infer that the linear shear modulus of an isotropic network of semiflexible chains, deforming affinely, is

$$G_0 = 6\rho_{\Lambda}k_{\rm B}T\left(\frac{\ell_{\rm p}^2\ell_0^2}{\ell_{\rm c}^4}\right) \approx 6\rho_{\Lambda}k_{\rm B}T\left(\frac{\ell_{\rm p}}{\ell_{\rm c}}\right)^2.$$
(51)

Those familiar with the classical theory of rubber elasticity will notice the similar expression in this case ($G_{\text{class}} = \rho_{\Lambda} k_{\text{B}} T$). The approach we outline here is, however, much more general and

may be taken to arbitrary nonlinear order, even for anisotropically ordered systems. For those, however, and more general for sparser systems, the assumption of non-affinity may break down. This is beyond the scope of the current lecture, but has been studied in considerable detail for both flexible rubbers and biopolymer gels.

5 Dynamic Response of Semiflexible Chains and Networks

To assess the dynamical response of the network, we first turn to the single filament again. We will look at dynamical response by analyzing the equilibrium fluctuations, and so will set the external force f to zero for now. Our point of departure is the semiflexible Hamiltonian Eq. (11), which we write as the *s*-integral of a Hamiltonian density functional h(r) [MLT⁻²]

$$\mathcal{H}_{SF} = \frac{\kappa}{2} \int_0^{\ell_c} \mathrm{d}s \left| \frac{\partial^2 \mathbf{r}}{\partial s^2} \right|^2 \equiv \int_0^{\ell_c} \mathrm{d}s \, \mathbf{h}(\mathbf{r}) \,. \tag{52}$$

Letting ζ [ML⁻¹T⁻¹] denote the transverse drag coefficient (per unit length) on a cylinder of radius *d* [L] and with length λ [L] moving through a medium with viscosity η [ML⁻¹T⁻¹]

$$\zeta = \frac{4\pi\eta}{\ln(\lambda/d)},\tag{53}$$

we find the (overdamped) Langevin equation that governs the dynamics of transverse fluctuations:

$$-\zeta \dot{\mathsf{r}}(s,t) - \frac{\delta \mathcal{H}_{SF}}{\delta \mathsf{r}} + \xi_{\mathsf{r}}(s,t) = 0, \qquad (54)$$

with the dot denoting derivation w.r.t. time, and $\xi_r(s,t)$ [MT⁻²] the fluctuating force per unit length. If we further let primes denote derivation w.r.t. *s*, the variational derivative is given by

$$\frac{\delta \mathcal{H}_{SF}}{\delta \mathbf{r}} = \frac{\mathrm{d}^2}{\mathrm{d}s^2} \left(\frac{\partial \mathbf{h}}{\partial \mathbf{r}''} \right) = \kappa \mathbf{r}''''(s, t) \,. \tag{55}$$

Taken together, the dynamic equation for r is thus

$$\dot{\mathsf{r}}(s,t) = -\left(\frac{\kappa}{\zeta}\right)\mathsf{r}'''(s,t) + \frac{\xi_{\mathsf{r}}(s,t)}{\zeta}\,. \tag{56}$$

This last equation is Fourier transformed once (see Appendix A), to go from real to q-space but retaining the time dependence (units: $r_q [L^2]$, $\xi_q [MLT^{-2}]$), to produce

$$\dot{\mathbf{r}}_q(t) = -\left(\frac{\kappa}{\zeta}\right)q^4\,\mathbf{r}_q(t) + \frac{\xi_q(t)}{\zeta}\,.\tag{57}$$

The quantity multiplying $r_q(t)$ on the right has dimensions $[T^{-1}]$, and defines a frequency ω_q $[T^{-1}]$:

$$\omega_q = \left(\frac{\kappa}{\zeta}\right) q^4 \,. \tag{58}$$

We now spell out in some detail how to solve the final Langevin equation, which reads (dropping the explicit t dependencies)

$$\dot{\mathbf{r}}_q = -\omega_q \mathbf{r}_q + \frac{\xi_q}{\zeta} \,. \tag{59}$$

A straightforward way to solve Eq. (59) is to use an integrating factor, multiplying the entire equation by $e^{\omega_q t}$. Doing so, and rearranging, gives

$$e^{\omega_q t} \left(\dot{\mathsf{r}}_q + \omega_q \mathsf{r}_q \right) = \frac{1}{\zeta} e^{\omega_q t} \xi_q \,. \tag{60}$$

If we define $\Omega_q \equiv e^{\omega_q t} \mathbf{r}_q \ [\mathsf{L}^2]$, we immediately recognize that the LHS of the previous equation is the time derivative $\dot{\Omega}_q$, which allows us to formally solve it by integrating

$$\dot{\Omega}_q = \frac{1}{\zeta} e^{\omega_q t} \xi_q \longrightarrow \Omega_q = \frac{1}{\zeta} \int_{-\infty}^t \mathrm{d}t_1 \, e^{\omega_q t_1} \, \xi_q(t_1) \,. \tag{61}$$

Thus, the value of r_q at some time t may be obtained by integrating the (weighted) noise as

$$\mathbf{r}_{q}(t) = \frac{1}{\zeta} \int_{-\infty}^{t} \mathrm{d}t_{1} \, e^{-\omega_{q}(t-t_{1})} \, \xi_{q}(t_{1}) \,. \tag{62}$$

Using this expression allows us to compute *any* correlation we are interested in, provided we know the correlative properties of the noise term. For now, let us assume this to be a white noise term with mean zero ($\langle \xi_q(t) \rangle = 0$), and amplitude A [M²L²T⁻³]:

$$\langle \xi_q(t_1)\xi_q^*(t_2)\rangle = \mathsf{A}\delta(t_1 - t_2)\,. \tag{63}$$

In a moment we will compute the value of A, but first let us compute the following correlation explicitly

$$\langle \mathbf{r}_{q}(t)\mathbf{r}_{q}^{\star}(0)\rangle = \frac{1}{\zeta^{2}} \int_{-\infty}^{t} dt_{1} \int_{-\infty}^{0} dt_{2} e^{-\omega_{q}(t-t_{1})} e^{\omega_{q}t_{2}} \langle \xi_{q}(t_{1})\xi_{q}^{\star}(t_{2})\rangle ,$$

$$= \frac{A}{\zeta^{2}} \int_{-\infty}^{0} dt_{2} e^{-\omega_{q}(t-t_{2})} e^{\omega_{q}t_{2}} ,$$

$$= \frac{A}{2\omega_{q}\zeta^{2}} e^{-\omega_{q}t} \quad (t > 0).$$
(64)

In going from the first to the second line, we have inserted Eq. (63) and have had to assume that t > 0. For t > 0, the t_1 -integral needs to be executed first to ensure the integration interval contains the peak of the δ -function. Of course, we could also assume that t < 0 in which case the t_2 integral has to be performed first. The result is similar, except for the fact that $t \rightarrow -t$ in the exponent. Taken together, and valid for all times, the result is

$$\langle \mathbf{r}_q(t) \mathbf{r}_q^{\star}(0) \rangle = \frac{\mathsf{A}}{2\omega_q \zeta^2} e^{-\omega_q |t|} \,. \tag{65}$$

To fix the value of A, we may compute this same correlator in a slightly different fashion, too. First, we Fourier transform Eq. (52) to obtain

$$\mathcal{H}_{SF} = \frac{\kappa}{2\ell_{\rm c}} \sum_{q} q^4 |\mathbf{r}_q|^2 \,. \tag{66}$$

By the equipartition theorem, each degree of freedom gets $\frac{1}{2}k_{\rm B}T$ in energy; the complex quantity r_q represents two degrees of freedom (one for the real part, one for the imaginary part), and therefore

$$\langle |\mathbf{r}_q|^2 \rangle = \langle \mathbf{r}_q(0) \mathbf{r}_q^*(0) \rangle = \frac{2k_{\rm B}T\ell_{\rm c}}{\kappa q^4} \,. \tag{67}$$

Comparing this result with Eq. (65), evaluated at t = 0, we find that the noise amplitude that yields the correct equilibrium correlations is

$$\mathsf{A} = 4k_{\mathrm{B}}T\ell_{\mathrm{c}}\zeta\,,\tag{68}$$

and thus that the full two-time correlation function is

$$\langle \mathsf{r}_q(t)\mathsf{r}_q^*(0)\rangle = \left(\frac{2\ell_{\rm c}}{q^4\ell_{\rm p}}\right) e^{-\omega_q|t|} \,. \tag{69}$$

For the dynamic response function, we shall be interested in fluctuations of the end-to-end length of the polymer. Recalling that we work at f = 0 (linear regime), the slack Δ [L] in a given conformation is given by

$$\Delta(t) = \frac{1}{2\ell_{\rm c}} \sum_{q} q^2 |\mathbf{r}_q(t)|^2 \,, \tag{70}$$

with equilibrium expectation value (at zero force)

$$\langle \Delta \rangle = \frac{\ell_{\rm c}}{6\ell_{\rm p}^2} \,. \tag{71}$$

In terms of the slack, the extension $\delta \ell(t)$ away from the equilibrium length $\ell_0 = \ell_c - \langle \Delta \rangle$ is given by

$$\delta \ell(t) = \langle \Delta \rangle - \Delta(t) \,. \tag{72}$$

We may now compute the zero-force autocorrelation $\phi(t)$ [L²] of the extension away from the equilibrium length:

$$\begin{aligned}
\phi(t) &\equiv \langle \delta\ell(t)\delta\ell(0) \rangle, \\
&= \langle \Delta(t)\Delta(0) \rangle - \langle \Delta \rangle^2,
\end{aligned}$$
(73)

where we use that $\langle \Delta(t) \rangle = \langle \Delta(0) \rangle = \langle \Delta \rangle$, the average being independent of time. By means of Eq. (70), we can compute the first term

$$\langle \Delta(t)\Delta(0)\rangle = \frac{1}{4\ell_{\rm c}^2} \sum_q q^4 \left\langle |\mathbf{r}_q(t)|^2 |\mathbf{r}_q(0)|^2 \right\rangle,$$

$$= \frac{1}{4\ell_{\rm c}^2} \sum_q q^4 \left\langle \mathbf{r}_q(t)\mathbf{r}_q^{\star}(t)\mathbf{r}_q(0)\mathbf{r}_q^{\star}(0) \right\rangle.$$

$$(74)$$

The four-point correlator appearing here may be simplified using Wick's theorem, which states that for variables x_i drawn from Gaussian distributions the following identity holds:

$$\langle x_1 x_2 x_3 x_4 \rangle = \langle x_1 x_2 \rangle \langle x_3 x_4 \rangle + \langle x_1 x_3 \rangle \langle x_2 x_4 \rangle + \langle x_1 x_4 \rangle \langle x_2 x_3 \rangle .$$
(75)

This reduces the expression for $\phi(t)$ to

$$\phi(t) = \frac{1}{4\ell_{\rm c}^2} \sum_q q^4 \langle \mathbf{r}_q(t) \mathbf{r}_q^{\star}(0) \rangle^2,$$

$$= \frac{1}{\ell_{\rm p}^2} \sum_q \left(\frac{1}{q^4}\right) e^{-2\omega_q |t|},$$
(76)

where in the last line we have used Eq. (69). Again, $\phi(t)$ may be determined directly by recording the fluctuations of the end-to-end length around its mean. The temporal Fourier transform of $\phi(t)$ is called the Power Spectral Density (PSD) $\langle |\delta \ell_{\omega}|^2 \rangle$ [L²T], and is given by

$$\langle |\delta\ell_{\omega}|^{2} \rangle = \int_{-\infty}^{\infty} dt \, \phi(t) \, e^{i\omega t} \,,$$

$$= \frac{1}{\ell_{\rm p}^{2}} \sum_{q} \left(\frac{1}{q^{4}}\right) \int_{-\infty}^{\infty} dt \, e^{-2\omega_{q}|t|} \, e^{i\omega t} \,,$$

$$= \frac{1}{\ell_{\rm p}^{2}} \sum_{q} \left(\frac{1}{q^{4}}\right) \left(\frac{4\omega_{q}}{\omega^{2} + 4\omega_{q}^{2}}\right) \quad (q = \frac{n\pi}{\ell_{\rm c}}, n = 1, 2, \dots).$$

$$(77)$$

Recasting this last expression slightly in terms of the fundamental frequency $\omega_1 = \omega_q n^{-4}$ and the fundamental wavenumber $q_1 = q n^{-1}$, we finally find

$$\langle |\delta \ell_{\omega}|^2 \rangle = \frac{1}{\omega_1 q_1^4 \ell_p^2} \sum_{n=1}^{\infty} \frac{1}{n^8 + (\omega/2\omega_1)^2} \,. \tag{78}$$

To connect this to mechanical response, we require the Fluctuation Dissipation Theorem. It allows to compute, from the PSD, the imaginary part of the complex response function α_{ω} [M⁻¹T²]—the function that measures the extension in response to an ω -periodic force

$$\delta\ell(\omega, t) = \alpha_{\omega} f(\omega, t) \,. \tag{79}$$

Note, that from this (as from the dimensions) it is clear that α_{ω} is the dynamic counterpart of the inverse spring constant. The FDT (without further proof here) states, that

$$\alpha_{\omega}^{\prime\prime} = \operatorname{Im}(\alpha_{\omega}) = \frac{\omega}{2k_{\rm B}T} \langle |\delta\ell_{\omega}|^2 \rangle .$$
(80)

Therefore, explicitly,

$$\alpha''_{\omega} = \frac{1}{k_{\rm B}Tq_1^4\ell_{\rm p}^2} \sum_{n=1}^{\infty} \frac{(\omega/2\omega_1)}{n^8 + (\omega/2\omega_1)^2} \,. \tag{81}$$

The real part $\alpha'_{\omega} = \text{Re}(\alpha_{\omega})$ may be found either by the Kramers-Kronig relation, or—if you've seen these before—by inspection. Either approach gives for the full complex quantity α_{ω} that

$$\alpha_{\omega} = \frac{1}{k_{\rm B} T q_1^4 \ell_{\rm p}^2} \sum_{n=1}^{\infty} \frac{1}{n^4 - i(\omega/2\omega_1)^2} \,. \tag{82}$$

This completes our calculation of the linear, dynamic response. At low frequencies, all dissipative processes will have relaxed and only the elastic response remains. In the limit $\omega \rightarrow 0$, the imaginary part in the denominator may be set to zero and the sum executed analytically, to yield

$$\alpha_0 = \frac{1}{k_{\rm B}Tq_1^4\ell_{\rm p}^2} \sum_{n=1}^{\infty} \frac{1}{n^4} = \frac{1}{90k_{\rm B}T} \left(\frac{\ell_{\rm c}^2}{\ell_{\rm p}}\right)^2 \,. \tag{83}$$

Indeed, we recognize the inverse of the linear spring constant k_{sp} derived in Eq. (28). For the high frequencies, conversely, we may pass to a continuum formulation and approximate α_{ω} as

an integral:

$$\alpha_{\omega} \approx \frac{1}{k_{\rm B}Tq_1^4 \ell_{\rm p}^2} \int_0^\infty dn \left(\frac{1}{n^4 - i(\omega/2\omega_1)^2}\right),$$

$$= \frac{1}{2\sqrt{2}k_{\rm B}T} \left(\frac{\ell_{\rm c}}{\ell_{\rm p}^2}\right) \left(\frac{2i\kappa}{\zeta}\right)^{3/4} \omega^{-3/4}.$$
(84)

Finally, the dynamic shear modulus $G(\omega)$ of a network may be obtained directly from α_{ω} via Eq. (45), by substituting the dynamic force-extension relation Eq. (79), $(f(\omega, t) = \alpha_{\omega}^{-1} \delta \ell(\omega, t))$ for the force f. The result, derived in completely analogous fashion to Eq. (49) for an isotropic network, deforming affinely, is that

$$G(\omega) = \frac{1}{15} \rho_{\Lambda} \alpha_{\omega}^{-1} \ell_{\rm c}^2 \,. \tag{85}$$

The corresponding high-frequency behavior $G(\omega) \sim \omega^{3/4}$ is characteristic for semiflexible polymer networks, and betrays the stiffness to bending of the constituent polymers.

Appendices

A Fourier transform convention

Let $f(\cdot)$ be a complex-valued function of a real continuous variable $x \in [0, L]$. Then, f may be decomposed into Fourier modes f_q as

$$f(x) = \frac{1}{L} \sum_{q} f_{q} e^{iqx} \qquad \left(q = \frac{n\pi}{L}; n = 1, 2, \cdots\right).$$
 (86)

The inverse transformation is given by

$$\mathbf{f}_q = \int_0^L \mathrm{d}x \, \mathbf{f}(x) \, e^{-iqx} \,. \tag{87}$$

With these conventions, the Kronecker delta is resolved as

$$\frac{1}{L} \int_{0}^{L} dx \, e^{i(q-q')x} = \delta_{qq'} \,, \tag{88}$$

and the Dirac delta function as

$$\frac{1}{L}\sum_{q}e^{i(x-x')q} = \delta(x-x').$$
(89)

The continuum limit of the q-sum, with these conventions, is to be taken as

$$\sum_{q} \to \frac{L}{\pi} \int_{0}^{\infty} \mathrm{d}q \,. \tag{90}$$

References

- [1] Doi, M., Soft Matter Physics (Oxford University Press, 2013)
- [2] Gittes, F., and F. C. MacKintosh. *Dynamic shear modulus of a semiflexible polymer network*, Phys. Rev. E **58**(2) R1241 (1998).
- [3] Storm, C., Pastore, J. J., MacKintosh, F. C., Lubensky, T. C., and Janmey, P. A. *Nonlinear elasticity in biological gels*, Nature, 435(7039), 191-194 (2005).

C 5 Motor Proteins and Cytoskeletal Filaments: from Motility Assays to Active Gels

T. Auth

Theoretical Soft Matter and Biophysics Institute of Complex Systems and Institute for Advanced Simulation Forschungszentrum Jülich GmbH

Contents

1	Introduction	2
2	Motility Assays	7
3	Active Gels	9
4	Conclusions and Outlook	11

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

The cytoskeleton of biological cells is a fascinating example for living matter. Besides its scaffolding function that provides mechanical stability, it drives a multitude of active processes– see also chapter C6–, such as cell motility [1, 2, 3, 4, 5, 6, 7], cell division [8, 9, 10, 11, 12, 13, 14], uptake processes [15, 16, 17, 18, 19, 20, 21, 22], and cytoplasmic streaming [23, 24, 25, 26]. Three of the main components are the polar cytoskeletal filaments, microtubules and actin, the molecular motors, kinesins, dyneins, and myosins, and crosslinking proteins [27].

In vivo experiments often depend on a large number of components, therefore the comparison of a physics-based bottom-up approach with experimental data may not be as direct as for well-controlled model systems. However, for example the experimental data on cytoplasmic streaming in oocytes resembles aspects of the active motion observed for motor-filament mixtures in circular confinement. A common finding in both, experiment and simulation, is for example an enhanced filament motion close to the confinement [26, 23]. Here, stable cortically anchored microtubules may serve as tracks for filament motion [24]. Figure 1 (a,b) shows an image of a *Drosophila* oocyte next to a computer simulation snapshot of a motor-filament mixture in circular confinement.

In vitro experiments using biological building blocks allow us to systematically investigate motor-filament systems. One specific filament-motor model system are so-called actin or microtubule motility assays. Here, the filaments are propelled on a bed of motors that are grafted to a substrate. Figure 1 (c) shows a schematic illustration for actin filaments bound to a microscope coverslip. A second model system are filament-motor mixtures confined to two dimensions, e.g. microtubules next to an water-oil interface. Figure 1 (d) shows a schematic illustration of a 2D microtubule bundle that is kept together by both, kinesin motor complexes and a depletion-mediated attraction. Because in this system the motors act between two filaments instead of between a filament and the substrate, this model system more closely resembles dynamics in a cytoskeleton.

From a physics point of view, the filaments can either be modeled as semiflexible polymers [28, 29, 30, 31], or-within the so-called Mikado model-as rigid rods [23, 32, 33, 34, 35]. For simplicity, we focus on rods to model filaments throughout this chapter. Motors and crosslinkers can be modeled as springs that attach to the filaments [23, 36, 31, 37]. While crosslinkers bind to fixed positions, the molecular motors move along the filaments; they transport cargo or slide filaments relative to each other. A major challenge for modeling, and in particular for cell-scale computer simulations, is the large size difference between the molecular motors and the cytoskeletal filaments. While typical stalk lengths of the motors are shorter than 100 nm, microtubules for instance can easily reach lengths up to several micrometers.

In Secs. 1 and 1, I introduce molecular motors and cytoskeletal filaments, respectively. In Sec. 2, I describe both results from experiments and computer simulations for motility assays with cytoskeletal filaments that are propelled on a substrate covered with molecular motors. In Sec. 3, I discuss motor-filament systems where motors crosslink and slide the filaments against each other. Finally, in Sec. 4, I summarise the different systems and provide an outlook for future research directions.

Motor Proteins

Motor proteins have a modular design. The motor heads bind to cytoskeletal filaments, lever arms act as force amplifiers, and tails connect the heads with cargo-binding domains-or with



Fig. 1: Motion of cytoskeletal filaments in vivo, in silico, and in vitro. (a) Confocal visualization of autofluorescent, endogeneous vesicles (red) and microtubules (green) in a Drosophila oocyte. The width of the image corresponds to approximately 60 μm. Reproduced from Ref. [26]. (b) A microtubule-motor mixture in circular confinement. The colors correspond to the orientation of the filaments, and the probe trajectories show the motion of single microtubules. Reproduced from Ref. [23]. (c,d) In-vitro model systems containing cytoskeletal filaments and molecular motors. (c) Schematic illustration of a high-density actin motility assay. The molecular motor HMM is immobilized on a coverslip and the filament motion is visualized by the use of fluorescently labelled reporter filaments with a ratio of labelled to unlabelled filaments of 1:200 to 1:320. Reproduced from Ref. [38]. (d) Schematic illustration of an extensile microtubule-kinesin bundle, a basic building block used for the assembly of active matter. Kinesin clusters exert inter-filament sliding forces, whereas depleting PEG polymers induce microtubule bundling. Reproduced from Ref. [39].

another set of heads. As an example, few kinesin motors are shown in Fig. 2. Numerous members of the three groups of myosins, kinesins, and dyneins are known whose members can vary appreciably: mammals, for instance, have genes for over 40 kinesins, 40 myosins, and more than a dozen dyneins [43]. According to their architecture, the motors walk along the filaments in different ways. Furthermore, motors can be both, non-processive and processive. While non-processive motors usually unbind after one step, processive motors, such as conventional kinesin, have two heads that 'hand-over-hand' can move several steps along a filament. A 'duty



Fig. 2: Kinesins. Surface features are rendered based upon atomic resolution structures when available and appear as smooth images for domains of unknown structure. The motor catalytic domains are displayed in blue, mechanical amplifiers in light blue, and tail domains implicated in cargo attachment in purple. Tightly associated motor subunits (light chains) are shown in green. Reproduced from Ref. [40].



Fig. 3: Load and ATP dependence of motor velocity. (a) Average bead velocity of an optical trap attachted to a motor, $v(\text{mean} \pm \text{s.e.m.})$, versus applied load for fixed ATP concentrations (red triangles, left axis, $5 \,\mu\text{M}$ ATP, N = 19 - 57; blue circles, right axis, $2 \,\mu\text{M}$ ATP, N = 37 - 87). The velocity point at (5.6 pN, $5 \,\mu\text{M}$ ATP) is likely to represent an overestimate because beads which stalled completely (v = 0) were indistinguishable from beads lacking active motors, and so were not included in the analysis. (b) Michaelis-Menten kinetics under load. Double logarithmic plot of the average bead velocity, $v(\text{mean} \pm \text{s.e.m.})$, versus ATP concentration for various loads (filled circles, $1.05 \pm 0.01 \,\text{pN}$, $N = 11 - 102 \,\text{runs}$; open circles, $3.59 \pm 0.03 \,\text{pN}$, $N = 8 - 79 \,\text{runs}$; diamonds, $5.63 \pm 0.06 \,\text{pN}$, $N = 19 - 58 \,\text{runs}$). Data were fitted to Michaelis-Menten curves (lines), $v = V_{\text{max}} [ATP] / ([ATP] + K_{\text{m}})$. Inset, fit parameters, V_{max} and K_{m} . Reproduced from Ref. [41].

ratio' can be used to describe the fraction of the time that a motor is attached to a filament [44]. Motors usually consume one molecule of ATP per step of typically few tenths of nanometers.



Fig. 4: Illustration of the dynamics of a motor protein in the simplest two-potential periodic continuum ratchet model in the absence of a load. The vertical arrow represents the input of chemical energy, for example, via the hydrolysis of an ATP molecule; this is followed by diffusion, a drop to the lower potential surface, and further diffusion. Reproduced from Ref. [42].

Because the hydrolysis of one molecule ATP to ADP releases about $20 k_B T$, we can estimate that a motor exerts a force of 1 - 10 pN. Indeed, typical experimentally measured motor forces are found to be in this range. If a load is applied to a motor, the velocity of the attached bead decreases and vanishes for loads of few pN, see Fig. 3 (a). In addition, a dependence of the motor velocity on the ATP concentration has been observed that can be described by a Michaelis-Menten kinetics [41], see Fig. 3 (b).

The mechano-chemical cycle of the motor can lead to stepping in both, forward and backward direction [45]. High backward loads can even induce processive backward stepping. For kinesins, a typical time of microseconds has been reported for each step, with dwell times inbetween the steps that depend on the ATP concentration. One physics-oriented approach to motorprotein motion are ratchets, i.e. periodic, but asymmetric free-energy landscapes [42, 46, 47]. Under the input of chemical energy, the motor can switch stochastically between both potentials. Theoretically the system can thus be described by a set of coupled Fokker-Planck equations. Contrary to purely Brownian ratchets, this coupled set of ratchets requires chemical energy, see Fig. 4.

Due to the mechanical part of the mechano-chemical cycle, motor motion may also be coupled mechanically. Collective effects have for example been studied experimentally for the cooperative extraction of membrane nanotubes by molecular motors [48]. Here, motors that have been individually attached to the membrane have been observed to cluster at the tips of tubes. Furthermore, several motors enlarge the distance that the cluster of motors can walk processively. Theoretical studies on the cooperative transport of cargo pulled by several motors show that a maximal number of only 7-8 kinesin motors that can be simultaneously attached to a cargo particle are sufficient to attain walking distances in the centimeter range [49]. For a tug-of-war situation, where groups of motors are oriented such that they can pull in different directions, both analytical calculations as well as computers simulations show that directed transport may emerge [50, 51].

Cytoskeletal filaments

The cytoskeletal filaments that are discussed in this chapter are actin and microtubules, see Fig. 1. Globular G-actin monomers have a diameter of 4 - 7 nm and can polymerize to form



Fig. 5: Actin and microtubule motility assays vs. self-propelled rods in two dimensions. (a) Above a critical density, wave-like structures are found in actin motility assays, compare Fig. 1. Reproduced from Ref. [38]. (b,c) Microtubule motility assays show either alignment or crossing events of microtubules. Reproduced from Ref. [57]. (d-f) Typical simulation snapshots from quasi two-dimensional simulations of rods: (d) isotropic phase for passive rods, (e) laning phase for high densities of self-propelled rods, and (f) giant cluster phase for self-propelled rods. The colors indicate the directions of the rods. Reproduced from Ref. [32].

polar, helical F-actin filaments with a diameter of about 7 nm. Here, the polarity affects both, walking of the motors on the filament as well as filament polymerization and depolymerization. In experiments, F-actin persistence lengths of $15 - 20 \,\mu\text{m}$ have been measured [52, 53]. Micro-tubules are formed by polymerization of α -tubulin and β -tubulin monomers that first dimerize and then polymerize to form hollow fibers that consist of 13 protofilaments linearly aggregated dimers. A microtubule has an outer diameter of about 24 nm, an inner diameter of about 12 nm, and a persistence length of thousands of micrometers [54, 53]. Thermal fluctuations of grafted microtubules, however, hint that shorter microtubules with lengths between 2.6 and 47.5 μ m can have significantly shorter persisence lengths have been reported for microtubule tips [56]. Further details on cytoskeletal filaments can be found in chapter C3.



Fig. 6: Penetrable self-propelled rods in two dimensions. (a) Schematic representation of the model of a self-propelled rod and coordinates used to characterize its position, orientation, propulsion force, and velocity. The rod is discretized into n_b beads to calculate the rod-rod interaction. (b) Potential profile of a rod along its long axis. Tics on the horizontal axis show the position of beads, separated from each other by r_{\min} . In our simulations, we use $n_b = 18$. (c) Crossing probability for two rods as a function of their crossing angle ϕ (as defined in the schematic in the inset). For each angle, 1000 simulations have been performed. The simulations are divided into 10 groups and the error bars are calculated as the standard deviation of the mean for these groups. The experimental data are taken from Ref. [57]. Reproduced from Ref. [32].

2 Motility Assays

In motility assays with cytoskeletal filaments, molecular motors are grafted to a microscope cover slide, see Fig. 1 (c). When polar cytoskeletal filaments are dispersed on these cover slides, they get propelled by the motors. Both actin and microtubule motility assays have been studied, microscopy images and corresponding computer simulation snapshots are shown in Fig. 5. Various phases of the system have been observed, such as an isotropic order of the rods, cluster formation, swirling, and band-like structures [57, 38]. Because the filaments in the experiments are attached to the substrate using flexible molecular motors, this is not a 2D system, but rather a system of filaments moving in a small slab above the cover slip. While the flexible and thin actin filaments have a high probability to cross each other, the thicker and more rigid microtubules have smaller crossing probabilities.

Actin and microtubule motility assays can be viewed as experimental realisations of self-propelled rods (SPR) [32, 31, 58, 59, 60]. Figure 5 therefore also displays simulations snapshots for self-propelled penetrable rods studied in two dimensions. In computer simulations, an isotropic phase has been observed below the Onsager transition for passive rods [32, 31], see Fig. 5 (d). Self-propelled rods at the same density form giant clusters, see Fig. 5 (e). At very high rod densities the so-called laning phase is found, where lanes of rods moving in opposite directions, see Fig. 5 (f).

In computer simulations, the filaments in a thin 3D slab can be modeled as penetrable rods that are discretized into beads that interact via a separation-shifted Lennard-Jones potential (SSLJ) [32, 61]

$$\phi(r) = \begin{cases} 4\epsilon \left[\left(\frac{\sigma^2}{\alpha^2 + r^2} \right)^6 - \left(\frac{\sigma^2}{\alpha^2 + r^2} \right)^3 \right] + \phi_0 & r \le r_{\text{cut}} \\ 0 & r > r_{\text{cut}} \end{cases}$$
(1)

Here, r is the distance between two beads, see Fig. 6. The parameter α characterizes the capping of the potential, and ϕ_0 shifts the potential to avoid a discontinuity at $r = r_{\text{cut}}$. $E = \phi(0) - \phi(r_{\text{cut}})$ is the potential energy barrier for overlap of two beads. Once E has been set to a certain value, we obtain $\epsilon = \alpha^{12} E/(\alpha^{12} - 4\alpha^6 \sigma^6 + 4\sigma^{12})$. The length $\alpha = \sqrt{2^{1/3}\sigma^2 - r_{\text{cut}}^2}$ is calculated by requiring the potential to vanish at the minimum of the SSLJ potential, $\sigma/r_{\text{cut}} = 2.5$, hence the potential is purely repulsive. The effective bead radius is $r_{\text{bead}} = r_{\text{cut}}/2$, which determines the effective rod thickness r_{cut} and the rod aspect ratio L/r_{cut} . The number n of beads per rod is chosen such that the rod-rod interaction potential is smooth enough that no interlocking between rods occurs, see Fig. 6.

For the Brownian dynamics simulations, we decompose the rod velocity into parallel and perpendicular components with respect to its axis, $\mathbf{v}_{rod,i} = \mathbf{v}_{rod,i,\parallel} + \mathbf{v}_{rod,i,\perp}$, see Fig. 6. In each simulation step, the velocities are calculated using

$$\mathbf{v}_{\mathbf{rod},i,\parallel}(t) = \frac{1}{\gamma_{\parallel}} \left(\sum_{j \neq i}^{N_{\mathbf{rod}}} \mathbf{F}_{ij,\parallel} + \xi_{\parallel} \mathbf{e}_{\parallel} + F_{\mathbf{rod}} \mathbf{e}_{\parallel} \right) , \qquad (2)$$

$$\mathbf{v}_{\mathbf{rod},i,\perp}(t) = \frac{1}{\gamma_{\perp}} \left(\sum_{j \neq i}^{N_{\mathbf{rod}}} \mathbf{F}_{ij,\perp} + \xi_{\perp} \mathbf{e}_{\perp} \right) , \qquad (3)$$

and

$$\omega_{\mathbf{rod},i}(t) = \frac{1}{\gamma_r} \left(\sum_{j \neq i}^{N_{\mathbf{rod}}d} M_{ij} + \xi_r \right) , \qquad (4)$$

where \mathbf{e}_{\parallel} and \mathbf{e}_{\perp} are unit vectors parallel and perpendicular to the rod axis, respectively. F_{rod} is the propulsion force for each rod. The friction coefficients are given by $\gamma_{\parallel} = \gamma_0 L_{rod}$, $\gamma_{\perp} = 2\gamma_{\parallel}$, and $\gamma_r = \gamma_{\parallel} L_{rod}^2/6$, where L_{rod} is the rod length [58]. The random values ξ_{\parallel} , ξ_{\perp} , and ξ_r for the forces in parallel and perpendicular direction and for the torque are drawn from Gaussian distributions with variances $\sigma_{rod}^2 L_{rod}$, $2\sigma_{rod}^2 L_{rod}$, and $\sigma_{rod}^2 L_{rod}^3/12$, respectively. For thermal noise, the variances are calculated using $\sigma_{rod}^2 = 2k_B T/\gamma_0 dt$.¹ Finally, \mathbf{F}_{ij} and M_{ij} are the force and torque from rod j to rod i, calculated using Eq. (1). Hydrodynamic interactions between the rods are largely screened, because of the nearby wall and the high rod density [62, 63, 64, 65, 66], and hence are neglected in the simulations.

There are three different types of energies in our systems: the thermal energy $k_{\rm B}T$, the propulsion strength $F_{\rm p}L$, and the energy barrier *E* due to bead-bead interactions. Dimensionless ratios can be used in order to characterize the importance of the different contributions. The Péclet number [32, 2],

$$Pe = \frac{LF_{p}}{k_{B}T},$$
(5)

is the ratio of propulsion strength to noise. The dimensionless ratio that compares the product of propulsion force and rod length with the energy barrier between rods is the penetrability coefficient [32], LE

$$Q = \frac{LF_p}{E}.$$
 (6)

¹In biological and synthetic SPR systems, the noise arises from the environmental noise, for example, from density fluctuations of signaling molecules for chemotactic swimmers or from motor activity. In this case, the noise does not have to be proportional to $k_{\rm B}T$ and the coupling between translational and rotational noise might be different.



Fig. 7: Active microtubule networks exhibit internally generated flows. An active microtubule network–in detail shown in Fig. 1–viewed on a large scale. Arrows indicate local bundle velocity direction. The scale bar corresponds to $80 \,\mu\text{m}$. Reproduced from Ref. [39].

A comparison of the simulation results for two rods with systematic experiments on crossing of microtubules reveals that our simulations well reproduce the dependence of the crossing probability on the angle under which the filaments collide, see Fig. 6. Here, the highest crossing probability is found if the filaments collide under a 90° angle, when the total energy barrier to cross is lowest. For smaller or larger collision angles, several beads may overlap simultaneously, such that the energy barrier increases. Although we find that several combinations of Q and Pe can lead to similar crossing probabilities, the parameters in the simulations can be adjusted in order to reproduce the experimentally determined two-microtubule interactions.

3 Active Gels

Filament-motor mixtures are active gels, i.e., soft materials with locally broken detailed balance [67]. An experimental model system is for example the microtubule-motor mixture shown in Fig. 1, where–contrary to the motility assays described in the previous section–the motors mutually propell the filaments relative to each other. The microscopic image in Fig. 7 shows both, the structure of the microtubule network that consists in this case of extensile microtubule-kinesin bundles, as well as the local bundle velocity.

Computer simulations that have been inspired by this experiment are the filament-motor mixtures in circular confinement shown in Ref. [23]. Here, short rod-like filaments are mixed with molecular motors and both structure and dynamics of the system is studied in detail. Interestingly, a strong correlaction between the properties of the molecular motors that connect the filaments and the structure and dynamics of the system is observed. In particular, dimeric and tetrameric motors have been simulated, see Fig. 8. While dimeric motors walk on one filament



Fig. 8: Molecular motor motion and corresponding structures of motor-filament mixtures. (a) Schematic showing the effects of tetrameric and dimeric motors on polar-aligned and antialigned filaments. Motor arms are shown to move from position (1) to position (2), in the direction of filament polarization, represented by the yellow marking at the filament tip. The gray representations show the initial attachment positions of motors. Active motor-arms that move on the filament direction of polarization are colored green, and immobile, anchored motorarms are colored red. Tetrameric motors have two motile arms on either cross-linking filament. Dimeric motors have an anchored arm and a motile arm. (A) Tetrameric motors crosslinking polar-aligned filaments induce small velocities. (B) Tetrameric motors cross-linking anti-aligned filaments induce larger velocities. (C) Correlated dimeric motors cross-linking polar-aligned filaments have the same effect as (D) dimeric motors cross-linking anti-aligned filaments. (E) Uncorrelated dimeric motors cross-linking polar-aligned filaments act antagonistically, (F) whereas uncorrelated dimeric motors cross-linking anti-aligned filaments act cooperatively. (b,c) Stationary configuration of rods with (b) dimeric motors and (c) tetrameric motors within circular confinement. Here, the filaments experience weak, depletion-mediated mutual attraction. The number of motors that crosslinks filaments at each time in the system is chosen to be equal to the number of filaments. The colors represent their orientation of the rods with respect to the radial direction from the center. Small black dots represent the position of motor heads. Reproduced from Ref. [23].

and bind to a fixed position on a second filament, tetrameric motors walk on both filaments.

Figure 8 (a) shows the mechanism by that either one or two motors slide two filaments relative to each other. Single tetrameric motors do not slide two parallel filaments, but do slide two antiparallel filaments. Single dimeric motors are always able to slide filaments. For two motors connecting two filaments, again tetrameric motors slide only antiparallel filaments. For dimeric motors, four different cases have to be distinguished, out of which in three cases sliding is predicted, while in one case the motors act antagonistically and only crosslink the filaments. In


Fig. 9: Translational mean-squared displacements of filament centers of mass for dimeric and tetrameric motors for different motor concentrations, N_m/N_f , where N_m is the number of motors and N_f is the number of filaments in the system. The mean-squared displacements are normalized using the squared rod length and lag time is normalized using the onset-of-activity time τ , i.e., the lag time where motor stresses are manifested in the MT dynamics. Reproduced from Ref. [23].

all cases, motors of both types mediate an effective attraction between the filaments in addition to a depletion-mediated attraction between the filaments that is present in the experiments.

Figure 8 (b,c) show simulation snapshots for polar, rod-like filaments in circular confinement, mixed with either dimeric or tetrameric motors. In both cases, the filaments directly at the confinement experience a larger depletion-mediated attraction with the confinement than the depletion-mediated filament-filament attraction, which can mimick cortically attached cytoskele-tal filaments that have been found for example in *Drosophila* oocytes [24]. Dimeric motors lead to the formation of large, polar-ordered domains, while tetrameric motors lead to the formation of polar ordered filament 'packages'.

The dynamics in filament-motor systems can for instance be quantified using the mean-squared displacement of the filaments. Figure 9 shows that thermal motion dominates at short time scales. For passive systems, the ratio of the mean-squared displacements and the lag time decreases with increasing lag time due to steric interactions between the filaments. For systems with motors and lag times that are longer than a threshold lag time, activity leads to superdiffusive filament motion. For dimeric motors, the mean-squared filament displacements are much larger than for tetrameric motors at the same motor densities.²

4 Conclusions and Outlook

Metabolically-driven living matter enables physicists to systematically study non-equilibrium systems. Filament-motor systems, such as the cytoskeleton of bioloical cells that provides structural stability and drives active processes, are also relevant to understand life at the level of its basic building blocks, the cells.

At the level of single molecular motors, approaches that take into account for both the me-

²For times longer than those where the mean-squared displacement reaches the confinement size, the ratio of the mean-squared displacements and the lag time decreases again with increasing lag time due to the maximum possible mean-squared displacement.

chanics and the chemical processes of the motors can model single-motor motion. Interesting physics-based aspects include typical forces that can be exerted by the motors, backstepping of motors in case of a mechanical backwards load, and cooperative interaction of several motors.

Motility assays where cytoskeletal filaments are propelled on layer of motors that are grafted to a substrate allow us to investigate structure formation and dynamics of 'self-propelled' cy-toskeletal filaments. For example cluster formation, wave-like structures, and swirls have been observed in experiments. The flexibility of the stalks of the motors makes motility assays systems where the filaments move in a thin slab above the substrate, not in a plane. Therefore the systems are three-dimensional thin slabs, which can be mimicked by simulating penetrable rods in two dimensions.

Two-dimensional systems of filaments and motors, such as microtubules next to an oil-water interface, are model systems that are much closer to the cytoskeleton in biological cells than motility assays. Here the motors can slide filaments relative to each other, such that persistent motion can emerge. The dynamics on the system level is determined by the filament properties, the confinement, and the properties of the molecular motors. For example, dimeric motors have been shown to lead to much more dynamic systems than tetrameric motors.

Towards a better understanding biological systems in vivo with physics-based approaches, the model systems described in this chapter have to be extended in various ways, see also chapters C6, D4, D6, D7, and E3. Biological cells, for instance, live in three dimensions instead of only two. Although two-dimensional models can reproduce motion of cytoskeletal filaments that resembles cytoplasmic streaming, a more accurate model for biological cells is expected to take into account for the third dimension. Furthermore, in cells the cytoskeleton is bounded by a deformable membrane instead of a hard wall, such that the confinement in the model systems should be able to react to stresses in the cytoskeletal network. Finally, cells interact with their environment and respond to both mechanical and chemical stimuli, as well as to the elastic properties of the environment. However, one suitable approach to address in particular complex biological systems is to use computer simulations based on continuum and coarsed-grained models.

Acknowledgements

This chapter is mainly based on the PhD projects of Masoud Abkenar [32] and Arvind Ravichandran [23], and benefitted substantially from the simulations performed by G. Vliegenthart. Helpful discussion with C. Abaurrea Velasco, G. Saggiorato, J. Elgeti, and G. Gompper are acknowledged.

References

- M. Nickaeen, I. L. Novak, S. Pulford, A. Rumack, J. Brandon, B. M. Slepchenko, and A. Mogilner, PLoS Comp. Biol. 13, e1005862 (2017).
- [2] C. Abaurrea Velasco, S. D. Ghahnaviyeh, H. N. Pishkenari, T. Auth, and G. Gompper, Soft Matter 13, 5865 (2017).
- [3] P. Sens and J. Plastino, J. Phys.: Condens. Matter 27, 273103 (2015).
- [4] L. Blanchoin, R. Boujemaa-Paterski, C. Sykes, and J. Plastino, Physiol. Rev. 94, 235 (2014).

- [5] E. L. Batchelder, G. Hollopeter, C. Campillo, X. Mezanges, E. M. Jorgensen, P. Nassoy, P. Sens, and J. Plastino, Proc. Natl. Acad. Sci. U.S.A. 108, 11429 (2011).
- [6] R. Sambeth and A. Baumgaertner, Phys. Rev. Lett. 86, 5196 (2001).
- [7] A. Mogilner and G. Oster, Biophys. J. 71, 3030 (1996).
- [8] E. Fischer-Friedrich, A. A. Hyman, F. Jülicher, D. J. Müller, and J. Helenius, Sci. Rep. 4 (2014).
- [9] H. Turlier, B. Audoly, J. Prost, and J.-F. Joanny, Biophys. J. 106, 114 (2014).
- [10] R. Tsukanov, G. Reshes, G. Carmon, E. Fischer-Friedrich, N. Gov, I. Fishov, and M. Feingold, Phys. Biol. 8, 066003 (2011).
- [11] J. Gregan, S. Polakova, L. Zhang, I. M. Tolić-Nørrelykke, and D. Cimini, Trends Cell Biol. 21, 374 (2011).
- [12] J. Pecreaux, J.-C. Röper, K. Kruse, F. Jülicher, A. A. Hyman, S. W. Grill, and J. Howard, Curr. Biol. 16, 2111 (2006).
- [13] G. Goshima, F. Nédélec, and R. D. Vale, J. Cell Biol. 171, 229 (2005).
- [14] I. M. Tolić-Nørrelykke, L. Sacconi, G. Thon, and F. S. Pavone, Curr. Biol. 14, 1181 (2004).
- [15] D. M. Richards and R. G. Endres, Proc. Natl. Acad. Sci. U.S.A. 113, 6113 (2016).
- [16] D. M. Richards and R. G. Endres, Biophys. J. 107, 1542 (2014).
- [17] D. T. Kovari, W. Wei, P. Chang, J.-S. Toro, R. F. Beach, D. Chambers, K. Porter, D. Koo, and J. E. Curtis, Biophys. J. 111, 2698 (2016).
- [18] L. Sanchez, P. Patton, S. M. Anthony, Y. Yi, and Y. Yu, Soft Matter 11, 5346 (2015).
- [19] F. Santoro, S. Dasgupta, J. Schnitker, T. Auth, E. Neumann, G. Panaitov, G. Gompper, and A. Offenhausser, ACS Nano 8, 6713 (2014).
- [20] M. Irmscher, A. M. de Jong, H. Kress, and M. W. Prins, J. R. Soc. Interface 10, 20121048 (2013).
- [21] S. Tollis, A. E. Dart, G. Tzircotis, and R. G. Endres, BMC Syst. Biol. 4, 149 (2010).
- [22] J. A. Champion and S. Mitragotri, Proc. Natl. Acad. Sci. U.S.A. 103, 4930 (2006).
- [23] A. Ravichandran, G. A. Vliegenthart, G. Saggiorato, T. Auth, and G. Gompper, Biophys. J. 113, 1121 (2017).
- [24] W. Lu, M. Winding, M. Lakonishok, J. Wildonger, and V. I. Gelfand, Proc. Natl. Acad. Sci. U.S.A. 113, E4995 (2016).
- [25] F. G. Woodhouse and R. E. Goldstein, Proc. Natl. Acad. Sci. U.S.A. 110, 14132 (2013).
- [26] S. Ganguly, L. S. Williams, I. M. Palacios, and R. E. Goldstein, Proc. Natl. Acad. Sci. U.S.A. 109, 15109 (2012).
- [27] B. Alberts, A. Johnson, J. Lewis, D. Morgan, and M. Raff, *Molecular Biology of the Cell* (Garland Science, New York, NY, 2015), 6th ed., ISBN 9780815344322.
- [28] T. Eisenstecken, G. Gompper, and R. G. Winkler, J. Chem. Phys. 146, 154903 (2017).
- [29] J. Pešek, P. Baerts, B. Smeets, C. Maes, and H. Ramon, Soft Matter 12, 3360 (2016).
- [30] R. E. Isele-Holder, J. Elgeti, and G. Gompper, Soft Matter 11, 7181 (2015).
- [31] P. Kraikivski, R. Lipowsky, and J. Kierfeld, Phys. Rev. Lett. 96, 258103 (2006).
- [32] M. Abkenar, K. Marx, T. Auth, and G. Gompper, Phys. Rev. E 88, 062314 (2013).
- [33] K. Kroy, Curr. Opin. Colloid Interface Sci. 11, 56 (2006).
- [34] J. Kierfeld, O. Niamploy, V. Sa-Yakanit, and R. Lipowsky, Eur. Phys. J. E 14, 17 (2004).
- [35] J. Wilhelm and E. Frey, Phys. Rev. Lett. 91, 108103 (2003).
- [36] R. Blackwell, O. Sweezy-Schindler, C. Baldwin, L. E. Hough, M. A. Glaser, and M. Betterton, Soft Matter 12, 2676 (2016).
- [37] F. Nédélec, J. Cell Biol. 158, 1005 (2002).
- [38] V. Schaller, C. Weber, C. Semmrich, E. Frey, and A. R. Bausch, Nature 467, 73 (2010).

- [39] T. Sanchez, D. T. Chen, S. J. DeCamp, M. Heymann, and Z. Dogic, Nature 491, 431 (2012).
- [40] R. D. Vale, Cell 112, 467 (2003).
- [41] K. Visscher, M. J. Schnitzer, and S. M. Block, Nature 400, 184 (1999).
- [42] A. B. Kolomeisky and M. E. Fisher, Annu. Rev. Phys. Chem. 58, 675 (2007).
- [43] M. Schliwa and G. Woehlke, Nature 422, 759 (2003).
- [44] J. Howard, Nature 389, 561 (1997).
- [45] N. J. Carter and R. Cross, Nature 435, 308 (2005).
- [46] F. Jülicher, A. Ajdari, and J. Prost, Rev. Mod. Phys. 69, 1269 (1997).
- [47] R. D. Astumian and M. Bier, Phys. Rev. Lett. 72, 1766 (1994).
- [48] C. Leduc, O. Campàs, K. B. Zeldovich, A. Roux, P. Jolimaitre, L. Bourel-Bonnet, B. Goud, J.-F. Joanny, P. Bassereau, and J. Prost, Proc. Natl. Acad. Sci. U.S.A. 101, 17096 (2004).
- [49] S. Klumpp and R. Lipowsky, Proc. Natl. Acad. Sci. U.S.A. 102, 17284 (2005).
- [50] M. J. Müller, S. Klumpp, and R. Lipowsky, Proc. Natl. Acad. Sci. U.S.A. 105, 4609 (2008).
- [51] F. Jülicher and J. Prost, Phys. Rev. Lett. 75, 2618 (1995).
- [52] A. Ott, M. Magnasco, A. Simon, and A. Libchaber, Phys. Rev. E 48(3), R1642 (1993).
- [53] F. Gittes, B. Mickey, J. Nettleton, and J. Howard, J. Cell Biol. 120(4), 923 (1993).
- [54] M. Elbaum, D. K. Fygenson, and A. Libchaber, Phys. Rev. Lett. 76(21), 4078 (1996).
- [55] F. Pampaloni, G. Lattanzi, A. Jonáš, T. Surrey, E. Frey, and E.-L. Florin, Proc. Natl. Acad. Sci. U.S.A. 103(27), 10248 (2006).
- [56] M. Van den Heuvel, S. Bolhuis, and C. Dekker, Nano Lett. 7(10), 3138 (2007).
- [57] Y. Sumino, K. H. Nagai, Y. Shitaka, D. Tanaka, K. Yoshikawa, H. Chaté, and K. Oiwa, Nature 483(7390), 448 (2012).
- [58] Y. Yang, V. Marceau, and G. Gompper, Phys. Rev. E 82(3), 031904 (2010).
- [59] A. Baskaran and M. C. Marchetti, Phys. Rev. E 77(1), 011920 (2008).
- [60] F. Peruani, A. Deutsch, and M. Bär, Phys. Rev. E 74(3), 030904 (2006).
- [61] M. E. Fisher and D. Ruelle, J. Math. Phys. 7(2), 260 (1966).
- [62] J. Elgeti, R. G. Winkler, and G. Gompper, Rep. Prog. Phys. 78, 056601 (2015).
- [63] J. Elgeti and G. Gompper, EPL (Europhys. Lett.) 101(4), 48003 (2013).
- [64] J. Elgeti, U. B. Kaupp, and G. Gompper, Biophys. J. 99(4), 1018 (2010).
- [65] J. Elgeti and G. Gompper, EPL (Europhys. Lett.) 85(3), 38002 (2009).
- [66] A. P. Berke, L. Turner, H. C. Berg, and E. Lauga, Phys. Rev. Lett. 101(3), 038102 (2008).
- [67] J. Prost, F. Jülicher, and J.-F. Joanny, Nat. Phys. 11(2), 111 (2015).

C 6 Hydrodynamics of the active cytoskeleton

Karsten Kruse Departments of Biochemistry and Theoretical Physics NCCR Chemical Biology University of Geneva

Contents

1	Introduction						
2	Generalized hydrodynamics of active gels						
	2.1 Hydrodynamic modes				3		
	2.2 The dissipation rate				4		
	2.3 The constitutive equations				4		
3	Retrograde flow in lamellipodia 5						
	3.1 Dynamic equations for an actin slab				5		
	3.2 The retrograde flow				7		
4	Contraction of a poroelastic active gel						
	4.1 The dynamic equations of a poroelastic active gel				8		
	4.2 Active contraction of a circular symmetric disk				9		
	4.3 The limit of small deformations				10		
	4.4 Large deformations				11		
5	Conclusions						
A	Discretization scheme for the numerical solution of the active poroelastic dynamic						
	equations				13		

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

From a physical point of view, the spontaneous emergence of flows and topological point defects through the internal conversion of chemical energy are among the most exciting features of the cytoskeleton and of filament-motor systems in general [1, 2]. In contrast to conventionally studied polymer solutions, these phenomena are not due to the application of an external field, but result from processes within the material. Specifically, the free energy released in the process of Adenosine-Triphosphate (ATP) hydrolysis generates conformational changes of motor proteins and affects filament assembly. The resulting force dipoles produce an "active stress", which drives the aforementioned processes.

Various physical approaches have been pursued to describe cytoskeletal dynamics. The most detailed descriptions are agent-based stochastic processes that are usually studied numerically. A prominent example of this approach is given by the Cytosim package, which includes a broad variety of cytoskeletal elements [3]. To give but two examples, it has been used to study the self-organization of microtubules and kinesins into asters and vortices [1] and to investigate the conditions under which an actomyosin network contracts [4]. In addition, stochastic simulations are developed with the specific purpose to study a concrete phenomenon, for example, the emergence of actin polymerization waves [5] or cortical actin length distributions [6].

Kinetic descriptions can be seen as mean-field theories of the stochastic processes just mentioned. They are continuum descriptions, where the expressions for the currents and the reactions are motivated by known molecular processes. However, they typically neglect molecular details and focus on the (hopefully) essential aspects of the molecular interactions and dynamics. Such approaches have been used to study possible mechanisms of stress generation in filament-motor networks [7, 8], the possible role of actin-polymerization waves for cell migration [9], and pattern formation in contractile rings [10].

Finally, phenomenological descriptions for macroscopic cytoskeletal dynamics have been developed [11]. Non-equilibrium thermodynamics, which is also known as generalized hydrodynamics, provides a framework for the systematic derivation of the corresponding dynamic equations if the material is close to thermodynamic equilibrium [12, 13]. This condition is expressed by the requirement that the system is locally at equilibrium. Non-equilibrium thermodynamics is purely based on symmetries and conservation laws. The material properties are captured by constitutive equations that are obtained from a systematic expansion of the currents in terms of thermodynamic forces. The latter are in turn expressed in terms of the state variables, such that a closed set of equations is obtained. Notably, the active stress is written as a linear function of the difference between the chemical potentials of ATP and its hydrolysis products. Hydrodynamic equations depend on a number of phenomenological constants, for example, the viscosities of simple fluids. Their values either have to be measured or can be related to molecular parameters by using a kinetic or agent-based description.

The hydrodynamics of active gels – or active gel theory – has been used to study generic physical properties of the cytoskeleton, for example, the dynamics of topological point defects or the spontaneous emergence of flows [11, 14, 15]. In addition, it has been employed to analyse subcellular and tissue dynamics. For example, the retrograde flow in lamellipodia [16], the formation and contraction of cytokinetic rings during cell division [17], or the spreading of epithelia during embryonic development were analyzed in this framework [18].

In this chapter, we will sketch the ideas underlying general hydrodynamics and use it to derive the dynamic equations of an active gel permeated by a solvent. We will then apply the equations to study a possible mechanism for generating the retrograde actin flow in lamellipodia. Finally, we will study the contraction dynamics of an active poroelastic gel.

2 Generalized hydrodynamics of active gels

Generalized hydrodynamics or non-equilibrium thermodynamics provides a general framework for describing systems that are close to thermodynamic equilibrium. In this context, "close to thermodynamic equilibrium" means that the system is locally at thermodynamic equilibrium, but not globally. As a consequence, one can define a free energy density and locally use the powerful concepts and tools of thermodynamics. In this section, we will sketch the application of this approach to active gels. For a detailed introduction into this topic see, for example, Refs. [12, 13]. A self-contained derivation of the hydrodynamic equations of a one-component polar active gel is given in Ref. [19].

2.1 Hydrodynamic modes

As a consequence of the condition that the system should be locally at thermodynamic equilibrium, only long-lived degrees of freedom are considered within non-equilibrium thermodynamics. These "hydrodynamic modes" satisfy the condition that their characteristic relaxation time decreases with an increase of the associated wave-number, $\tau \sim k^{-2}$. An example is provided by a particle density c that obeys the diffusion equation, $\partial_t c = D\partial_x^2 c$. A Fourier-transform in space leads to $\frac{d}{dt}c_k = -Dk^2c_k$, such that $c_k \propto e^{-Dk^2t}$ and $\tau = D^{-1}k^{-2}$.

Typically, hydrodynamic modes result either from conservation laws as in the above example or from a broken continuous symmetry. For active gels, one typically considers the particle numbers of the gel and the solvent (cytosol), of ATP, ADP, and P_i , and momentum as the conserved quantities. The possible macroscopic orientational order of the actin network is a hydrodynamic mode emerging from a broken continuous symmetry. In the following, we will only consider isotropic cases and do not account for a polar or nematic order parameter. Furthermore, we will consider only systems coupled to a heat bath, such that temperature is constant and energy is not conserved. Finally, we will assume that the concentrations of ATP and of its hydrolysis products are constant in space and time. This is probably appropriate for many reconstituted systems *in vitro*, in case they are endowed with an ATP-regeneration system or for living cells that metabolize nutrients to create ATP. Even if there were spatial or temporal fluctuations in these concentrations, they would be irrelevant as long as the availability of ATP is not rate limiting.

Under these conditions, the relevant dynamic equations for the hydrodynamic modes are

$$\partial_t \rho_{\rm gel} + \partial_\alpha \rho_{\rm gel} v_{\rm gel} = 0 \tag{1}$$

$$\partial_t \rho_{\rm sol} + \partial_\alpha \rho_{\rm sol} v_{\rm sol} = 0 \tag{2}$$

$$\partial_t g_\alpha + \partial_\beta \sigma_{\alpha\beta}^{\text{tot}} = 0. \tag{3}$$

In the two continuity equations for the gel and solvent mass densities, ρ_{gel} and ρ_{sol} , respectively, \mathbf{v}_{gel} and \mathbf{v}_{sol} denote the local gel and solvent velocity. In the balance equation for the total momentum $\mathbf{g} = \rho_{gel}\mathbf{v}_{gel} + \rho_{sol}\mathbf{v}_{sol}$, the total momentum flux tensor σ^{tot} is the same as the mechanical stress tensor. For simplicity, we have set all source and sink terms equal to zero, which amounts to neglecting any bulk exchange between the gel and the solvent and to assuming that there are no bulk external forces. This does not exclude the application of forces

or exchange of matter at the surfaces, which are accounted for by boundary conditions. Finally, we have used in these equations Einstein's summation convention for identical indices, for example, $\partial_{\alpha}v_{\alpha} = \sum_{\alpha=1}^{3} \partial_{\alpha}v_{\alpha}$. Cellular processes typically occur at low Reynolds number, such that inertial terms can eventually be neglected in the momentum balance equation; it then expresses force balance. However, for the time being, the inertial terms are kept to couple the mechanical stress tensor to the dynamics, as we will see in the next section.

To close the dynamic equations, one still needs expressions for the currents. In the framework of non-equilibrium thermodynamics, pairs of conjugated generalized fluxes and forces are identified by considering the entropy production rate. At constant temperature T on can equivalently consider the dissipation rate. The fluxes are then expanded up to linear order in terms of the forces.

2.2 The dissipation rate

At constant temperature, the dissipation rate Θ can be expressed in terms of the free energy:

$$\dot{\Theta} = -\frac{\mathrm{d}}{\mathrm{d}t} \int \mathrm{d}\mathbf{r} \left\{ \frac{g_{\alpha}^2}{2\rho} + f\left(n_{\mathrm{gel}}, n_{\mathrm{sol}}, n_{\mathrm{ATP}}, n_{\mathrm{ADP}}, n_{\mathrm{P}}\right) \right\}.$$
(4)

Here, f is the free energy density and n_i denotes the particle number density of species i with $\rho_{\text{gel}} = m_{\text{gel}}n_{\text{gel}}$ and $\rho_{\text{sol}} = m_{\text{sol}}n_{\text{sol}}$, where m_{gel} and m_{sol} are, respectively, the molecular masses of the gel and the solvent. Here, all components are assumed to be liquids or gases. The cytoskeletal filament network really is a viscoelastic material, such that the free energy density should also depend on the strain tensor u, see below, when the case of a poroelastic gel is treated. The general case of a viscoelastic active gel has been considered in Refs. [20, 21]. For the cytoskeleton in live cells, though, the stress relaxation time is on the order of 10 s. On longer times scales, which are relevant for many cellular processes, it is appropriate to describe the cytoskeleton as a viscous fluid. Returning to Eq. (4), exchanging the time derivative and the spatial integration, using the conservation laws (1)-(3) as well as the Gibbs-Duhem relation $dP = \sum_i n_i d\mu_i$, where P is the hydrostatic pressure and $\mu_i = \partial f / \partial n_i$ are the chemical potentials of the various components, we finally arrive at

$$\dot{\Theta} = -\int \mathrm{d}\mathbf{r} \left\{ \sigma^d_{\alpha\beta} v_{\alpha\beta} + j_\alpha \partial_\alpha \bar{\mu} + r\Delta \mu \right\},\tag{5}$$

where only the relevant terms have been kept. In this expression, σ^d is the deviatory stress with components $\sigma_{\alpha\beta}^d = \sigma_{\alpha\beta}^{\text{tot}} - P\delta_{\alpha\beta}$, where we have now neglected inertial terms and where $\delta_{\alpha\beta} = 1$ if $\alpha = \beta$ and zero otherwise. The components of the symmetric part of the strain rate tensor v are $v_{\alpha\beta} = (\partial_{\alpha}v_{\beta} + \partial_{\beta}v_{\alpha})/2$, where v is the center-of-mass velocity, $(\rho_{\text{sol}} + \rho_{\text{gel}}) \mathbf{v} = \rho_{\text{gel}}\mathbf{v}_{\text{gel}} + \rho_{\text{sol}}\mathbf{v}_{\text{sol}}$. The diffusion current j is defined through $\rho_{\text{gel}}\mathbf{v}_{\text{gel}} = \rho_{\text{gel}}\mathbf{v} + m_{\text{gel}}\mathbf{j}$ and the reduced chemical potential $\bar{\mu} = \mu_{\text{gel}}/m_{\text{gel}} - \mu_{\text{sol}}/m_{\text{sol}}$. Finally, r denotes the ATP-hydrolysis rate and $\Delta\mu = \mu_{\text{ATP}} - \mu_{\text{ADP}} - \mu_{\text{P}}$.

2.3 The constitutive equations

From Equation (5), we identify the generalized fluxes σ^d , **j**, and r as well as the respective conjugated generalized forces v, $\nabla \overline{\mu}$, and $\Delta \mu$. Expressing the fluxes up to linear order in terms

of the forces, we obtain the following constitutive equations

$$\sigma_{\alpha\beta}^{d} = 2\eta \left(v_{\alpha\beta} - \frac{1}{d} v_{\gamma\gamma} \delta_{\alpha\beta} \right) + \nu v_{\gamma\gamma} \delta_{\alpha\beta} - \zeta \Delta \mu \delta_{\alpha\beta} \tag{6}$$

$$j_{\alpha} = -\gamma \partial_{\alpha} \bar{\mu} \tag{7}$$

$$r = \Lambda \Delta \mu + \zeta v_{\gamma\gamma}.$$
(8)

The first two terms in Eq. (6) are the usual contributions to the stress of a viscous fluid resulting from pure shear and pure contractile/extensile flows with the corresponding shear viscosity η and the bulk viscosity ν . Furthermore, $\delta_{\alpha\beta}$ is the Kronecker symbol and equates to 1 for $\alpha = \beta$ and to 0 otherwise. Other fluxes and forces of the same tensorial order are coupled by phenomenological constants. Specifically, γ is related to a diffusion constant and Λ determines the hydrolysis rate for a given difference $\Delta \mu$ in the chemical potentials. Most interestingly in the present context, there is a cross-term relating $\Delta \mu$ to the deviatory stress and consequently the rate of strain tensor to the ATP-hydrolysis rate. The coupling coefficient ζ is, up to a sign, the same for both terms as imposed by the Onsager relations.

3 **Retrograde flow in lamellipodia**

Cells crawling on a solid substrate extend a flat protrusion at the leading edge, the lamellipodium. At the leading edge, the membrane is often pushed forward by polymerizing actin. The cell body follows through actin network contraction induced by myosin motors. It has been observed that relative to the substrate, the actin network in the lamellipodium flows backwards from the leading edge. In this section, we will analyze a cartoon version of the lamellipodium to show that this retrograde flow is a generic consequence of contractile active stresses in the actin network. We will neglect permeation of the solvent and consider an effective one-component description of the gel [21, 22]. This description had been introduced in [23].

3.1 **Dynamic equations for an actin slab**

We consider an actin slab of fixed height h moving on a substrate that coincides with the (x, y)plane. Actin polymerizes at x_r at velocity v_p and depolymerizes at x_l at velocity v_d into the direction of positive x, see Fig. 1. The gel assembly dynamics breaks the isotropy of the system by making it globally polar. However, locally the material properties are still isotropic. As we are interested in the behavior on long time scales, we assume the gel to be purely viscous as in the previous section. For simplicity, we will take it to be infinitely compressible.

The only relevant conservation law of the problem is momentum conservation in form of the force balance equation. Given the conditions exposed above, the constitutive equation for the total stress reads

$$\sigma_{\alpha\beta}^{\text{tot}} = \frac{\eta}{2} \left(\partial_{\alpha} v_{\beta} + \partial_{\beta} v_{\alpha} - \partial_{\gamma} v_{\gamma} \delta_{\alpha\beta} \right) + \nu \partial_{\gamma} v_{\gamma} \delta_{\alpha\beta} - \zeta \Delta \mu.$$
(9)

In the following, we will assume translational invariance in the y-direction. Consequently, the dynamic quantities only depend on x and z.



Fig. 1: Simplified description of a lamellipodium. It is represented by a slab of an active gel with constant height h. It polymerizes with velocity v_p at $x = x_r$ and depolymerizes with velocity v_d at $x = x_l$. The steady state velocity of the slab is U.

Let us focus on the steady state and look at force balance on a small strip of gel of width Δx . In the *x*-direction we then get

$$-\int_0^h \sigma_{xx}^{\text{tot}}(x,z) \, dz + \int_0^h \sigma_{xx}^{\text{tot}}(x+\Delta x,z) \, dz \tag{10}$$

$$-\int_{x}^{x+\Delta x} \sigma_{xz}^{\text{tot}}(x',0) \ dx' + \int_{x}^{x+\Delta x} \sigma_{xz}^{\text{tot}}(x',h) \ dx' = 0.$$
(11)

Note that we have assumed that there are no external bulk forces present. The coupling to the environment is taken into account through the boundary conditions on the stress tensor σ^{tot} . We assume a free boundary at the top of the slab, z = h, and friction with the substrate at the bottom, z = 0. These boundary conditions, which imply $\sigma_{xz}^{\text{tot}}(x, 0) = \xi v_x(x, 0)$, where ξ is an effective friction coefficient, and $\sigma^{\text{tot}}(x, h) = 0$, together with the definition

$$\sigma(x) := \frac{1}{h} \int_0^h \sigma_{xx}^{\text{tot}}(x, z) \, dz \tag{12}$$

yield

$$-h\sigma(x) + h\sigma(x + \Delta x) - \xi v(x)\Delta x = 0.$$
(13)

Here, we have used $v(x) = v_x(x, 0)$. In the limit $\Delta x \to 0$ we arrive at

$$\frac{d}{dx}\sigma = \frac{\xi}{h}v.$$
(14)

As the height h is constant and the gel is infinitely compressible, we do not need to consider force balance in z-direction.

From the constitutive equation (9) we obtain

$$\sigma = \tilde{\eta} \frac{d}{dx} v - \zeta \Delta \mu. \tag{15}$$

This completes the definition of the dynamic equations.



Fig. 2: Steady state of the actin slab with $L/\lambda = 10$ and $\zeta \Delta \mu < 0$. Shown are the profiles of the stress (a) and the velocity (b) along the active gel slab.

3.2 The retrograde flow

Combinig Equations (14) and (15) we get

$$\lambda^2 \frac{d^2}{dx^2} \sigma - \sigma = \zeta \Delta \mu, \tag{16}$$

where the characteristic length λ is determined by $\lambda^2 = \eta h/\xi$. The general solution is given by $\sigma(x) = Ae^{x/\lambda} + Be^{-x/\lambda} - \zeta \Delta \mu$. We now use the slab's reference frame, such that $x_l = 0$ and $x_r = L$. The integration constants A and B are fixed by the boundary conditions $\sigma(0) = \sigma(L) = 0$ to be

$$A = \zeta \Delta \mu \frac{\mathrm{e}^{-L/\lambda}}{1 + \mathrm{e}^{-L/\lambda}} \tag{17}$$

$$B = \zeta \Delta \mu \frac{1}{1 + \mathrm{e}^{-L/\lambda}} \tag{18}$$

yielding

$$\sigma(x) = \zeta \Delta \mu \left\{ \frac{\cosh \frac{2x-L}{2\lambda}}{\cosh \frac{L}{2\lambda}} - 1 \right\}.$$
 (19)

For the velocity, which is still measured relative to the substrate, we then have

$$v = \frac{h}{\xi} \frac{d}{dx} \sigma \tag{20}$$

$$=\frac{h\zeta\Delta\mu}{\xi\lambda}\frac{\sinh\frac{2x-L}{2\lambda}}{\cosh\frac{L}{2\lambda}}.$$
(21)

For $\zeta \Delta \mu < 0$ we thus obtain a retrograde flow at the leading edge of the slab and an anterograde flow at the trailing edge, see Fig. 2. Note that in absence of active processes, $\Delta \mu = 0$, there are no flows generated.

Let U denote the constant velocity at which the slab is moving. This implies

$$v_{\rm p} + v(L) = U \tag{22}$$

$$v_{\rm d} + v(0) = U.$$
 (23)

These conditions fix the system length L and the velocity U. Explicitly, we have

$$U = \frac{1}{2} \left(v_{\rm p} + \mathrm{d} \right) \tag{24}$$

$$L = 2\lambda \operatorname{artanh} \frac{\xi \lambda \left(v_{\rm p} - v_{\rm d} \right)}{2h\zeta \Delta \mu}.$$
(25)

This example shows how the coupling of the chemical energy released during ATP-hydrolysis to mechanical stresses leads to flows of the actin cytoskeleton. Similarly, in confining channels, it can lead to spontaneous laminar [15, 19] and more complex flow patterns [24].

4 Contraction of a poroelastic active gel

In this section, we will study the contraction of an active gel that is permeated by a solvent. Indeed, some works suggest the presence of poroelastic effects in cellular cytoskeletal dynamics [25, 26]. These notably comprise the diffusion of stress and the generation of a solvent flow by a contracting actin network. Such solvent flows might play an important role for distributing proteins and other biomolecules that can, for example, be used for assembling organelles or for signalling purposes. In contrast to the previous examples, we will now consider an elastic active gel.

4.1 The dynamic equations of a poroelastic active gel

Consider an isotropic elastic gel that is permeated by a solvent. There are three conserved quantities in this problem, the gel mass, the solvent mass, and momentum. The respective densities ρ_{gel} and ρ_{sol} of the gel and the solvent evolve according to the continuity equations

$$\partial_t \rho_{\rm gel} + \partial_\alpha \rho_{\rm gel} \dot{u}_\alpha = 0 \tag{26}$$

$$\partial_t \rho_{\rm sol} + \partial_\alpha \rho_{\rm sol} v_\alpha = 0. \tag{27}$$

Here, **u** is the displacement field that describes deformations of the gel from its relaxed state in absence of activity, such that $\dot{\mathbf{u}} \equiv \partial_t \mathbf{u}$ is the gel deformation velocity, and **v** the solvent velocity field. To first order in the displacement field, the density ρ_{gel} can be expressed in terms of the initial gel density ρ_0 at t = 0 and the displacement field **u** through

$$\rho_{\rm gel} = \rho_0 \left(1 - \partial_\alpha u_\alpha + \mathcal{O}(u^2) \right), \tag{28}$$

which solves Eq. (26) up to first order in u. Furthermore, the combined gel-solvent system is incompressible, $\rho = \rho_{gel} + \rho_{sol} = const$, such that the gel volume fraction $\phi \equiv \rho_{gel} / \rho$ obeys

$$\partial_{\alpha} \left(\phi \dot{u}_{\alpha} + (1 - \phi) \, v_{\alpha} \right) = 0. \tag{29}$$

The material properties are determined by the constitutive equations for the mechanical stress in the system, which enter the momentum conservation equation. This condition reduces again to force balance, because we consider overdamped dynamics, such that inertial terms can be neglected. It determines the gel displacement and solvent velocity fields, such that the dynamics is fully specified. The force balance condition can be written separately for the gel and for the solvent. Explicitly,

$$-\partial_{\beta}\sigma_{\alpha\beta}^{\rm sol} = \gamma \left(\dot{u}_{\alpha} - v_{\alpha} \right) \tag{30}$$

$$-\partial_{\beta}\sigma_{\alpha\beta}^{\text{gel}} = -\gamma \left(\dot{u}_{\alpha} - v_{\alpha} \right). \tag{31}$$

Here, γ accounts for the friction between the gel and the solvent and has units of a viscosity divided by a length squared. In general, it changes with the gel volume fraction ϕ , but we neglect this dependence for simplicity. The constitutive equations for the gel and solvent stresses result from linear stress-strain and stress-strain rate relations [20, 21]

$$\sigma_{\alpha\beta}^{\rm sol} = 2\eta v_{\alpha\beta} + P\delta_{\alpha\beta} \tag{32}$$

$$\sigma_{\alpha\beta}^{\text{gel}} = K u_{\gamma\gamma} \delta_{\alpha\beta} + 2\mu \left(u_{\alpha\beta} - \frac{1}{3} \delta_{\alpha\beta} u_{\gamma\gamma} \right) - \zeta \Delta \mu \delta_{\alpha\beta}. \tag{33}$$

The hydrostatic pressure acts as a Lagrange multiplier and is determined by the incompressibility condition. The elastic properties of the dry gel in absence of activity are captured by the constant bulk and shear moduli K and μ . We neglect any direct coupling of the activity, that is, $\Delta\mu$, to the solvent flow.

4.2 Active contraction of a circular symmetric disk

Consider a circular symmetric disk of an active gel embedded in a viscous solvent. We use cylindrical coordinates and neglect any dependence on the angular coordinate θ and the height coordinate z. The relevant components of the stress in the gel are then given by

$$\sigma_{rr}^{\text{gel}} = K \left(\partial_r u_r + \frac{u_r}{r} \right) + 2\mu \partial_r u_r - \zeta \Delta \mu \tag{34}$$

$$\sigma_{\theta\theta}^{\text{gel}} = K \left(\partial_r u_r + \frac{u_r}{r} \right) + 2\mu \frac{u_r}{r} - \zeta \Delta \mu.$$
(35)

From these expressions, we obtain the dynamic equation

$$-\gamma \left(\dot{u}_r - v_r \right) = -\left(\nabla \cdot \sigma^{\text{gel}} \right)_r \tag{36}$$

$$= -\left[\partial_r \sigma_{rr}^{\text{gel}} + \frac{1}{r} \left(\sigma_{rr}^{\text{gel}} - \sigma_{\theta\theta}^{\text{gel}}\right)\right]$$
(37)

To determine the radial component of the solvent velocity v_r , we use the continuity equation for the total density ρ expressed in terms of the gel volume fraction

$$\frac{1}{r}\partial_r \left(r\phi \dot{u}_r + r\left(1 - \phi\right) v_r \right) = 0.$$
(38)

We consider a situation, where the gel is embedded in a circular recipient, such that the total material flux at the boundary of the system equals zero. Hence, the continuity equation yields

$$v_r = -\frac{\phi}{1-\phi}\dot{u}_r.$$
(39)

The solvent flow is always directed opposite to the gel displacement velocity. For a gel volume fraction $\phi > 1/2$, the solvent flow is faster than the gel displacement velocity.

According to Equation (28), the gel volume fraction can be expressed in terms of the initial gel volume fraction ϕ_0 and the radial gel displacement field u_r . Using this relation and expression (39) for v_r we can rewrite the dynamic equation (37) in a form that only depends on the displacement vector field u_r :

$$\dot{u}_r = \frac{1-\phi}{\gamma} \partial_r \left[(K+2\mu) \left\{ \partial_r u_r + \frac{u_r}{r} \right\} - \zeta \Delta \mu \right].$$
(40)

Finally, we have to specify the boundary conditions. In the center, clearly u_r (r = 0) = 0. Furthermore, we assume that there are no external forces applied to the border of the gel and the corresponding boundary conditions read σ_{rr}^{gel} $(r = R + u_r(R)) = 0$ and $\sigma_{\theta\theta}^{\text{gel}}$ $(r = R + u_r(R)) = 0$. Here, R is the initial and $R + u_r(R)$ the current radius of the gel. The homogenous active stress $-\zeta\Delta\mu$ only contributes at the boundary. Consequently, we can absorb it into a boundary condition and solve

$$\dot{u}_r = \frac{1-\phi}{\gamma} \partial_r \left[(K+2\mu) \left\{ \partial_r u_r + \frac{u_r}{r} \right\} \right]$$
(41)

with $(K+2\mu) \left\{ \partial_r u_r + \frac{u_r}{r} \right\} \Big|_{r=R} = \zeta \Delta \mu.$

In the following, we will scale space with the initial radius R, time with $\gamma R^2/(K + 2\mu)$, and stresses with $K + 2\mu$. Consequently, the system dynamics depends only on two dimensionless parameters: the initial gel volume fraction ϕ_0 and the dimensionless active stress $\zeta \Delta \mu$, for which we have kept the original notation. Anticipating that the solution of the dynamic equation for $\zeta \Delta \mu < 0$ describe contraction of the gel, we note that the scaling implies that the maximal contraction velocity scales inversely proportionally with the initial system size whereas the characteristic relaxation time scales with R^2 . The latter is consistent with observations of contracting actomyosin networks *in vitro* [27].

The steady state, $\dot{u}_r = 0$ is given by

$$u_r^{(0)}(r) = \frac{\zeta \Delta \mu}{2} r,\tag{42}$$

such that the final strain is proportional to the activity. The gel density as well as the stress are homogenous in the steady state: the gel volume fraction equals $\phi_0(1-\zeta\Delta\mu)$ and the total stress vanishes.

4.3 The limit of small deformations

In general, the dynamic equation (41) is difficult if not impossible to solve analytically. In case the overall deformation of the gel is small, changes in the gel volume fraction ϕ can be neglected. In the dynamic equation (41), we can then set $\phi = \phi_0$ leaving us with a linear equation:

$$\dot{u}_r = (1 - \phi_0) \,\partial_r \left(\partial_r u_r + \frac{u_r}{r} \right). \tag{43}$$

The boundary condition is $\left(\partial_r u_r + \frac{u_r}{r}\right)\Big|_{r=1} = \zeta \Delta \mu$. The radial deformation field is thus determined by the diffusion equation, which depends only through the boundary condition on the active stress $\zeta \Delta \mu$.

To solve the linearized dynamic equation, we write $u_r(r,t) = u_r^{(0)}(r) + w(r,t)$. The function w then also obeys Eq. (43), but now with the boundary condition $\left(\partial_r w + \frac{w}{r}\right)\Big|_{r=1} = 0$. Using the separation *ansatz* $w(r,t) = T(t)\Psi(r)$. We obtain

$$T = e^{-(1-\phi_0)\omega t} \tag{44}$$

$$\frac{\mathrm{d}}{\mathrm{d}r}\left(\Psi' + \frac{\Psi}{r}\right) = -\omega\Psi,\tag{45}$$

where the prime indicates derivation with respect to r. The solution to the second equation is given by $J_1(\sqrt{\omega}r)$, where J_1 is the Bessel function of the first kind. The possible values of ω are determined by the boundary conditions, that is,

$$\sqrt{\omega}J_1'(\sqrt{\omega}) + J_1(\sqrt{\omega}) = 0.$$
(46)

Under this condition, the Bessel functions are pairwise orthogonal:

$$\int_0^1 J_1(\sqrt{\omega}r) J_1(\sqrt{\omega'}r) \ r \mathrm{d}r = \delta_{\omega\omega'}.$$
(47)

The general solution to Eq. (45) is then given by

$$w(r,t) = \sum_{j=1}^{\infty} a_j e^{-(1-\phi_0)\omega_j t} J_1(\sqrt{\omega_j}r),$$
(48)

where ω_j denote the possible solutions of Eq. (46) and where the values a_j , j = 1, ... are determined by the initial condition $w(r, t = 0) = w_{\text{init}}(r) \equiv u_{r,\text{init}} - u_r^{(0)}$ through

$$a_{j} = \int_{0}^{1} w_{\text{init}}(r) J_{1}(\sqrt{\omega_{j}}r) r \, \mathrm{d}r / \int_{0}^{1} J_{1}^{2}(\sqrt{\omega_{j}}r) r \, \mathrm{d}r.$$
(49)

For sufficiently small values of $\zeta \Delta \mu$, the system is always in the limit of small deformations and its dynamics is completely specified by the above solution, see Fig. 3a. Eventually, the relaxation occurs with a critical time scale $\omega_1^{-1} (1 - \phi_0)^{-1}$, where ω_1 is the smallest non-zero positive solution to Eq. (46).

4.4 Large deformations

In the general case of arbitrary radial deformations, the dynamic equation (41) is nonlinear and we have to resort to numerical solutions. For the numerical solution of the dynamic equations we use an explicit forward Euler method. The discretization scheme is detailed in App. A.

For small activities, the solution of the full equation agrees well with the linearized equation (43), see 3a. In the general case, initially the contraction is too slow, see Fig. 3b. The numerical solution shows that contraction starts always from the boundary, see Fig. 4. This was to be expected, because the active stress is unbalanced only there. The displacement and hence the stress in the gel then propagate towards the center of the disk according to the diffusion equation (41) with a state dependent diffusion constant. Together with the generation of a solvent flux, this is one of the hall marks of poroelasticity.



Fig. 3: Radius of a contracting poroelastic disk as a function of time for small activity $\zeta \Delta \mu = -0.1$ (a) and large activity $\zeta \Delta \mu = -1$ (b). Shown are the solution of the linearized dynamic equation (blue solid line) and of the full dynamic equation that is numerically integrated (black dashed line). The initial gel volume fraction is $\phi_0 = 0.01$.



Fig. 4: The gel volume fraction as a function of space and time for the same parameters as in *FIg. 3b.*

5 Conclusions

In this chapter, we have sketched the derivation of the hydrodynamic equations for isotropic active gels and discussed two examples of spontaneously emergent gel flows. These flows were due to contractile stresses generated by active processes that are driven by the hydrolysis of ATP. Since we restricted attention to isotropic gels without any polar order, the active stress only entered the description through the boundary conditions. In the general case, the active stress is anisotropic and its components change with changing polar (or nematic) order. In that case, much richer spontaneous flow patterns can be observed [2, 24]. They are tightly connected to the emergence of topological point defects. For high enough activity, one can even observe low Reynolds number turbulence, that is, spatiotemporal chaos.

Even though the generalized hydrodynamic equations for active gels are *a priori* only valid on macroscopic length scales, this approach has been successfully applied to describe subcellular dynamics, notably, the formation and contraction of actomyosin rings during cell division [17]. It will be interesting to see other applications in the future, in particular, through coupling active gel theory to cell signaling pathways. Since it is based on symmetries rather than molecular processes, active gel theory can and has been applied to tissue dynamics. In particular, in the context of embryonic development, a number of interesting phenomena wait to be analyzed in this framework.

Appendices

A Discretization scheme for the numerical solution of the active poroelastic dynamic equations

In this appendix, we explicit the discretization scheme, we used to solve the dynamic equation (41) for the radial displacement field. We introduce a dynamic lattice carrying the gel. Its sites are given by $i\Delta r + u_i$, i = 0, ..., N with $N = 1/\Delta r$. Here u_i is the radial displacement vector at $i\Delta r$. The gel volume fraction and the stresses are associated with the bonds between the lattice sites. We will, respectively, denote them by $\phi_{i-\frac{1}{2}}$, $\sigma_{rr,i-\frac{1}{2}}$, and $\sigma_{\theta\theta,i-\frac{1}{2}}$ for the bond between sites i and i - 1.

For calculating the gel volume fraction for a given displacement field, we do not take the approximate form given in Eq. (28), but the exact form that is obtained from equating the gel mass in a volume element before and after deformation. It yields for i = 2, ..., N

$$\phi_{i-\frac{1}{2}} = \phi_0 \frac{(2i-1)\,\Delta r^2}{\left[(2i-1)\,\Delta r + u_i + u_{i-1}\right]\left[\Delta r + u_i - u_{i-1}\right]} \tag{50}$$

and

$$\phi_{\frac{1}{2}} = \phi_0 \frac{\Delta r^2}{\left(\Delta r + u_1\right)^2}.$$
(51)

For the stresses, Eqs. (34) and (35) lead to

$$\sigma_{rr,i-\frac{1}{2}} = \left(\frac{u_i - u_{i-1}}{\Delta r}\right) + \frac{u_i + u_{i-1}}{(2i-1)\Delta r} - \zeta \Delta \mu$$
(52)

$$\sigma_{\theta\theta,i-\frac{1}{2}} = \left(\frac{u_i - u_{i-1}}{\Delta r}\right) + \frac{u_i + u_{i-1}}{(2i-1)\Delta r} - \zeta\Delta\mu$$
(53)

for $i = 2, \ldots, N$ and

$$\sigma_{rr,\frac{1}{2}} = \frac{u_1}{\Delta r} + \frac{u_1}{\Delta r} - \zeta \Delta \mu \tag{54}$$

$$\sigma_{\theta\theta,\frac{1}{2}} = \frac{u_1}{\Delta r} + \frac{u_1}{\Delta r} - \zeta \Delta \mu \tag{55}$$

$$\sigma_{rr,N+\frac{1}{2}} = \sigma_{\theta\theta,N+\frac{1}{2}} = 0,$$
(56)

where we have used the dimensionless form introduced in the main text. Then, the time evolution is given by

$$\dot{u}_{i} = \left(1 - \frac{\phi_{i-\frac{1}{2}} + \phi_{i+\frac{1}{2}}}{2}\right) \left\{ \frac{\sigma_{rr,i+\frac{1}{2}} - \sigma_{rr,i-\frac{1}{2}}}{\Delta r + \frac{u_{i+1} + u_{i-1}}{2}} + \frac{1}{2} \frac{\sigma_{rr,i+\frac{1}{2}} + \sigma_{rr,i-\frac{1}{2}} - \sigma_{\theta\theta,i+\frac{1}{2}} - \sigma_{\theta\theta,i-\frac{1}{2}}}{i\Delta r + u_{i}} \right\}$$
(57)

for i = 1, ..., N if we set $u_0 = u_{N+1} = 0$.

Finally, we use mass conservation, Eq. 39, to obtain the radial solvent velocity. The corresponding discretized field $v_{i-\frac{1}{2}}$ gives the solvent velocity at $(i-\frac{1}{2})\Delta r$ for $i = 1, \ldots, N$. It can be computed directly from $\phi_{i-\frac{1}{2}}$ and \dot{u}_i . Explicitly

$$v_{i-\frac{1}{2}} = -\frac{\phi_{i-\frac{1}{2}}}{1-\phi_{i-\frac{1}{2}}}\frac{\dot{u}_{i-1}+\dot{u}_i}{2}$$
(58)

for i = 1, ..., N. Outside of the gel, the solvent velocity vanishes, $v_{i-\frac{1}{2}} = 0$ for $i \ge N + 1$.

References

- [1] F. J. Nedelec, T. Surrey, A. C. Maggs, and S. Leibler, Nature 389, 305 (1997).
- [2] T. Sanchez, D. T. N. Chen, S. J. DeCamp, M. Heymann, and Z. Dogic, Nature 491, 431 (2012).
- [3] F. Nedelec and D. Foethke, New J. Phys. 9, 427 (2007).
- [4] J. M. Belmonte and F. Nedelec, Elife 5, 941 (2016).
- [5] A. E. Carlsson, Phys. Rev. Lett. 104, 228102 (2010).
- [6] M. Fritzsche, C. Erlenkämper, E. Moeendarbary, G. Charras, and K. Kruse, Sci. Adv. 2, e1501337 (2016).
- [7] K. Kruse and F. Jülicher, Phys. Rev. Lett. 85, 1778 (2000).
- [8] K. Kruse and F. Jülicher, Phys. Rev. E 67, 051913 (2003).
- [9] K. Doubrovinski and K. Kruse, Phys. Rev. Lett. 107, 258103 (2011).
- [10] V. Wollrab, R. Thiagarajan, A. Wald, K. Kruse, and D. Riveline, Nat. Comm. 7, 11860 (2016).
- [11] J. Prost, F. Jülicher, and J. F. Joanny, Nat. Phys. 11, 111 (2015).
- [12] S. R. de Groot and P. Mazur, *Non-Equilibrium Thermodynamics* (Dover Publications Inc, New York, 1985).
- [13] P. M. Chaikin and T. C. Lubensky, *Principles of condensed matter physics*, Cambridge University Press (2000).
- [14] K. Kruse, J.-F. Joanny, F. Jülicher, J. Prost, and K. Sekimoto, Phys. Rev. Lett. 92, 078101 (2004).
- [15] R. Voituriez, J. F. Joanny, and J. Prost, EPL (Europhys. Lett.) 70, 404 (2007).
- [16] K. Kruse, J. F. Joanny, F. Jülicher, and J. Prost, Phys. Biol. 3, 130 (2006).
- [17] H. Turlier, B. Audoly, J. Prost, and J.-F. Joanny, Biophys. J. 106, 114 (2014).
- [18] M. Behrndt, G. Salbreux, P. Campinho, R. Hauschild, F. Oswald, J. Roensch, S. W. Grill, and C.-P. Heisenberg, Science 338, 257 (2012).
- [19] S. Fürthauer, M. Neef, S. W. Grill, K. Kruse, and F. Jülicher, New J. Phys. 14, 023001 (2012).
- [20] J.-F. Joanny, F. Juelicher, K. Kruse, and J. Prost, New J. Phys. 9, 422 (2007).
- [21] A. C. Callan-Jones and F. Jülicher, New J. Phys., 093027 (2011).
- [22] J. F. Joanny, K. Kruse, J. Prost, and S. Ramaswamy, Eur. Phys. J. E 36, 52 (2013).
- [23] F. Jülicher, K. Kruse, J. Prost, and J. F. Joanny, Phys. Rep. 449, 3 (2007).
- [24] M. Neef and K. Kruse, Phys. Rev. E 90, 052703 (2014).
- [25] G. T. Charras, J. C. Yarrow, M. A. Horton, L. Mahadevan, and T. J. Mitchison, Nature 435, 365 (2005).
- [26] E. Moeendarbary, L. Valon, M. Fritzsche, A. R. Harris, D. A. Moulding, A. J. Thrasher, E. Stride, L. Mahadevan, and G. T. Charras, Nat. Mater. 12, 253 (2013).
- [27] I. Linsmeier, S. Banerjee, P. W. Oakes, W. Jung, T. Kim, and M. P. Murrell, Nat. Comm. 7, 12615 (2016).

C 7 Synapses

C. Karus, C. R. Rose Institute of Neurobiology Heinrich Heine University Düsseldorf

Contents

1	Intr	oduction	2	
2	Neu	rons and electrical activity	2	
	2.1	General structure of neurons	3	
	2.2	Neuronal membranes generate electrical signals	4	
3	Synapses: connecting and computing			
	3.1	The glutamatergic synapse	7	
	3.2	The GABAergic synapse	7	
5	Glia	al cells: an additional level of complexity		
3	Glia	al cells: an additional level of complexity		
	5.1	Microglia	12	
	5.2	Macrogna	12	
	5.5 5.4	Astrocytes and neurons form the tripartite synapse V^+ buffering	12	
	5.4	Clutemate untake and neurometabolic coupling	14	
	5.5	Giutamate uptake and neurometabolic coupling	13	
6	Con	clusion	16	
Ref	erence	°S		

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Nowadays, we live in an interconnected and very complex world. Only very rarely do we realize that the concept of being interconnected and constantly online has already been developed by nature millions of years ago. In our bodies, cells of our nervous system constantly receive and integrate information which is then sent to direct and remote neighbors or to other organs of the body. The main computation unit is the central nervous system (CNS), consisting of the brain and spinal cord. The peripheral nervous system (PNS) provides sensory cells as well as 'data highways', transferring sensory information to the CNS (afferent fibres) or motor and vegetative output from the CNS to the target organs (efferent fibres).

Brain function comprises various levels of information processing from perception of the environment via sensory systems over muscle control and spatial movement (motor functions) up to more complex tasks (higher functions) such as generation of circadian rhythms, conciousness, emotions, speech, learning and memory. All of these different processes are performed by the basic elements of the brain: neurons and glial cells. The human brain contains about 100 billion individual neurons and an equally high number of glial cells, but these cells do not act individually. Instead, they are functionally connected to each other by chemical and/or elecrical synapses. Each neuron can form thousands of synapses to other neurons, and also receives a large number of inputs from other cells. The structural and functional complexity of the neuronal network is thus extraordinary. Complexity is further increased by the involvement of different types of glial cells, which modulate and regulate synaptic properties. Moreover, this network is not static, but rather changes dynamically. Even in the mature brain, new synaptic connections are constantly formed or strengthened while others are depressed or even removed, a process which is the cellular basis for learning and memory.

This lecture provides an introduction into how neurons communicate with each other at synapses. Moreover, it presents some basic knowledge on the cellular mechanisms of learning and memory formation. Finally, we show that neurons are not the only players important for brain function. Rather, they are complemented and regulated by glial cells.

Neuronal function and basic synaptic mechanisms as presented in this lecture comply with current textbook knowledge. Therefore, the corresponding text is not referenced throughout, rather the reader is referred to standard textbooks such as [1, 2], from which the information presented was taken. This also largely holds true for information given on basic properties of glial cells. Here, readers may consult [3, 4]. Information that transcends basic textbook content is referrenced as it is discussed in the text.

2 Neurons and electrical activity

Evolution resulted in animal nervous systems of high complexity culminating in the brains of cetaceans, great apes and humans. Nevertheless, the neurons as basic functional units of highly evolved brains have a lot in common with those of simple organisms such as the nematode *Caenorhabditis elegans* or the leech *Hirudo medicinalis*. Neurons were first described and their morphology analyzed in more detail in the mid-late 19th century, when suitable labeling techniques such as the Golgi silver staining were developed. Wilhelm Gottfried von Waldeyer, a German anatomist, introduced the term "neuron" in 1881. Still, in the 19th century the supporters of the "reticularistic" view claimed that brain cells were fused together to form a continuum (syncytium), while their opponents postulated that each neuron forms an individual entity (cell theory). Finally, the latter opinion prevailed and the "neuron doctrine" was established, proposing that the cellular processes of neurons represent the

morphological substrates of signal processing and are not in direct continuation with neighboring cells. This concept, however, requires a mechanism that transfers signals from an individual neuron to the next. The existance of a hypothetical connective structure, the so-called *synapse*, was postulated by Charles Sherrington in 1897, but it took more than half a century until synaptic contacts could be first visualized by electron microscopy. Only recently, it became clear that there is also a second active network in the brain, represented by electrically coupled (and thus directly connected) networks of astrocytes, partly confirming basic ideas of the reticular theory.

2.1 General structure of neurons

As also discussed in chapter D3, neurons are morphologically polarized cells. This means that they extend distinct, specialized processes, called dendrites and axons, from their cell body (Fig. 1). Depending on the brain region and specific cell type, neurons posess from a few to up to hundreds of dendrites, each of which can branch several times. The entirety of dendrites is often referred to as dendritic tree. The dendrites bear the postsynaptic sites at which the neuron receives its input from other cells. Electrical signals generated at synapses spread anterogradly (and usually passively) along the dendrites towards the soma. Axons are the output structures of neurons and typically originate at or close to the cell body. If the excitatory synaptic input onto dendrites is strong enough and the depolarisation reaches the threshold for opening of fast voltage-gated sodium channels, axons generate action potentials at their initial segments. These propagate actively along the axon and its colaterals towards the presynaptic terminals, which oppose the postsynaptic sites of target neurons (see Fig.2). In sum, neurons receive electrical input, integrate the resulting signals and may then generate an output signal.



Fig. 1: General structure of neurons -A) Schematic textbook view of a neuron, showing the cell body from which numerous dendrites (light blue) as well as a single axon (beige) extend. The axon terminals contact other neurons at synapses. B) Fluorescence

image of a live pyramidal neuron in the mouse brain (hippocampus). The cell was dialyzed with a fluorescent dye employing whole-cell patch-clamp (see patch pipette attached to the soma). Note that this cell has two dendritic trees, one extending downwards, the other one upwards. The small-diameter axon is barely visible, running in parallel to the top border of the image.

2.2 Neuronal membranes generate electrical signals

The cellular mechanisms for the generation of resting membrane potentials and electrical signaling in neurons are described in detail in chapter D3. Here, we only provide a brief overview.

Cell membranes consist of a phospholipid bilayer separating the cytosol from the extracellular millieu. The bilayer is not permeable for large, polar or charged molecules. This enables cells to control the ionic composition of their cytoplasm, which differs considerably from that of the extracellular space (ECS). For example, the intracellular potassium ion (K⁺) concentration is about 100 mM, and the sodium ion (Na⁺) concentration is ~15 mM. In the ECS, the K⁺ concentration is much lower ([K⁺]_o~5 mM) while [Na⁺]_o is high (~150 mM). These differences in intracellular and extracellular concentrations are generated and maintained by the continuous activity of a plasma membrane ion pump, the Na⁺-K⁺-ATPase (NKA), which exports 3 Na⁺ in exchange for 2 K⁺.

Transmembrane ion channels generate a water-filled passage across the plama membrane and thereby open a possible route for ion permeation in between the two compartments. In neurons (and glial cells), mainly K⁺-selective ion channels are open at rest, enabling the transmembrane diffusive movement of K⁺. Based on the chemical gradient for K⁺ set by the NKA, there is an initial net efflux of K⁺ through these channels, leaving unbalanced negative charges behind until the chemical driving force is balanced by the generated negative polarization of the membrane. The electrical potential, at which the chemical and electrical forces acting upon a given ion are in balance and no net current flows, is called the equilibrium potential (E_{Ion}), alternative names are "reversal" or "Nernst" potential. It can be calculated by the Nernst equation:

$$E_{Ion} = \frac{R * T}{Z * F} * \ln \frac{[Ion]_o}{[Ion]_i}$$

For K^+ , the equilibrium potential E_K equals about -80 mV (at 37°C). This value roughly represents the resting membrane potential of many neurons (-60 to -70 mV) and macroglial cells (-80 to -90 mV).

Active neurons undergo rapid changes in the ion permeability of their membrane following opening or closing of ion channels, resulting in their de- or hyperpolarization. In general, when describing voltage changes in neurons, the term "excitation" describes a depolarizing shift of the membrane potential that makes action potential generation more likely. In contrast, "inhibition" is defined as a process that makes action potential generation less likely, which often is induced by a hyperpolarization of the membrane. "Subthreshold" signals do not reach the threshold for induction of action potentials, "suprathreshold" signals result in the action potentials firing.

There are two basic types of transient electrical signals in central neurons: small changes (a few mV) in membrane voltage generated locally at synapses, which spread passively (electrotonically) along the dendritic tree towards the soma. Action potentials, in contrast, represent large (~80-90 mV) and rapid changes in membrane voltage, arising at specialized

initial segments of axons and spreading actively along the latter. The generation of action potentials is initiated by a voltage-dependent, massive, transient opening of Na⁺ channels. As a result, the membrane permeability for Na⁺ is increased, surpassing the K⁺ permeability. Because E_{Na} is about +60 mV, Na⁺ flows into the cells and depolarizes the membrane. Na⁺ flux rapidly ceases due to voltage-dependent channel inactivation and delayed efflux of K⁺ then repolarizes the membrane towards the previous resting potential.

Taken together, electrical signalling of neurons is based on defined, regulated changes in the permeability of their cell membrane to specific ions. This results in ion currents, which may change the membrane potential. The electrical signals propagate from their site of origin throughout the cell; signals generated at synapses spread anterogradly along the dendrites towards the soma and action potentials generated at axon initial segements propagate along the axon towards presynaptic terminals.

3 Synapses: connecting and computing

Neurons are functionally connected by synapses which provide the machinery for the controlled and rapid transmission of electrical signals from cell to cell. There are two basic types of synapses in the brain: electrical and chemical synapses (Fig.2).

Electrical synapses enable a direct charge transfer in between cells. They are far less common in neurons than chemical synapses, but represent the basis for the formation of astrocyte networks. At an electrical synapse, the membranes of the connected cells are in very close vicinity and form so-called gap junctions, membrane areas which contain proteins called connexons (Fig. 2A). Connexons form ion channels (hemichannels) and each of these hemichannels directly opposes a connexon in the neighboring membrane. This results in the formation of a continuous pore in between the two cells (called gap junction channel), allowing the direct transfer of ions (and thus, current) and other small molecules. Electrical synapses transmit electrical signals with essentially no delay, but – because signals spread electrotonically - with a reduction in rise time and amplitude. Cells connected via gap junctions are often referred to as electrically coupled. Most electrical synapses function bidirectionally and e.g. serve in synchronization of discharge patterns of neuronal networks. In astrocytes, gap junctions not only mediate the transfer of ions and charge but also enable the distribution of metabolites such as glucose in the network to high-activity regions.

Chemical synapses represent the major, canonic route of signal transmission between neurons of the brain. While glial cells are not known to be coupled to each other by chemical synapses, they can still perceive and react to signals released from active synapses (cf. Fig.5). Moreover, certain subtypes of glial cells have even been described to form quasi "postsynaptic" elements of chemical synapses with neurons. Finally, glial cells (astrocytes) can release chemical messengers that are sensed by neurons (so-called "gliotransmitters"; see Fig.5), adding to the complexity and diversity of chemical synapses.

Classical chemical synapses work unidirectionally (although there is some feedback communication) and consist of distinct, highly specialized pre- and postsynaptic elements, formed by the presynaptic and postsynaptic neurons, respectively (Fig.2B). Chemical synapses are structurally and functionally separate entities with their opposing membranes about 20 nm apart. The presynapse is formed by the terminals of axons (also called presynaptic terminal). It receives electrical signals (action potentials) via the axon and transforms them into chemical signals. This done by a voltage-dependent release of a

chemical messenger, called neurotransmitter. A given synapse usually uses one specific (main) neurotransmitter "x" and is therefore called "x"-ergic.



Fig. 2: Synapses - A) Electrical synapses are formed by gap junctions. Gap junctions are membrane areas connected via gap junction channels. Each gap junction channel consists of two apposed connexon hemichannels, thereby forming a continuous pore between the interior of coupled cells. Gap junctions enable the bidirectional passage of ions (currents) and small molecules. B) Chemical synapses pass information unidirectionally. They consist of a presynaptic terminal (at the end of an axon) and postsynaptic site, separated by a synaptic cleft. The presynapse contains neurotransmitter vesicles, while the postsynapse exposes specific neurotransmitter receptors towards the synaptic cleft.

The postsynaptic elements, typically located along the dendritic tree and cell body of the postsynaptic neuron, receive those chemical signals and transform them back to electrical signals again. This is realized by special receptor proteins in the plasma membrane, which bind a given neurotransmitter with high affinity. Neurotransmitter receptors are either ion channels themselves, and their activation can therefore directly induces ion currents and changes in membrane voltage of the postsynaptic cell, these are called "ionotropic receptors". Alternatively, neurotransmitter receptors are coupled to intracellular signaling cascades, resulting in the generation of an intracellular second messenger, these are called "metabotropic receptors". Activation of metabotropic receptors may cause a secondary change in ion channel conductance and thereby generate a delayed voltage signal in the postsynaptic cell.

The basic machinery and functional principles of chemical synaptic transmission are illustrated in the following chapters, using an excitatory (glutamatergic) and and inhibitory (GABAergic) synapse as examples (cf. Fig. 3). However, it is important to emphazise that there are many more transmitter substances than glutamate and GABA, not all of them strictly exitatory or inhibitory. Next to classical neurotransmitters, neuromodulators add to the complexity. All of these mechanisms are prerequisits for the many specialized brain functions and to enable the adjustment of synaptic properties to current demands, a process called "synaptic plasticity". Such plasticity is required to enable intricate cognitive processes such as learning and memory.

3.1 The glutamatergic synapse

The amino acid glutamate is the most common excitatory neurotransmitter in the CNS. In the presynaptic terminal, glutamate is synthetised from glutamine by the enzyme glutaminase and transported into intracellular vesicles (small compartments enclosed by membrane) via the vesicular glutamate transporter "vGLUT" (Fig. 3). These neurotransmitter vesicles store the glutamate at high concentration until an action potential arrives at the terminal via the axon. A fraction of the vesicles is already attached to the (inner) presynaptic cell membrane opposing the synaptic cleft by a specialized protein complex, called "SNARE complex" and ready to fuse (fusion pool). Invasion of the presynaptic terminal by an action potential causes the opening of voltage-dependent calcium channels (Ca_v), resulting in the influx of calcium. The SNARE complex contains the calcium-sensitive protein synaptotagmin, which then changes its conformation causing the neurotransmitter vesicles to fuse with the cell membrane and to release glutamate into the synaptic cleft.

Excitatory synapses are usually located on so-called dendritic spines (see highlighted synapse in Fig.1A and Fig.2A), forming small postsynaptic sub-compartments at dendrites. Glutamate released from the presynaptic terminal diffuses towards the postsynaptic membrane which displays glutamate receptors orientated towards the extracellular space. Binding of glutamate to ionotropic receptors results in the direct opening of an cation-selective ion pore in the receptor. Three types of ionotropic glutamate receptors are known and named after agonist substances: AMPA receptors (α-amino-3-hydroxy-5-methyl-4-isoxazolepropionic acid receptor), NMDA receptors (N-methyl-D-aspartic acid receptor), and kainate receptors. Ionotropic glutamate receptor are mainly permeable for sodium and potassium, resulting in a depolarisation of the postsynaptic cell. This depolarization shift is called EPSP (excitatory postsynaptic potential; term independent of causing neurotransmitter), the underlying ion current is called EPSC (exitatory postsynaptic current). In the CNS, single EPSPs are usually too weak to overcome the action potential firing threshold, but synaptic potentials sum up spatially and temporally so that action potentials can be are generated upon repeated synaptic activity.

NMDA receptors, unlike the majority of AMPA receptors, are additionally permeable to calcium and their activation therefore results in the generation of an intracellular calcium signal. Moreover, and again in contrast to AMPA receptors, NMDA receptors require more than glutamate binding to open. At negative resting membrane potentials, their ion pore is blocked by magnesium ions, preventing the generation of ion current even with glutamate bound. This block is relieved upon a depolarisation of the postsynaptic membrane, as e. g. achieved by additional strong activation of AMPA receptors. NMDA receptors, therefore, represent coincidence detectors, generating electrical signals and calcium influx only with both pre-synaptic activity (resulting in release of glutamate) and with post-synaptic activity (activation of AMPA receptors resulting in a strong postsynaptic depolarisation).

As a final step in glutamatergic transmission, the action of glutamate is terminated by its uptake into neighboring astrocytes (cf. Fig.5). This essential step is described in more detail in chapter 4.5.

3.2 The GABAergic synapse

In neuronal networks, excitation needs to be balanced and controlled by opposing inhibitory signals. The most common inhibitory neurotransmitter in the brain is GABA (gamma-aminobutyric acid), which is synthetised in the presynaptic terminals (Fig.3B) from glutamate by decarboxylation through the enzyme glutamate decarboxylase (GAD). GABA is then



stored in vesicles and released upon depolarisation of the presynaptic terminal by action potentials.

glutamatergic synapse

GABAergic synapse

Fig. 3: Exemplary chemical synapses - A) Excitatory glutamatergic synapse - glutamate is synthetized in the presynaptic terminal and packed into vesicles via vesicular glutamate transporters (vGLUT). An action potential depolarizes the presynaptic membrane, causing the opening of voltage-gated calcium channels (Ca_v). The resulting calcium increase triggers the fusion of vesicles and the release of glutamate into the synaptic cleft. The postsynaptic element contains ionotropic receptors for glutamate (AMPA and NMDA receptors). AMPA receptors open upon glutamate binding and allow passage of Na⁺ and K⁺, resulting in a postsynaptic depolarisation (EPSP). NMDA receptors, which require an additional depolarisation to open their ion pore, are permeable to Na⁺, K⁺ and Ca²⁺. Glutamate binding causes generation of an EPSP.B) Inhibitory GABAergic synapse - GABA is synthetized from glutamate and transported into vesicles via vesicular GABA transporters (vGAT). Ionotropic GABA_A receptors at the postsynaptic membrane are permeable to CI⁻. Their opening results in an efflux of CI⁻ and a postsynaptic hyperpolarisation (IPSP).

The postsynaptic membrane of GABAergic synapses is equipped with GABA receptors. Ionotropic GABA_A receptors are Cl⁻ permeable ion channels. In mature neurons, the Nernst potential for Cl⁻ is slightly more negative than the resting membrane potential. Opening of GABA_A receptors therefore leads to the efflux of Cl⁻ and stabilizes the membrane potential at a negative value. The resulting change in membrane voltage is called an IPSP (inhibitory postsynaptic potential), the underlying current called IPSC (inhibitory postsynaptic current). Like EPSPs, IPSPs are spatially and temporally summed up in the postsynaptic neuron. They drive the membrane to voltages negative to the action potential threshold and thereby dampen neuronal excitability.

4 LTP – a change in synaptic strength as cellular basis for learning and memory

The mature brain only has a very limited capacity for the formation of new neurons. Learning and memory formation is thus largely based on changing the properties of existing neuronal networks. Such changes are brought about by altering presynaptic or postsynaptic properties of a given synapse. This can be achieved by changing presynaptic transmitter release or the number and type of postsynaptic transmitter receptors and/or by adapting the number of synapses in a network.

In 1949, Donald Hebb proposed that changes in neuronal networks are induced by and depend on the concommittant activity of neurons functionally connected to each other. In brief, his ideas are often summarized in the sentence: "*neurons that fire together, wire together*". Nowadays, this principle is widely accepted and processes based on it known as "Hebbian learning".

A brain region, in which neuronal networks exhibit major characteristics of "Hebbian learning", is the hippocampus. The hippocampus is an evolutionarily old structure of the cortex which plays a central role in spatial learning. Its neuronal ciruitry is well characterized and basically consists of a trisynaptic excitatory network. The third synaptic connection within this network consists of presynaptic terminals of glutamatergic axons, called "Schaffer collaterals" and postsynaptic neurons in the "CA1" area. These specific synapses show an activity-dependent change in their properties called "long-term potentiation" (LTP). LTP describes a persistant increase in the amplitude of EPSPs, induced by defined patterns of pre-(and post) synaptic activity of the involved neurons.

To experimentally analyze the properties of synaptic connections between Schaffer collaterals and CA1 neurons, the Schaffer collaterals (afferent axons) are electrically stimulated to fire an action potential, causing them to release glutamate. This glutamate activates postsynaptic ionotropic glutamate receptors of the AMPA type on CA1 neurons, resulting in an EPSP, which can be recorded using electrophysiological techniques. If the stimulation of afferent axons is done at a low frequency, e. g. every 10 seconds (0.1 Hz), the amplitude of the induced EPSPs is rather stable over time, implying that the synaptic connection is stable and reliable.

If axons are stimulated at 100 Hz for 1 or 2 seconds, the resulting EPSPs sum up and reach the threshold for induction of action potentials in the postsynaptic cell. Notably, however, even after termination of this high-frequency activation, EPSP amplitudes in response to single stimuli at 0.1 Hz stay significantly increased as compared to the situation before the strong stimulus. This long-lasting, activity-dependent potentiation of the synaptic response is called "LTP". It is maintained for many hours, if not days; the synapse thus "remembers" its strong use.

After its original description in the 1970'ies, LTP has has long been regarded as an interesting experimental phenomenon that neuronal networks can undergo under artificial stimulation conditions. Only about 40 years later, it was shown that this phenomenon is induced by and the basis of spatial learning in behaving animals. Nowadays, it is the most widely accepted model for cellular learning and memory also for the intact brain.

A wealth of studies has paved the way for our understanding of the molecular mechanisms of LTP induction. Under conditions of regular, relatively sparse activity, glutamate release from a presynaptic terminal mainly activates AMPA receptors, inducing subthreshold EPSPs and postsynaptic depolarisations that are not strong enough to relieve the magnesium block of

NMDA receptors. If presynaptic axons fire at high frequencies, the resulting EPSPs sum up, eventually enabling activation of NMDA in addition to AMPA receptors. This additional current not only increases the EPSP amplitude and duration further. It also results in an NMDA receptor-mediated influx of calcium into the postsynaptic cell. Notably, this happens only, if presynaptic cells release glutamate and postsynaptic cells are strongly depolarized, emphazising the role of NMDA receptors as coincidence detectors as mentioned above.

NMDA receptor-mediated postsynaptic calcium signals trigger changes which strenghten the future postsynaptic response, including an increase in AMPA receptor conductance and the insertion of additional AMPA receptors at existing synapses. Finally, this can even lead to the growth of new synapses, driven by Ca^{2+} dependent alteration of gene experession. This so-called "late" LTP involves more permanent structural changes such as the enlargement of existing synapses and formation of new synaptic contact sites.

5 Glial cells: an additional level of complexity

Brain function not only depends on neurons, but also requires glial cells. Glial cells are represented by macroglia, which develop from ectodermally-derived neural cell lineages and are closely related to neurons. Microglia, the other big class of glial cells, are of mesodermal origin, and represent a cell type not directly related to neurons nor macroglia (Fig. 4).

Although glial cells were already discovered and described in the 19th century (first mentioned under the name neuroglia by Rudolph Virchow in 1858), they have been largely neglected and regarded as mere connective tissue for more than a century (glia is greek for glue). A major reason for this was that - in contrast to neurons - they generally do not undergo strong electrical signalling. Therefore, they were regarded as essentially non-responsive cells during the initial periods of electrophysiology-dominated neurophysiological studies. Throughout the last decades, intensified research efforts on glial cells and new imaging methods have, however, demonstrated that glia undergo vivid signalling, which is mainly based on chemical second messengers such as calcium. These studies have established that glia are vital for brain function, both for homeostatic control and metabolic supply of neurons as well as for synapse formation and information processing in neural networks. For basic information and as references to the following paragraphs the reader may consult textbooks like [3, 5].

Fig. 4: A) Schematic overview of different types of glial cells showing microglial cells, oligodendrocytes and astrocytes in addition to neurons. B) Immunostaining of astrocytes in the mouse hippocampus. The cytoskeleton protein GFAP (glial fibrillary acidic protein) is labelled in yellow and cell nuclei in blue (image: L. Müller-Thomsen, Institute of Neurobiology, HHU Düsseldorf). C) Single astrocyte in the mouse hippocampus dialysed with fluorescent dye (image: C.R. Rose, Institute of Neurobiology, HHU Düsseldorf). D) Immunostaining of MBP (myelin basic protein), labeling oligodendrocytes in cell culture. Oligodendrocytes were grown in culture dishes with artificial scaffolding to mimic the presence of axons (see[6]). Astrocytes are labelled via GFAP staining, nuclei are counterstained with (DAPI: 4,6-diamidino-2-phenylindole) (image: J. Jadasz, Department of Neurology, HHU Düsseldorf). E) Microglial cells in mouse hippocampal slice, immunostained for the specific marker protein Iba1 (ionized calcium-binding adapter molecule 1). In addition, the astrocytic glutamate transporter protein GLAST is labelled (image: A.E. Schreiner, Institute of Neurobiology, HHU Düsseldorf).











5.1 Microglia

Microglial cells (Fig. 4E) originate from the mesoderm (the germinal layer producing e.g. bones, muscles and blood cells) and act as resident phagocytes of the brain. During development of the brain, microglia help to adjust the number of neurons and synapses by a so-called pruning of surplus connections and extranumerary cells. Later, in the adult brain, microglia continuously monitor their environment, including synaptic connections, with their highly motile processes. They seem to directly participate in the activity-dependent adjustment of the synapses and also release signal molecules that have significant impact on neuronal communication.

Under pathological conditions, injury or disease, microglia become activated and direct their processes or even migrate completely towards the site of injury. There, they phagocytose and thereby remove parts of dead cells and extracellular debris. Moreover, they appear to form a protective barrier that isolates diseased from healthy tissue, a property also known for astrocytes (see below).

5.2 Macroglia

The group of macroglial cells, which derives from the ectodermal germinal layer, is further subdived into oligodendrocytes, NG-2 cells, astrocytes and ependymal cells. The latter group includes the cells that line the ventricular walls of the brain and the central canal of the spinal cord. There, they are critically involved in the production and composition control of cerebrospinal fluid. NG-2 cells are precursor cells for oligodendrocytes. Their functions in the healthy adult brain are not completely understood.

Oligodendrocytes and astrocytes (Fig. 4D) are present throughout the CNS. Oligodendrocytes are a major component of white matter tracts. They form large, densely packed membrane sheats (myelin) that are tightly wrapped around axons increasing their membrane resistance. Myelination dramatically speeds up the velocity of action potential conduction, which is required to provide fast information transfer along axons. Degeneration of the myelin sheath - as it occurs in degenerative diseases such as multiple sclerosis and leukodystrophies - is a detrimental state for the CNS. While myelination can sometimes be partly restored, a permanent loss of myelin often results in the degeneration of axons and a complete loss of function.

Astrocytes are a essential for ion and transmitter homeostasis at synapses (see also Fig.5). Moreover, they are part of the blood-brain barrier and involved in neurometabolic coupling, that is in the adaptation of brain glucose uptake and metabolism to neuronal needs. Finally, they are more and more recognized as being an active part of brain information processing. Because of their manyfold functions at chemical synapses, astrocytes will be discussed in further detail below.

5.3 Astrocytes and neurons form the tripartite synapse

Astrocytes occur in distinct morphologies and display different functional properties depending on the developmental stage, brain region and health status. Radial glia cells are precursors of astrocytes which are mainly present throughout brain development. These bipolar cells also serve as neuronal progenitor cells as well as migration scaffolding e. g. for development of cortical layers. In adulthood, retinal Müller glia and cerebellar Bergmann glia retain radial glia-like morphologies.

In the white matter tracts, astrocytes exhibit a fibrous appearance, bearing several long processes with which they contact blood vessels or nodes of Ranvier along axons. Protoplasmic astrocytes of grey matter areas are highly ramified cells with a multitude of very

fine processes in contact with blood vessels or synapses (as shown in Fig.4). Blood vessels are in fact nearly completely enwrapped by astrocyte processes called perivascular endfeet. Fine astrocyte processes contacting neuronal synapses, called perisynaptic processes, have become subjects of great interest within the last two decades. They are equipped with the molecular machinery to control the extracellular millieu in the perisynaptic area. A task fulfilled by astrocytes and one of the first astrocyte properties known, is their uptake of K^+ released by active neurons (see Fig.5). Moreover, astrocytes play a vital role in the re-uptake and removal of neurotransmitters, a process especially critical for glutamate (cfp. Figs.3A and 5).

In addition, astrocytes can sense and respond to synaptic activity. This is because they also express receptors for neurotransmitters. Once activated, these can initiate intracellular signalling cascades e. g. involving the generation of calcium signals. Activity-related glial calcium signalling can induce the release of so-called gliotransmitters, among them glutamate, D-serine or ATP that act back on neighboring neurons. To account for this dynamic interplay at synapses, involving information transfer and communication between three individual cells (presynaptic, postsynaptic neuron and astrocyte), the term "tripartite synapse" was established. It implies that astrocytes are an important, active signalling components of synapses.

In contrast to the majority of neurons, astrocytes are generally coupled to each other by gap junctions and thereby form a syncytium of interconnected cells. This coupling enables the direct transfer of ions (e.g. Na⁺), metabolites and signalling molecules like calcium between cells. Astrocytes thus form a second, defined cellular network in the brain, that functions separately from that of synaptically-connected neurons. Gap junctional coupling has been shown to mediate the transfer of metabolites to active brain regions. Moreover, it promotes the removal of K^+ from the extracellular space into astrocytes by shuttling it to more remote regions, a process called "spatial buffering" (see below).



Fig. 5: Tripartite glutamatergic synapse - Astrocytes take up glutamate from the extracellular space and replenish neuronal glutamine pools. In addition, astrocytes express transmitter receptors (TM receptors) which initiate intracellular signalling cascades and eventually may cause the release of gliotransmitters. Moreover, astrocytes match local activity levels and energy demand to their metabolism (neurometabolic coupling). Finally, astrocytes take up potassium released by neurons and, thereby, regulate network excitability.

5.4 K⁺ buffering

As pointed out above (cp. 2.2), cell membranes (including neuronal and glial membranes at rest) are mainly permeable for K⁺. Accordingly, their resting membrane potential is close to the equilibrium potential K⁺, which can be calculated by the Nernst equation. With a typical K⁺ distribution of 100 mM K⁺ inside and 5 mM K⁺ outside cells, E_K is ~-80 mV (at 37°C). The dependence of neuronal (and even more so glial) membrane potentials on the E_K also implies that changes in K⁺ directly result in changes in the former. Because this is so, extracellular K⁺ concentrations have to be controlled tightly.

Notwithstanding, neurons release substantial amounts of K^+ during action potentials through voltage-gated delayed rectifier channels. Moreover, at glutamatergic synapses, K^+ is released following opening of (K^+ -permeable) ionotropic glutamate receptors (AMPA and/or NMDA receptors) at the postsynaptic site. Because the extracellular space only provides a very limited volume, these ion fluxes cause significant changes in the extracellular K^+ concentration [7]. Experiments using K^+ sensitive microelectrodes in vertebrate brain,
however, revealed that extracellular K^+ only increases to a maximum level of 10-12 mM under physiological conditions, independent from the strength and duration of activity. This so-called K^+ ceiling level [8, 9] is general property of heathly brain tissue, and only overcome under pathological conditions, e. g. during epileptic seizures. Experimentally increasing $[K^+]_o$ was indeed found to result in uncontrolled discharge activity of the neuronal network [10] and is frequently employed as a model for induction of epileptiform discharges *in situ* and *in vivo*.

Under physiological conditions, control of extracellular K^+ concentration and maintanance of the ceiling level are mainly mediated by net uptake and active clearance of K^+ by astrocytes [11] (Fig.5). Two mechanisms have been proposed to contribute to astrocyte K^+ clearance. The first is spatial buffering and limited to situations in which the K^+ increase occurs locally. This enables uptake of K^+ through inward-rectifier K^+ channels (mainly $K_{ir}4.1$), which then diffuses intra- and intercellularly to inactive regions, where it is relased again. Spatial buffering has been demonstrated and elegantly decribed in the retina (there also called "spatial siphoning"), where Müller glial cells take up K^+ in the inner plexiform layer and then release it into the vitreous humor [12].

Astrocytes also take up K^+ by the NKA. In addition to the ubiquitously expressed α 1-subunit of the NKA, which is already saturated with K^+ at resting levels [13], astrocytes express the NKA α 2-subunit [14]. These have a lower affinity for external K^+ and therefore increase pumping activity with increases in $[K^+]_o$. Moreover, astrocytes express NKCC1 (Na⁺-K⁺-Cl⁻ cotransporter 1) which can contribute to net uptake of K^+ [15].

Net uptake of ions changes the osmolarity of the cytoplasm and causes water influx resulting in cell swelling. Such swelling of the cells and accordingly a decrease in extracellular space volume has been observed in *in vivo* experiments [16]. Volume changes can massively impact neuronal excitability and are counteracted by homeostatic mechanisms, namely regulatory volume decrease [17].

5.5 Glutamate uptake and neurometabolic coupling

As a component of the tripartite synapse (Fig.5), astrocytic perisynaptic processes are in close spatial association with the pre- and postsynaptic elements. This proximity enables astrocytes to bind and trap neurotransmitters released by active neurons. At glutamatergic synapses, astrocytes are responsible for the fast and efficient removal of the transmitter glutamate from the ECS (Fig. 5). Thereby astrocytes shape excitatory synaptic transmission, limiting postsynaptic receptor activation and diffusion of glutamate to adjacent synapses (a phenomenon termed "spill-over") [18]. Moreover, the removal of glutamate is the key mechanism preventing glutamate-induced overexcitation and resulting cell damage [19].

Astrocyte glutamate uptake maintains the extracellular glutamate concentrations at very low levels and ensures rapid recovery to these levels after its synaptic release [20]. As diffusion is a nondirectional process, glutamate not only diffuses towards the postsynaptic receptors but also towards the borders of the synaptic site. Here, the membranes of astrocytic perisynaptic processes are enriched with high-affinity glutamate transporter proteins [19, 21]. Five subtypes of glutamate transporters have been described. In the human brain, they are termed excitatory amino acid transporters (EAATs) 1-5. Astrocytes typically express EAAT1 and 2, named GLAST and GLT-1, respectively in the murine brain [19].

All of these glutamate transporters employ the same transport mechanism. To import glutamate against its electrochemical gradient, the transporters exploit the energy stored in the transmembrane gradients for Na⁺ and K⁺. To import one glutamate molecule, three Na⁺ ions enter and one K⁺ ion leaves the cell. In addition, one proton enters the cytoplasm. Accordingly, the transport generates a net inward current, as four positive and one negative

charge enter while only one positve charge leaves the cell (net influx of two positive charges). This uptake current can be recorded with electrophysiological techniques, providing a semiquantitative measure of glutamate transport activity of astrocytes. There is evidence that glutamate transporter activity can affect the functionality of other electrogenic transporters such as GABA transporters in astrocytes (for example increased inhibition by reversal of GABA transport direction [22, 23]), thereby functioning as feedback and modulating synaptic activity on a broader scale.

Based on its transport properties, glutamate uptake into astrocytes transiently increases the astrocytic Na⁺ concentration which in turn stimulates the activity of the NKA [24]. It has been calculated that the import of a single glutamate molecule from the synaptic cleft requires 1.5 molecules of ATP to restore transmembrane ion gradients. This stimulates astrocyte metabolism to replenish ATP levels by breakdown of glycogen and/or increased uptake of glucose from the blood. Cellular glycoloysis and production of lactate is increased. This lactate is released by astrocytes and taken up as metabolic substrate by neurons, a phenomenon called astrocyte-neuron lactate shuttle [25]. It has been proposed that this represents a basic principle of neurometabolic coupling (Fig.5), that is the rapid adaptation of cellular metabolism to the activity level of neurons [24, 26].

Hence, synaptic functionality is tightly connected to astrocytic regulation of glutamate levels and energy supply. This further highlights that glutamatergic synapses are tripartite systems including neuronal and glial components as equally important players.

6 Conclusion

Synapses are the functional element of information-transfer from one neuron to another. The major type of synapses are chemical synapses which transform an electrical input signal into a chemical signal (presynaptic neurotransmitter release) and back to an electrical one (postsynaptic EPSP/IPSP generation). The most important neurotransmitters in the CNS are glutamate (excitatory) and GABA (inhibitory).

Chemical synaptic transmission inherits various possibilities for modifications and adaptations and the chemical synapse is the plastic and dynamically changing element of the learning brain.

Contrary to standard textbook presentations, recent work has provided evidence that glial cells, especially astrocytes, are more than passive bystanders of synaptic activity doing the housekeeping. Rather, astrocytes are now viewed as a important functional parts of the tripartite synapse. E.g. astrocytes influence cellular excitability by K⁺ buffering, they shape synaptic transmission via neurotransmitter uptake and release of gliotransmitters. Furthermore, astrocytes match activity levels of neurons to their metabolic supply.

References

[1] P.M. Bear M. F., Connors B.W., Neuroscience - Exploring the Brain, 4 ed., Wolters Kluwer2015.

[2] A.G.J. Purves D., Fitzpatrick D., Hall W. C., LaMantia A.-S., White L. E., Neuroscience, 5 ed., Sinauer2012.

[3] A. Verhratsky, A. Butt, Glial Neurobiology, John Wiley and Sons Ltd2007.

[4] H. Kettenmann, Ransom, B., Neuroglia, Oxford University Press2012.

[5] B.A. Barres, M.R. Freeman, B. Stevens, Glia, Cold Spring Harbor Laboratory Press2015.

[6] J.J. Jadasz, L. Tepe, F. Beyer, I. Samper Agrelo, R. Akkermann, L.S. Spitzhorn, M.E. Silva, R.O.C. Oreffo, H.P. Hartung, A. Prigione, F.J. Rivera, J. Adjaye, P. Kury, Human mesenchymal factors induce rat hippocampal- and human neural stem cell dependent oligodendrogenesis, Glia 66(1) (2018) 145-160.

[7] I. Dietzel, U. Heinemann, G. Hofmeier, H.D. Lux, Stimulus-induced changes in extracellular Na+ and Cl- concentration in relation to changes in the size of the extracellular space, Exp Brain Res 46(1) (1982) 73-84.

[8] U. Heinemann, H.D. Lux, Ceiling of stimulus induced rises in extracellular potassium concentration in the cerebral cortex of cat, Brain Res 120(2) (1977) 231-49.

[9] G.G. Somjen, Ion regulation in the brain: implications for pathophysiology, Neuroscientist 8(3) (2002) 254-67.

[10] S.F. Traynelis, R. Dingledine, Potassium-induced spontaneous electrographic seizures in the rat hippocampal slice, J Neurophysiol 59(1) (1988) 259-76.

[11] W. Walz, Role of astrocytes in the clearance of excess extracellular potassium, Neurochem Int 36(4-5) (2000) 291-300.

[12] E.A. Newman, D.A. Frambach, L.L. Odette, Control of extracellular potassium levels by retinal glial cell K+ siphoning, Science 225(4667) (1984) 1174-5.

[13] G. Blanco, R.W. Mercer, Isozymes of the Na-K-ATPase: heterogeneity in structure, diversity in function, Am J Physiol 275(5 Pt 2) (1998) F633-50.

[14] K.J. Sweadner, Overlapping and diverse distribution of Na-K ATPase isozymes in neurons and glia, Can J Physiol Pharmacol 70 Suppl (1992) S255-9.

[15] L. Hertz, J. Xu, D. Song, E. Yan, L. Gu, L. Peng, Astrocytic and neuronal accumulation of elevated extracellular K(+) with a 2/3 K(+)/Na(+) flux ratio-consequences for energy metabolism, osmolarity and higher brain function, Front Comput Neurosci 7 (2013) 114.

[16] I. Dietzel, U. Heinemann, G. Hofmeier, H.D. Lux, Transient changes in the size of the extracellular space in the sensorimotor cortex of cats in relation to stimulus-induced changes in potassium concentration, Exp Brain Res 40(4) (1980) 432-9.

[17] K.T. Kahle, J.M. Simard, K.J. Staley, B.V. Nahed, P.S. Jones, D. Sun, Molecular mechanisms of ischemic cerebral edema: role of electroneutral ion transport, Physiology (Bethesda) 24 (2009) 257-65.

[18] A.V. Tzingounis, J.I. Wadiche, Glutamate transporters: confining runaway excitation by shaping synaptic transmission, Nature reviews. Neuroscience 8(12) (2007) 935-47.

[19] N.C. Danbolt, Glutamate uptake, Prog Neurobiol 65(1) (2001) 1-105.

[20] Y. Zhou, N.C. Danbolt, GABA and Glutamate Transporters in Brain, Front Endocrinol (Lausanne) 4 (2013) 165.

[21] N. Zerangue, M.P. Kavanaugh, Flux coupling in a neuronal glutamate transporter, Nature 383(6601) (1996) 634-7.

[22] L. Heja, G. Nyitrai, O. Kekesi, A. Dobolyi, P. Szabo, R. Fiath, I. Ulbert, B. Pal-Szenthe, M. Palkovits, J. Kardos, Astrocytes convert network excitation to tonic inhibition of neurons, BMC Biol 10 (2012) 26.

[23] P. Unichenko, A. Dvorzhak, S. Kirischuk, Transporter-mediated replacement of extracellular glutamate for GABA in the developing murine neocortex, Eur J Neurosci 38(11) (2013) 3580-8.

[24] C.R. Rose, J.Y. Chatton, Astrocyte sodium signaling and neuro-metabolic coupling in the brain, Neuroscience 323 (2016) 121-34.

[25] L. Pellerin, P.J. Magistretti, Sweet sixteen for ANLS, J Cereb Blood Flow Metab 32(7) (2012) 1152-66.

[26] P.J. Magistretti, Role of glutamate in neuron-glia metabolic coupling, Am J Clin Nutr 90(3) (2009) 875S-880S.

D 1 Physics and modelling of intracellular diffusion

Svyatoslav Kondrat Institute of Physical Chemistry Department of Complex Systems Warsaw, Poland

Contents

1	Introduction	2
2	Macromolecular self-diffusion 2.1 How to characterize diffusion?	2 3
3	How to model diffusion?3.1Langevin equation3.2Ermak-McCammon equation3.3Hydrodynamic interactions3.4Software packages	4 5 6 7 8
4	How to measure diffusion?	8
5	Effect of crowding on long-time diffusion coefficients5.1Experiments.5.2Simulations.	9 9 10
6	Anomalous subdiffusion	11
7	Effect of crowder's composition on diffusion	12
8	Conclusions and outlook	12

Lecture Notes of the $49^{\rm th}$ IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.



Fig. 1: *Crowded environment inside living cells.* (*left*) An artist's view of the cytoplasm. *Picture taken from [4]. (right) Snapshot from the Brownian dynamics simulations of Ref. [5]. Various macromolecules are depicted by different colors.*

1 Introduction

An important difference between a typical (non-biological) laboratory system and a living cell is an enormous variety and amount of metabolites and macromolecules diffusing, interacting and reacting in cells. Laboratory-scale macromolecular experiments often deal with volume fractions of the order of a few percents, while the volume taken up by the macromolecules in a biological cell is of the order of 20% to 50% [1] (Figure 1). The effects of such a crowded environment on macromolecular structure and reactivity inside living cells were first analysed in the early 80s [2, 3], but its importance for diffusion and reactions had not been appreciated until recently. For instance, Ellis wrote, concluding his 2001 review article [4], that crowding "should become a routine variable to study" and continued suggesting journals to "reject manuscripts on the grounds that this important variable has not been controlled."

Most physicochemical processes proceed differently in a biologically crowded environment: Diffusion slows down enormously and reactions may occur with different rates, in particular due to the reduced diffusion. Understanding diffusion in dense biological systems is, therefore, of critical importance for life sciences and for designing biotechnological applications. In this Lecture I will introduce the notion of Brownian motion (or self-diffusion) and we will discuss briefly how to describe, model and measure the diffusion properties (we shall restrict our attention to translational diffusion, however). The focus will be on physics and simulations, with a particular emphasis on the effects important for crowded, biologically relevant systems. Various aspects of crowding are also discussed in Chap. D2.

In this short Chapter it has not been possible to account for all relevant approaches and phenomena important for intracellular diffusion. I thus refer the interested readers to a few reviews [6–9], and note that this contribution shall merely be considered as a humble introductory note into the physics and modelling of diffusion in biophysically crowded environments inside living cells.

2 Macromolecular self-diffusion

Brownian motion or self-diffusion is the random motion of macromolecules (or other objects) suspended in a fluid. Such motion results from the collision of macromolecules with the fast-moving smaller molecules in a suspension, such as water and metabolites. This process shall be



Fig. 2: Normal and anomalous diffusion. Schematic of the behaviour of the mean-square displacement (MSD) with time. D_s and D_l are short-time and long-time self-diffusion coefficients, respectively, τ_B is the timescale of momentum relaxation, and τ_s and τ_l are the crossover times from the short-time and to the long-time normal regimes.

distinguished from a related process of diffusion, sometimes called transport diffusion, which is the natural motion of molecules from a region of high concentration to a region of low concentration (more precisely, high and low chemical potentials). Here, we shall deal exclusively with self-diffusion of macromolecules, particularly in a crowded environment inside living cells, but we shall briefly mention the transport diffusion in conclusions (Section 8). For brevity, however, we shall use self-diffusion and diffusion interchangeably where it does not lead to confusion.

2.1 How to characterize diffusion?

One way to quantify self-diffusion is to look at the mean square displacement (MSD) of a diffusing particle (the mean displacement is obviously zero):

$$MSD(t) = \langle [\boldsymbol{r}(t) - \boldsymbol{r}(0)]^2 \rangle, \tag{1}$$

where $\langle \cdots \rangle$ means *ensemble* averaging, r(t) is the molecule's position at time t and t = 0 is the beginning of observation. In practice the number of trajectories is often very limited, and *time* averaging within a small time frame is added to improve statistics (provided the system behaves ergodically, which is not always the case [10, 11]).

At very short times, MSD ~ t^2 due to momentum relaxation ($t < \tau_B \sim m/\gamma$, where m is the particle mass and γ the friction coefficient; see Fig. 2). This time scale is commonly omitted when discussing Brownian motion.

For a dilute system, the MSD typically behaves linearly with time for $t > \tau_B$; then

$$D = \lim_{t \to \infty} \frac{\text{MSD}(t)}{2dt}$$
(2)

defines the diffusion coefficient (here d is the dimensionality). The diffusion coefficient is one of the most important quantitative measures of molecular diffusion. It is analogous to the velocity in the classical mechanics and tells us the (approximate) distance a particle can travel in time t,

which is $l \approx \sqrt{Dt}$. The diffusion coefficient can be extracted from simulations (Section 3) and assessed experimentally (Section 4).

In a crowded system (such as the cytoplasm) the situation is more complex. In this case the MSD does not always behave linearly with time, and one can distinguish short-time and long-time self-diffusion coefficients, D_s and D_l , respectively, as schematically depicted in Fig. 2. In between these two *normal* regimes, *i.e.*, for $\tau_s < t < \tau_l$, there is a region of *anomalous* diffusion. The crossover time to the anomalous diffusion can be roughly estimated from the mean inter-macromolecular distances, assuming the anomalous regime commences on the time scale associate with the macromolecular collisions [12]. For macromolecules of the same radius *a* such simple considerations give

$$\tau_s \approx \frac{a^2}{D_s} \left([4\pi/(3\eta)]^{1/3} - 2 \right)^2,$$
(3)

where η is the packing fraction of macromolecules. Equation (3) shows that τ_s increases with the size of a macromolecule (D_s decreases with a), and it decreases monotonically with increasing the volume fraction. Specifically, $\tau_s \to \infty$ as $\eta \to 0$, as one may expect, and τ_s vanishes at $\eta_{\text{max}} = \pi/6 \approx 0.54$, manifesting that at high volume fractions a macromolecule encounters other macromolecules on the length scale of its own dimension; note that D_s depends on η due to hydrodynamic interactions (which, in general, contain many-body far and near-field contributions), but here we neglect this dependence for simplicity. This simple estimate turns out to be in a good agreement with simulations, at least in some range of volume fractions [12]. Taking as an example a = 2.5nm and the diffusion coefficient $D_s = 6.6$ Å²/ns (which corresponds to the tRNA/triphosphate isomerase), we find $\tau_s \approx 55$ ns for volume fraction 20% (density ≈ 5 mM). In the anomalous regime, one can define a time-dependent diffusion coefficient

$$D(t) = \frac{\text{MSD}(t)}{2dt} = \Gamma t^{\alpha - 1},$$
(4)

where α is the exponent of anomaly and Γ is the generalized transport (or anomalous diffusion) coefficient. For a crowded system, such as the cytoplasm, $\alpha < 1$; this is termed anomalous *subdiffusion*; in the opposite case, when $\alpha > 1$, it is called anomalous superdiffusion.

The crossover time to the long-time normal diffusion, τ_l , is not easy to estimate, but experiments and simulations suggest that τ_l is in the range from tens of microseconds to milliseconds [12– 14]. It is also possible that $\tau_l \to \infty$ [15], or that the system size is too small (*viz.*, the linear size $L < \sqrt{D_l \tau_l}$), in which case the long-time normal regime is not observed.

In addition to MSD and two diffusion coefficients, there are other quantities characterizing diffusion, such as rotational diffusion coefficients, mean dwell times, time autocorrelation functions, *etc.* However, for the present discussion, it will be sufficient to consider MSD and $D_{s,l}$.

3 How to model diffusion?

Ideally, one would like to perform '*ab initio*' simulations of a whole system, in order to approximate the reality as close as possible. Within the classical approach, this amounts to solving Newtonian equations of motions for all molecules in a system (for details see Chap. A10). For realistic biological systems, however, such simulations are dramatically expansive computationally. Although molecular dynamics simulations of crowded cytoplasm-like environments

do exist [9], the time scales covered so far are of the order of a few nanoseconds to at most microsecond, which is often below both τ_s and τ_l .

Another way to deal with such systems is to separate short and long time scales, and to consider 'averaged' equations of motions on longer time scales. This leads to the well-known Langevin equation, which we discuss below.

3.1 Langevin equation

In a biophysical context, a natural method of dividing a system into short and long time scales is to consider the motion of *macromolecules* and to average out over the motion of small molecules, such as water, ions and (some) metabolites. Then the Langevin equation for Nspherically-symmetrical macromolecules is $(i = 1, \dots, N)$

$$m_i \frac{d^2 \boldsymbol{r}_i}{dt^2} = -\nabla_i W(\{\boldsymbol{r}_k\}) - \sum_{j=1}^N \gamma_{ij} \frac{d\boldsymbol{r}_j}{dt} + \sum_{j=1}^N \sigma_{ij} \boldsymbol{\xi}_j(t),$$
(5)

where m_i is the mass of molecule *i*, *W* is the *effective* interaction potential, which takes into account solvent-mediated interactions, and $\nabla_i = \partial/\partial r_i$ is the Nabla operator, so that $F_i = -\nabla_i W$ is the effective force acting on molecule *i*.

Now, γ_{ij} is the $N \times N$ friction matrix (of $d \times d$ blocks, where d is the dimensionality) and ξ_i is the *stochastic* or fluctuating force, both due to the presence of solvent. The stochastic force is the white noise with zero mean, $\langle \xi_i \rangle = 0$, which is uncorrelated,

$$\langle \boldsymbol{\xi}_i(t)\boldsymbol{\xi}_j(t')\rangle = 2\delta_{ij}\delta(t-t'). \tag{6}$$

The coupling constants σ_{ij} , *i.e.*, the magnitudes of the random forces, are related to the friction matrix by the fluctuation-dissipation theorem

$$\gamma_{ij} = (k_B T)^{-1} \sum_k \sigma_{ik} \sigma_{kj},\tag{7}$$

where k_B is the Boltzmann constant and T kinetic temperature, as usual. Relation (7) can be obtained by integrating eqn (5) once and calculating the average of the velocity squared, $\langle v^2 \rangle$, and then employing the equipartition theorem $m_i \langle v_i^2 \rangle/2 = dk_B T/2$.

The Newtonian equations of motions (no solvent) are reproduced by setting $\sigma_{ij} = 0$ (and hence $\gamma_{ij} = 0$) and, since W is the effective potential, by replacing W by the interaction potential in the absence of solvent.

At low Reynolds numbers, *i.e.*, for a slowly moving, viscous fluid, momentum relaxations can be omitted at long time scales ($t \gg \tau_B$, Fig. 2). This means $dv_i/dt \ll \sum_j (\gamma_{ij}/m_i)v_j$, where $v_i = dr_i/dt$ is the velocity, and we obtain

$$\frac{d\mathbf{r}_{i}}{dt} = -(k_{B}T)^{-1}\sum_{j}D_{ij}\nabla_{j}W(\{\mathbf{r}_{k}\}) + (k_{B}T)^{-1}\sum_{jk}D_{ik}\sigma_{kj}\boldsymbol{\xi}_{j}(t),$$
(8)

where $D_{ij} = k_B T(\gamma^{-1})_{ij}$ is the diffusion matrix, which is inverse of the friction tensor. Equation (8) is the main equation for Brownian motion.

3.2 Ermak-McCammon equation

Instead of applying the strong friction limit, as above, one can integrate the Langevin equation directly on the time scales larger than the momentum relaxation ($t \gg m_i/\gamma_{ii}$). Assuming that the diffusion matrix is position dependent, but varies slowly in space, it is possible to obtain for the particle displacement [16]

$$\Delta \boldsymbol{r} = \boldsymbol{r}(\Delta t) - \boldsymbol{r}(0) = \Delta t \sum_{j} \nabla_{j} D_{ij}^{(0)} - \Delta t (k_{B}T)^{-1} \sum_{j} D_{ij}^{(0)} \nabla_{j} W^{(0)} + \boldsymbol{R}_{i}(\Delta t), \quad (9)$$

where the upper index in $D_{ij}^{(0)}$ and $W^{(0)}$ indicates that the diffusion matrix and the potential must be evaluated at t = 0 (in case they depend on time and/or positions of the macromolecules); $\mathbf{R}_i(\Delta t)$ is a random displacement averaged over time Δt , which satisfies [16]

$$\langle \mathbf{R}_i(\Delta t)\mathbf{R}_j(\Delta t)\rangle = 2D_{ij}^{(0)}\Delta t.$$
 (10)

Equations (9) and (10) have been derived by Ermak and McCammon [16] and can be viewed as a finite-difference (propagator) scheme for Brownian dynamics simulations. This scheme is equivalent to the first-order Euler algorithm for ordinary differential equations, and has been extended by Iniesta and de la Torre [17] to the second order Runge-Kutta approach by taking into account the second-order corrector step.

The Ermak and McCammon [16] derivation sets naturally the limits on the time step Δt . Indeed, on the one hand, it is clear that Δt must be greater than the timescale of the momentum relaxation, *i.e.*, $\Delta t \gg m_i/\gamma_{ii}$. On the other hand, Δt must be sufficiently small such that the interaction potential and the friction matrix can be considered constant during the time step Δt . For a single spherical particle in a homogeneous media eqn (9) simplifies to

$$\Delta \boldsymbol{r} = \boldsymbol{R}_i(\Delta t) = \sqrt{2D_0\Delta t} \,\boldsymbol{x},\tag{11}$$

where x is a random vector satisfying the Gaussian distribution, so that $\langle \mathbf{R}^2(\Delta t) \rangle = 2D_0\Delta t$. Using now eqn (11) to calculate particle's trajectory, and performing a sufficient number of independent simulations to gather enough statistics, we can calculate the MSD and extract the diffusion coefficient D using eqn (2) to confirm that it coincides with the diffusion coefficient D_0 used as an input in eqn (11).

For N spherical macromolecules in a homogeneous medium one has

$$\Delta \boldsymbol{r} = -\Delta t (k_B T)^{-1} \sum_j D_{ij} \nabla_j W^{(0)} + \boldsymbol{R}_i (\Delta t).$$
(12)

The random displacement is

$$\boldsymbol{R}_{i} = \sqrt{2\Delta t} \sum_{j} B_{ij} \cdot \boldsymbol{x}_{j}, \qquad (13)$$

where x_j is a Gaussianly distributed random vector, the dot means a convolution over d components of x_j (note that B is a $N \times N$ matrix of $d \times d$ blocks), and

$$D_{ij} = \sum_{k} B_{ik} B_{kj}.$$
(14)

Thus, in order to calculate the random force one needs to take a 'square root' of the diffusion matrix. Since essentially any B that satisfies eqn (14) can be used in eqn (12), one often takes the standard Cholesky decomposition [18]. In the case discussed, however, the matrix D is diagonal hence the 'decomposition' can be performed in a straightforward way and only once before a simulation starts (assuming that the diffusion constants do not change in time). However, this changes when the hydrodynamic interactions are taken into account. We address this in the next section.

3.3 Hydrodynamic interactions

As we have discussed, in Brownian dynamics simulations the water, metabolites and other small molecules inside a cell are not taken into account explicitly, but effectively present a viscous environment for diffusion of macromolecules. A macromolecule moving in this viscous medium excites a long-range flow that affects other molecules, resulting in an effective interaction between the macromolecules.

In order to account for such hydrodynamic interactions, one needs to solve the Navier-Stokes equation for the solvent for current positions and velocities of macromolecules each simulation step. This can be done by calculating the viscous stress γ_h exerted on a macromolecule by other macromolecules and noticing that the force acting on this macromolecule is $\gamma_h v$, where v is its velocity. This permits us to incorporate the hydrodynamic interactions into the diffusion matrix $D_h = (k_B T) \gamma_h^{-1}$. Assuming that such a viscous stress is a superposition of independent contributions from each macromolecule, and further neglecting the effects related to the macromolecule sizes, which shall be valid for large distances, it is possible to obtain for macromolecules of the same radius a [19, 20]

$$D_{ij} = \frac{k_B T}{8\pi\nu r_{ij}} \left\{ I + \frac{\boldsymbol{r}_{ij} \otimes \boldsymbol{r}_{ij}}{r_{ij}^2} + \frac{2a^2}{r_{ij}^2} \left(\frac{1}{3}I - \frac{\boldsymbol{r}_{ij} \otimes \boldsymbol{r}_{ij}}{r_{ij}^2} \right) \right\},\tag{15}$$

where $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$ is a vector connecting the centers of macromolecules *i* and *j*, \otimes denotes tensor product, *I* is the 3 × 3 unit matrix, and ν is the fluid viscosity; the diagonal components are $D_{ii} = (k_B T / 6\pi\nu a)I$. Equation (15) is known as the Rotne-Prager-Yamakawa (RPY) tensor. It can be extended to macromolecules of different sizes [21], and, formally, to distances $r_{ij} < 2a$ [19, 20]. At short distances, however, also the lubrication and many-body forces must be taken into account [8]. To avoid this complication, a macromolecule can be approximated by small spheres placed on its surface [22–25]; then the RPY tensor between these small spheres is valid in a wide range of macromolecular separations, but the tensor size increases proportionally to the number of spheres taken to approximate the macromolecules.

In any case, at *each* Brownian dynamics (BD) simulation step the diffusion matrix must be factorized in order to obtain the hydrodynamically correlated random displacements (see eqn (13) and (14)). In addition, since such hydrodynamic forces are long-ranged (eqn (15)), it is necessary to take into account not only the macromolecules present in a computational box, but also all their periodic images (in simulations of a bulk system). This can be done by using the so-called Ewald summation [26]. All this increases the computational cost of BD simulations significantly. For instance, the frequently used Cholesky decomposition, which we have already mentioned, scales as N^3 , where N is the number of macromolecules; the pair-wise inter-macromolecular interactions scale as N^2 , and the Ewald summation as N^2 or $N \log(N)$, depending on the method used; thus the hydrodynamics are the bottleneck of BD simulations. Their computational cost can be lowered by using Chebyshev approximation for the diffusion matrix developed by Fixman [27], but it still scales as $N^{2.5}$ [7]. Geyer and Winter [28] have proposed an approximate N^2 algorithm, which is based on a certain Ansatz (expansion) for the random force, with the unknown parameters determined approximately from the appropriate variance-covariance relation akin of eqn (6). It shows a good agreement with the standard approaches, at least for the tested systems [29, 30]. There are also other important methods and improvements [31–35]. Of particular interest is the Stockesian dynamics [31], which allows to take into account both long-range many body and short-range lubrication forces. Discussion of these methods is beyond the scope of this Chapter, but various approaches to hydrodynamics are discussed in Chap. A10.

3.4 Software packages

There are a few open-source packages that offer ready-to-use software codes for Brownian dynamics (BD) simulations. The Brownian and Langevin dynamics have been implemented in the standard Gromacs [36] and LAMMPS [37] simulation packages. However, Gromacs does not not seem to include the hydrodynamic forces, essential for crowded systems, as we shall see (Section 5.2). In LAMMPS [37], the hydrodynamics are implemented using lattice-Boltzmann approach. Brownmove [38] is an implementation of the Geyer and Winter N^2 algorithm [28]. BDpack [39] is a software package that implements the recently introduced matrix-free method of Saadat and Khomami [35] for dilute and semi-dilute polymeric solutions. BD_BOX [40] is probably most versatile and highly optimized Brownian dynamics simulation package for rigid and flexible molecules [41], but it takes into account only the long-range hydrodynamics. The GPU-optimized HOOMD-blue simulation toolkit [42] can also run BD simulations, and the hydrodynamics can be added via their plug-in system (*e.g.*, with RPY tensor as in Ref. [43]). The short and long range hydrodynamics have recently been implemented in ESPResSo simulation package [44], but its performance has not been optimized (Christian Holm, personal communication).

Concluding, it seems that a well-developed and optimized open-source software package for BD simulations, which supports both long-range and short-range hydrodynamics, is not currently available.

4 How to measure diffusion?

Although the focus of this Chapter is physics and modelling of intracellular diffusion, in order to understand better the connection with experiments, we shall briefly mention the three main techniques developed for measuring diffusion properties. For details of these (fluorescentbased) methods see Chap. A3.

Perhaps conceptually the simplest (but technically advanced) is the *single-particle tracking* (SPT). Such experiments amount to introducing a fluorescent dye into traced macromolecules and video-recording their diffusion. The MSD can be extracted from the recorded videos by analysing tracer trajectories, similarly as in simulations. Typically the spatial resolution of SPT is of the order of a few nanometres, with the time resolution on the scale of milliseconds, but high-speed tracking techniques exist allowing for the resolution of tens of microseconds [45, 46].

Probably the most frequently used method is the *fluorescent correlation spectroscopy* (FCS). Similarly as SPT, it relies on labeling tracers by a fluorescent dye, but in FCS the fluctuating



Fig. 3: *Effect of crowding on diffusion from in vitro experiments. Diffusion coefficient of a rhodamine green as a function of the crowder concentration (Ficoll-70). The inset shows the same data in log scale. The figure shows the results of a FCS study. Reproduced from Ref. [52].*

fluorescent light intensity is measured, rather than the trajectories. The measured intensity can be related to the time autocorrelation function, from which the diffusion coefficient and the anomaly exponent can be extracted by fitting to the analytical expressions [47, 48]. The advantage of FCS is the time resolution, which can be of the order of microseconds.

The *fluorescent recovery after photobleaching* (FRAP) is similar to FCS. Here, however, a small region in a sample with fluorescing macromolecules is initially bleached by an intense laser pulse, and the fluorescent light intensity (rather than fluctuations) is monitored as the bleached region recovers due to the diffusion of the fluorescent macromolecules from the outside of this region [49–51].

All three methods are applicable *in vivo* by in-cell expression of fluorescent macromolecules, frequently a green fluorescent protein (GFP).

5 Effect of crowding on long-time diffusion coefficients

5.1 Experiments

Experiments indicate that the long-time diffusion coefficients of macromolecules are dramatically reduced in the crowded environments inside living cells. For instance, the diffusion coefficient of GFP *in vivo* has been measured to be 10 - 15 times lower than in a dilute solution [53, 54].

A more systematic analysis can be carried out *in vitro*, where the concentration of crowders can be easily controlled. Figure 3 shows the results of a FCS study of the long-time diffusion of rhodamine green in a concentrated solution of Ficoll-70 as crowders [52]. It has been found that the diffusion coefficient D_w decreases exponentially with the crowder concentration C, $viz., D_w = D_w^0 \exp\{-aC\}$, where a is a fitting parameter and D_w^0 is the diffusion coefficient in infinite dilution. Interestingly, this work has also demonstrated that the reduction in diffusion is comparable for large macromolecules and for small solutes [52] (the plot not reproduced here).



Fig. 4: *Effect of crowding on diffusion from Brownian dynamics simulations.* (*left*) *Diffusion constant of GFP in the cytoplasm for different models (steric; steric and electrostatic; steric, electrostatic and Lennard-Jones with different interaction parameters). The arrow points out to the value of the depth of the Lennard-Jones interactions producing the observed 10-fold decrease of the diffusion coefficient. Reproduced from Ref. [5]. (right) Hydrodynamic interactions describe correctly the reduction of the diffusion coefficient in the cytoplasm (full circles), while the Lennard-Jones (van der Wals) interactions may overestimate it due to clustering (open rectangles). Thin vertical line denotes GFP. Reproduced from Ref. [55]*

5.2 Simulations

From a modeling perspective, it shall be clear that at least short-range repulsive interactions between macromolecules must be take into account. McGuffee and Elcock [5] have shown, however, that hard core (steric) interactions alone are not sufficient to reproduce the observed decrease of the diffusion coefficients. These authors considered a cytoplasm model consisting of 50 different types of macromolecules (Fig. 1), and studied the effect of steric, electrostatic and van der Waals interactions on diffusion. Using Brownian dynamic simulations, as described in Section 3.2 (without hydrodynamic interactions), they showed that steric and electrostatic interactions are not sufficient to explain the observed slow-down of diffusion. However, the long-time self-diffusion coefficient turns out to be sensitive to the van der Waals interactions, and it seems possible to match the simulated D_l for GFP with the diffusion coefficient measured in experiments by varying the depth of the Lennard-Jones potential within the physically reasonable range (Fig. 4).

However, Ando and Skolnick [55] have argued that van der Waals interactions may lead to clustering and overestimate the slow-down for larger macromolecules. In contrast, the hydrody-namic interactions (with short and long-range contributions) adequately reproduce the observed reduction of the diffusion coefficient (particularly of GFP) in the cytoplasm (Fig. 4).

Another important aspect of Ando and Skolnick's work [55] is that the details of macromolecule's structure seem to be of negligible importance for diffusion. This has been shown by comparing directly the translational diffusion coefficients in a molecularly-shaped system and in a system where the macromolecules were modeled as spherical particles.



Fig. 5: Anomalous subdiffusion from in vitro experiments. The exponent of anomalous subdiffusion, α , as a function of the crowder (obstacle) concentration. (left) Exponent α for streptavidin in dextran of different sizes. (right) Exponent α for streptavidin and GFP in dextran 276.5 kDa, and for fluorescein and FITC-Dextran in 401.3 kDa dextran. The figure shows the results of a FCS study. Reproduced from Ref. [56].

6 Anomalous subdiffusion

So far we have discussed diffusion focusing on its long-time normal behaviour. However, there is a 'significant body of evidence' indicating that *in vivo* diffusion is anomalous on extended time scales. Weiss et al. have even proposed to use the anomaly exponent (α in eqn (4)) as a measure of how crowded a system is [57]. In their experiments, Weiss et al. studied the diffusion of dextran inside HeLa cells by FCS and found the anomaly exponent in the range between $\alpha = 0.73$ and $\alpha = 0.79$. Similar values have been obtained in other experiments [58, 59].

As the structural properties of the cell interior are not well known and easily controllable, a more systematic analysis can be carried out *in vitro*, which allows to control the crowder's concentration and size. A FCS study shows [56], in particular, that the anomaly exponent α decreases with increasing the crowder's volume fraction ϕ and saturates at $\alpha \approx 0.75$ at high ϕ (Fig. 5); it can be well fitted by $\alpha(\phi) = \alpha_1 + \exp(-\phi/\phi_0)$, where ϕ_0 and α_1 are fitting parameters. The saturation is likely due to the entanglement of polymer chains of crowders (dextran in this case) at high concentrations so that effectively the tracer diffuses inside a cross-linked polymer network, which gives $\alpha = 3/4 \approx \alpha_1$ [56]. Interestingly, the anomalous diffusion shows the same features for globular tracers (streptavidin and GFP), but the diffusion is normal or only slightly anomalous for polymer tracers (fluorescein and FITC-Dextran, Fig. 5), suggesting a different physics governing their diffusion [60, 61].

While normal diffusion is characterized by the universal Gaussian distribution, anomalous subdiffusion is non-universal and can be due to a variety of reasons [15, 62–64]. Detailed discussion of all possible mechanisms is out of scope of this Lecture; we shall only mention that the frequently observed mechanism, particularly in eukaryotes, is fractional Brownian motion, which corresponds to diffusion in viscoelastic-like medium [65]. Macromolecular trapping with varying trapping times, which is described by the so-called continuous time random walk model, has been demonstrated to take place on short time scales [59].



Fig. 6: *Effect of composition on diffusion.* Snapshots from Brownian dynamics simulations of (a) dense two-component system and (b) model cytoplasm. The volume fraction occupied by macromolecules is approximately 19% in both cases. (c) Comparison of the mean-square displacement (MSD) for a macromolecule of size 2.5nm in the model cytoplasm and in a two-component system. The cytoplasm has a volume fraction of 19.2%, and the two-component system 18.6% and 11.6%; the concentration of smaller molecules is 0.875. Reproduced from Ref. [12]

7 Effect of crowder's composition on diffusion

Living cells are characterized by constant changes of the cell constituents, causing the variation of the *relative concentrations* of macromolecules. To understand how this affects diffusion, we have performed [12] Brownian dynamics simulations of a two-component system (Fig. 6a), where the effect of composition can be studied more systematically, and compared these results with the results for a model cytoplasm (Figure 6b).

Figure 6c demonstrates that diffusion depends sensitively on the relative concentrations of macromolecules. In particular, the (middle-time) diffusion in the cytoplasm is faster than in the two component system with the same packing fraction. Similarly, the diffusion is comparable in the cytoplasm and in the two component system (the latter with the volume fraction 12%), although the cytoplasm is over 50% 'more crowded' than the two-component system.

These results raise two interesting questions. Firstly, artificial crowders (such as dextran) are often used in *in vitro* experiments (Sections 5.1 and 6) to mimic the environment inside living cells. However, Fig. 6c suggests that such nearly monodisperse crowders may in fact be inadequate for this purpose. Secondly, the volume fraction is frequently used as a measure of crowdedness, while Fig. 6c shows that it *cannot* serve as a unique measure of how crowded a system is. Indeed, diffusion depends sensitively on molecular composition, and can in fact be faster in systems with higher volume fractions.

8 Conclusions and outlook

We have discussed the main principles of how to model macromolecular self-diffusion, and briefly reviewed one of the most important manifestations of the crowded world inside living cells, which is a dramatic slow-down of the macromolecular diffusion as well as its anomalous behaviour. From a modelling perspective, the steric and hydrodynamic forces seem to play the key role in the reduction of the diffusion coefficient [55]. However, the atomistic details of macromolecules do not seem to be of any significant importance for transport diffusion [55],

which shall alleviate the computational burden of many Brownian dynamics simulations. Interestingly, the macromolecular composition affects considerably the macromolecular diffusion, which can be faster in systems with higher volume fractions [12].

Many questions remain unanswered, however. I shall only briefly mention some of them.

- Diverse experimental results have been reported for diffusion on long time scales. While some experiments show the slow-down of the normal diffusion [52–54], strong evidence exists in support of the anomalous subdiffusion [56, 57, 59, 65]. Whether and when the anomalous diffusion turns into the normal regime is not clear. It is important to note in this context that the characteristic times of the anomalous diffusion reported in the experiments are of the order of up to a few seconds; the Brownian dynamics (BD) simulations, discussed in this Lecture, are currently unable to deal with such long time scales. On the other hand, the temporal resolution of a typical experiment is too low to access the time scales of BD simulations.
- We have discussed self-diffusion, which is relevant to macromolecules, since they are typically present in cells in very low 'copy-numbers'. However, the concentration of smaller molecules (metabolites) is often relatively high, in which case the *transport* diffusion becomes a relevant process. The transport and self-diffusion can differ dramatically, but this topic has not been touched in the biophysical context.
- Spectacular behaviour has recently been reported by Parry et al. [66], who observed a dramatic slow-down of diffusion in cells with suppressed metabolic activity (*i.e.*, no or few reactions taking place). The origin of this effect is not yet well understood, but it points out to an interesting inter-dependence of *in vivo* reactions and diffusion. Clearly, diffusion controls the rate of (diffusion-limited) reactions, but it turns out that also reactions can influence the intracellular diffusion in a dramatic way.
- I hope I have convinced you that the intracellular diffusion is important and interesting; however, it is essentially the reactions that make Life. Incorporating reactions in the Brownian dynamics simulations is a difficult and computationally expansive task. Although a number of multiscale and coarse-grained approaches have been introduced [67–70], a well-developed reliable framework does not seem to exist. The development of such a framework, which would allow for spatially-resolved whole-cell simulations, will likely be the focus of future research activities. In combination with the advanced experimental studies this will bring new discoveries and a better understanding of the Physics of Life.

Acknowledgment

I am grateful to A. Cherstvy (Uni Potsdam) and M. Długosz (Warsaw University) for a critical reading of this Chapter and for fruitful comments and suggestions.

References

- [1] S. B. Zimmerman and S. O. Trach, J Mol Biol. 5, 599 (1991).
- [2] A. P. Minton, Biopolymers 20, 2093 (1981), ISSN 1097-0282.

- [3] A. B. Fulton, Cell **30**, 345 (1982).
- [4] R. J. Ellis, Curr. Opin. in Struct. Biology 11, 114 (2001).
- [5] S. R. McGuffee and A. H. Elcock, PLoS Comput Biol 6, e1000694 (2010).
- [6] A. H. Elcock, Curr. Opin. Struct. Biol. 20, 196 (2010).
- [7] M. Długosz and J. Trylska, BMC Biophys. 4, 3 (2011).
- [8] J. Skolnick, J. Chem. Phys. 145, 100901 (2016).
- [9] M. Feig, I. Yu, P. Wang, G. Nawrocki, and Y. Sugita, J. Phys. Chem. B 121, 8009 (2017).
- [10] S. K. Ghosh, A. G. Cherstvy, and R. Metzler, Phys. Chem. Chem. Phys. 17, 1847 (2015).
- [11] S. K. Ghosh, A. G. Cherstvy, D. S. Grebenkov, and R. Metzler, New J. Phys. 18, 013027 (2016).
- [12] S. Kondrat, O. Zimmermann, W. Wiechert, and E. v. Lieres, Phys. Biol. 12, 046003 (2015).
- [13] E. Vilaseca, I. Pastor, A. Isvoran, S. Madurga, J. L. Garcés, and F. Mas, Theor. Chem. Acc. 128, 795 (2011).
- [14] P. M. Blanco, M. V. J. L. Garcés, S. Madurga, and F. Mas, Entropy 19, 105 (2017).
- [15] F. Höfling and T. Franosch, Rep. Prog. Phys. 76, 046602 (2013).
- [16] D. L. Ermak and J. A. McCammon, J. Chem. Phys. 69, 1352 (1978).
- [17] A. Iniesta and J. G. de la Torre, J. Chem. Phys. 92, 2015 (1990).
- [18] G. H. Golub and C. F. van Loan, Matrix Computations (Baltimore: Johns Hopkins, 1996).
- [19] J. Rotne and S. Prager, J. Chem. Phys. 50, 4831 (1969).
- [20] H. Yamakawa, J. Chem. Phys. 53, 436 (1970).
- [21] J. G. de la Torre and V. A. Bloomfield, Biopolymers 16, 1747 (1977).
- [22] M. Długosz and J. M. Antosiewicz, J. Phys. Chem. B, 119, 8425 (2015).
- [23] J. W. Swan and G. Wang, Physics of Fluids 28, 011902 (2016).
- [24] M. Długosz and J. M. Antosiewicz, J. Phys. Chem. B 120, 7114 (2016).
- [25] J. M. Antosiewicz, K. Kamiński, and M. Długosz, J. Phys. Chem. B 121, 8475 (2017).
- [26] C. W. J. Beenakker, J. Chem. Phys. 85, 1581 (1986).
- [27] M. Fixman, Macromolecules 19, 1204 (1986).
- [28] T. Geyer and U. Winter, J. Chem. Phys. 130, 114905 (2009).
- [29] T. Geyer, BMC Biophysics 4, 7 (2011).
- [30] R. R. Schmidt, J. G. H. Cifre, and J. G. de la Torre, J. Chem. Phys. 135, 084116 (2011).
- [31] A. J. Banchio and J. F. Brady, J. Chem. Phys. 118, 10323 (2003).
- [32] T. Ando, E. Chow, and J. Skolnick, J Chem. Phys. 139, 121922 (2013).
- [33] T. Ando, E. Chow, Y. Saad, and J. Skolnick, J. Chem. Phys. 137, 064106 (2012).
- [34] A. Saadat and B. Khomami, J. Chem. Phys. 140, 184903 (2014).
- [35] A. Saadat and B. Khomami, Phys. Rev. E 92, 033307 (2015).
- [36] Gromacs, http://www.gromacs.org/.
- [37] LAMMPS, http://lammps.sandia.gov/.
- [38] T. Geyer, *Brownmove*, URL http://gepard.bioinformatik. uni-saarland.de/services/brownmove.
- [39] BDPack, URL http://amir-saadat.github.io/BDpack/.
- [40] M. Długosz and P. Zielinski, BD_BOX, URL http://www3.cent.uw.edu.pl/ ~mdlugosz/downloads.html.
- [41] M. Dlugosz, P. Zielinski, and J. Trylska, J. Comput. Chem. 32, 2734 (2011).
- [42] HOOMD-blue, URL http://glotzerlab.engin.umich.edu/hoomd-blue/.
- [43] Z. Varga, G. Wanga, and J. Swan, Soft Matter 11, 9009 (2015).
- [44] ESPResSo, URL http://espressomd.org/wordpress/.
- [45] A. Kusumi, C. Nakada, K. Ritchie, K. Murase, K. Suzuki, H. Murakoshi, R. S. Kasai,

J. Kondo, and T. Fujiwara, Annu. Rev. Biophys. Biomol. Struct. 34, 351 (2005).

- [46] W. J. Greenleaf, M. T. Woodside, and S. M. Block, Annu. Rev. Biophys. Biomol. Struct. 36, 171 (2007).
- [47] O. Krichevsky and G. Bonnet, Rep. Prog. Phys. 65, 251 (2002).
- [48] S. T. Hess, S. Huang, A. A. Heikal, and W. W. Webb, Biochemistry 41, 697 (2002).
- [49] E. A. J. Reits and J. J. Neefjes, Nature Cell Biology 3, E145 (2001).
- [50] J. Lippincott-Schwartz, E. Snapp, and A. Kenworthy, Nature Reviews Molecular Cell Biology 2, 444 (2001).
- [51] A. S. Verkman, in *Biophotonics, Part A* (Academic Press, 2003), vol. 360 of *Methods in Enzymology*, pp. 635 648.
- [52] E. Dauty and A. S. Verkman, J. Mol. Recognit. 17, 441 (2004).
- [53] M. B. Elowitz, M. G. Surette, P.-E. Wolf, J. B. Stock, and S. Leibler, J. Bacteriol. 181, 197 (1999).
- [54] M. C. Konopka, I. A. Shkel, S. Cayley, M. T. Record, and J. C. Weisshaar, J. Bacteriol. 188, 6115 (2006).
- [55] T. Ando and J. Skolnick, Proc. Natl. Acad. Sci. USA 107, 18457 (2010).
- [56] D. S. Banks and C. Fradin, Biophys. J. 89, 2960 (2005), ISSN 0006-3495.
- [57] M. Weiss, M. Elsner, F. Kartberg, and T. Nilsson, Biophys. J 87, 3518 (2004).
- [58] H. Engelke, D. Heinrich, and J. O. Rädler, Phys. Biol. 7, 046014 (2010).
- [59] J.-H. Jeon, V. Tejedor, S. Burov, E. Barkai, C. Selhuber-Unkel, K. Berg-Søorensen, L. Oddershede, and R. Metzler, Phys. Rev. Lett. 106, 048103 (2011).
- [60] Y. Wang, C. Li, and G. J. Pielak, J. Am. Chem. Soc. 132, 9392 (2010).
- [61] J. Shin, A. G. Cherstvy, and R. Metzler, New J. Phys. 16, 053047 (2014).
- [62] T. Geyer, J. Chem. Phys. 137, 115101 (2012).
- [63] D. Ernst, J. Köhler, and M. Weiss, Phys. Chem. Chem. Phys. 16, 7686 (2014).
- [64] Y. Meroz and I. M. Sokolov, Physics Reports 573, 1 (2015).
- [65] D. Ernst, M. Hellmann, J. Köhler, and M. Weiss, Soft Matter 8, 4886 (2012).
- [66] B. R. Parry, I. V. Surovtsev, M. T. Cabeen, C. S. O'Hern, E. R. Dufresne, and C. Jacobs-Wagner, Cell 156, 183 (2014).
- [67] D. C. Wylie, Y. Hori, A. R. Dinner, and A. K. Chakraborty, J. Phys. Chem. B 110, 12749 (2006).
- [68] G. Kalantzis, Comput. Biol. Chem. 33, 205 (2009).
- [69] M. Flegg, J. Chapman, and R. Erban, J. R. Soc., Interface 9, 859 (2011).
- [70] S. Kondrat, O. Zimmermann, W. Wiechert, and E. v. Lieres, Eur. Phys. J. E. 39, 11 (2016).

D 2 Macromolecular Crowding: Colloidal Suspensions as Models for Biological Systems

P. R. Lang, Y. Liu ICS-3 Soft Condensed Matter Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Introduction		2
2	Basics of important theoretical approaches		3
	2.1 Asakura-Oosawa and Vrij theory of depletion interaction		3
	2.2 Free volume theory		4
	2.3 Scaled particle theory		5
3	Crowding and colloidal phase behavior		6
	3.1 Phase behavior of hard sphere suspensions		6
	3.2 Depletion interaction and colloidal phase behavior		9
4	Reaction equilibria and kinetics		10
	4.1 Equilibrium constants		10
	4.2 Reaction kinetics		12
5	Conclusions		14

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Many biological systems are containing high contents of various macromolecular and colloidal species, they may thus be considered as an aqueous suspension with a high volume fraction of dispersed material. Human blood contains roughly 45 % of cellular material (erythrocytes, leucocytes and thrombocytes) accompanied by about 6-8 % of proteins. The cytoplasm of cells hosts a total volume fraction of the order of 20% of of macromolecules comprising mainly proteins, polysaccharides, DNA and RNA and the nucleus of eucaryotic cells is besides harboring the genome, cramped with bio-polymers of different type and function. Although the concentration of an individual species may be quite low, like thrombocytes and leucocytes, which make less than one percent of the blood volume, the total amount of suspended material can be very large. Such systems are usually called "crowded".

From the perspective of polymer and colloid physics, it is evident that the properties of suspensions and the contained macromolecules will be strongly influenced by their crowded environment. Nevertheless for long time, the standard procedure in the field of biochemistry was to characterize and describe, e. g. proteins on the basis of in vitro experiments carried out in dilute solutions [1]. While the importance of crowding and the resulting excluded volume effects for cell biology was recognized almost four decades ago [2, 3], the significance of crowding for conformation and functioning of biopolymers and entire cells have gained increasing interest only over the last fifteen to twenty years [4, 5] nowadays resulting in several hundreds of scientific papers per year [6].

At mass contents in the range of several hundred mg/mL biological macromolecules do not only interact by site specific interactions but they mutually affect each other by non-specific interactions as van der Waals attraction, hydrophobic and electrostatic interactions as well as excluded volume effects, where the latter are often dominating at physiological conditions. Excluded volume effects usually also referred to as crowding or macromolecular crowding, embrace all properties of macromolecules which are due to the fact that they have less volume available as compared to infinitely dilute suspensions. Under these constraints polymer conformations, diffusion rates, reaction equilibria and kinetics may be drastically different from those occurring in dilute systems. Proteins react to their environment by adjusting their size and shape adopting states spanning from unfolded chains, which resemble ideal polymer coils, to more or less globular assemblies, showing characteristics of nano-particles. In recent years, the impact of crowding onto folding and stability of proteins has been investigated theoretically [7, 8, 9], by computer simulations [10, 11, 12] and experimentally [13, 14, 15, 16].

At sufficiently high concentrations, excluded volume effects can lead to conformational changes and phase separations, which may be leading to the formation of compartments, which coexist in the nuclei of eucaryotic cells despite the absence of separating membranes [17, 18, 19]. Further, inert crowder particles can suppress or promote protein-ligand binding depending on the size and shape ratio between the reaction partners and the product [20, 21].

Beyond the realms of cell biology, concentrated suspensions of biological macromolecules play an important role in various fields. In food industry concentrated mixtures of proteins and polysaccharides are applied in attempts to stabilize water in water emulsions for oil free food formulations [22]. The crowding induced short range order of α -, β - and γ crystalline in the eye lens is of paramount biomedical importance, since it is essential for the lens' transparency [23, 24, 25]. And last but not least the so called depletion effect [26, 27, 28] is applied in protein crystallization for their characterization by x-ray diffraction [29, 30].



Fig. 1: Illustration of excluded volume and depletion interaction in a suspension of colloids (big spheres) and depletants (small spheres). The gray areas represent depletion zones, the sum of which together with the volume of all colloids constitutes the excluded volume for the small spheres. Whenever depletion zones overlap (blue area), the excluded volume is reduced by the negative overlap volume and the entropy of the system is increased.

extend [31, 32, 33]. We will refrain from addressing this subject, as it is discussed in the dedicated chapter D.1 of this book, and even reviewing the plethora of the recent work on static properties in detail is beyond the scope of this contribution. We will rather focus on two major subjects, chosen by personal bias, where we can provide qualitative explanations for some generic features on the the basis of relatively simple approaches from soft matter physics.

2 Basics of important theoretical approaches

2.1 Asakura-Oosawa and Vrij theory of depletion interaction

In any concentrated or crowded suspension, a large part of the available volume is not accessible to particles due to the presence of other particles. Let us consider a colloidal suspension consisting of two types of particles: large spheres and small spheres (depletants), as illustrated in Fig.1. The small spheres can not enter the depletion zones around the large spheres. The sum of the big spheres' volume and the sum of the depletion zones' volumes constitute the volume from which the small spheres are excluded. When two large spheres approach each other, their depletion zones overlap and the total excluded volume is reduced. In other words, the total available volume for small spheres to move freely is increased. The resulting gain of entropy for the system will effectively generate an attractive force between large particles. The scale of the depletion force or potential can be estimated by the Asakura-Oosawa-Vrij theory [26, 27, 28].

For the situation of two kinds of differently sized spheres with radii R_c and r_d the interaction potential between two large spheres can be easily calculated on the bases of geometrical considerations and the assumption that the small spheres thermodynamically behave as an ideal gas

(see for example [34, 35])

$$\frac{\Phi(h)}{k_B T} = \begin{cases} -\rho_d \frac{\pi}{6} (2r_d - h)^2 \left(3R_c + 2r_d + \frac{h}{2} \right) & \text{for } h \le 2r_d \\ 0 & \text{for } h > 2r_d \end{cases} ..$$
(1)

Here h is the surface-to-surface separation between two large spheres, ρ_d is the number density of the small spheres, k_B is Boltzmann's constant and T is the absolute temperature.

If the depletant species is not a sphere but an ideal polymer coil, eq.1 can be used as a good approximation to calculate the pair potential between the large spheres, if the small spheres' diameter is replaced by the thickness of the depletion zone δ , which depends on the ratio of the polymer size over the sphere size, $q = R_q/R_c$. Here R_q is the polymer radius of gyration and

$$\delta = \frac{R_g}{q} \left[\left(1 + \frac{6q}{\sqrt{\pi}} + 3q^2 \right)^{1/3} - 1 \right]$$
(2)

Crowded biological systems are not categorically different from the case of colloidal suspensions. Thus biological macromolecules are subjected to depletion forces which might effect their tendency to form assemblies and their phase behavior.

2.2 Free volume theory

In order to predict phase transitions in crowded systems, a treatment of excluded volume effects on a pair level, as discussed in the preceding section will not be sufficient. The general conditions for two or more phases coexisting are that the chemical potentials of all components in all phases are equal and that the pressure in the coexisting phases is the same. To calculate these properties in systems consisting of colloidal particles and a depleting polymer, the so-called free volume theory [36, 37] was successfully applied. The theory starts from the calculation of the semi-grand potential of a system of N_c colloids and N_d polymers

$$\Omega(V, T, N_c, \mu_d) = F(V, T, N_c, N_d) - \mu_d N_d, \tag{3}$$

where μ_d is the chemical potential of the depletant and $F(V, T, N_c, N_d)$ is the system's Helmholtz free energy at constant composition. We can rewrite this equation as

$$\Omega(V, T, N_c, \mu_d) = F_0(V, T, N_c) - \int N_d(\mu_d) d\mu_d,$$
(4)

making use of the thermodynamic relation $(\partial \Omega / \partial \mu_d)_{V,T,N_c} = -N_d$, where $F_0(V,T,N_c)$ is the Helmholtz free energy of the pure colloidal system. Since expressions for the latter are available (see e. g. [34]) we only need to calculate N_d . To this end the system is allowed to osmotically equilibrate with an (infinitely large) reservoir of a depletant solution, via a hypothetical membrane, which is permeable for solvent and depletant molecules, but not for the colloids. By this equilibrium, the chemical potential of the polymer is the same in the reservoir as in all coexistent phases in the system, although the polymer may partition among the phases at different concentrations. Thus

$$\mu_d = \mu_d^0 + k_B T \ln \frac{N_d}{V_{free}} \tag{5}$$

and

$$\mu_d = \mu_d^0 + k_B T \ln \rho_d^R \tag{6}$$

where V_{free} is the total system volume minus the excluded volume and ρ_d^R is the depletant number density in the reservoir. By equating eqs. 5 and 6 we obtain an expression for the number of depletants in the system

$$N_d = \rho_d^R V_{free}.\tag{7}$$

Using the Gibbs-Duhem relation $\rho_d^R d\mu_d = d\Pi^R$ for the osmotic pressure in the reservoir, we can solve the integral in eq. 4 to get

$$\Omega(V, T, N_c, \mu_d) = F_0(V, T, N_c) - \Pi^R V_{free}^0,$$
(8)

where we further introduced the main approximation of the theory, by equating the free volume with the free volume in the pure colloid system, V_{free}^0 . This implies that the depletants behave as an ideal gas, not causing any excluded volume themselves. Under this assumption, the pressure is also given by the ideal gas analogue $\Pi^R = \rho_d^R k_B T$.

Since for certain size ratios q multiple overlap between depletion zones is possible, their overlapping volumes are not pairwise additive, and the calculation of the free volume is all but trivial. Here we use again the fact that the chemical potential of the depletant is fixed. We can write it as

$$\mu_d = \mu_d^0 + k_B T \ln \frac{N_d}{V} + W \tag{9}$$

where W is the reversible work needed to insert an additional depletant particle into the system. Equating this formulation with eq.5, keeping in mind the approximation $V_{free} = V_{free}^0$, we obtain the fraction of the free volume as

$$\alpha = \frac{V_{free}^0}{V} = \exp\left\{-W/k_BT\right\}.$$
(10)

Now the only missing ingredient is the work W, which can be estimated from scaled particle theory as outlined in the next section.

2.3 Scaled particle theory

The scaled particle theory [38] was originally developed to derive expressions for the pressure and the chemical potential of hard sphere fluids by relating them to the reversible work, i. e. the change in free energy, required to introduce an additional sphere to the system, which for our purpose shall have a radius, which is equivalent the thickness of the depletion zone δ . The work is calculated by expanding the radius of the sphere from zero to its final value, i. e. here the radius is $\lambda\delta$ with $0 \le \lambda \le 1$. In limit of very small λ it is assumed that there is no overlap of depletion zones, thus the free volume fraction is

$$\alpha = \frac{V - N_c 4\pi (R_c + \lambda \delta)^3 / 3}{V} \tag{11}$$

which together with eq.10 results

$$\lim_{\lambda \to 0} W(\lambda) = -k_B T \ln\left(1 - \rho_c \frac{4\pi}{3} (R_c + \lambda \delta)^3\right),\tag{12}$$

where ρ_c is the colloid number density in the system. In the other extreme, for very large depletants the work is approximately the volume work needed to create a cavity which is just large enough to accommodate the particle at the given pressure, thus

$$W = \frac{4\pi}{3} (\lambda \delta)^3 \Pi.$$
(13)

In the scaled particle theory the work is calculated by expanding the $\lambda \to 0$ limit as a series in λ and adding the value for the large depletant case

$$W(\delta) = W(\lambda \to 0) + \frac{\partial W(\lambda \to 0)}{\partial \lambda} \lambda + \frac{1}{2} \frac{\partial^2 W(\lambda \to 0)}{\partial \lambda^2} \lambda^2 + \frac{4\pi}{3} (\lambda \delta)^3 \Pi.$$
(14)

After some tedious algebra a complicated but analytical expression for α in dependence of q and the colloid volume fraction ϕ_c is obtained

$$\alpha = (1 - \phi_c) \exp\left\{-U(\phi_c)\right\} \tag{15}$$

with

$$U(\phi_c) = aQ + bQ^2 + cQ^3$$

$$Q = \phi_c / (1 - \phi_c)$$

$$a = 3q + 3q^2 + q^3$$

$$b = 9q^2 / 2 + 3q^3$$

$$c = 3q^3$$
(16)

This is used to calculate the semi-grand potential, the chemical potential of the colloids and the pressure in the system.

3 Crowding and colloidal phase behavior

In soft matter physics, it is well known that polymer solutions, colloidal suspensions as well as combinations of both may undergo phase transitions of various types, if the solute content is sufficiently increased. In the field of structural biology this effect was long used to achieve protein crystallization for their characterization by x-ray diffraction [29, 30]. The most prominent case of such a phase transition in the field of biology is probably the separation into a protein rich and a protein poor phase in the eye lens cytoplasm, as this is the cause for cataract formation [39, 40, 41, 25]. In cell biology it is speculated that liquid/liquid phase separation is the mechanism by which compartments coexist in the nuclei of eucaryotic cells despite the absence of separating membranes [17, 18, 19]. Therefore, in this chapter we will discuss the fundamentals of phase transitions and coexistence in simple soft matter systems, to set the stage for the much more complex situation in biological environments.

3.1 Phase behavior of hard sphere suspensions

A suspension of spherical particles in a solvent, interacting solely via their excluded volume is termed a hard sphere fluid, if the particle number density, ρ_c , is low. In the very high dilution limit, the system may be regarded as the colloidal analogue to an atomic ideal gas, i. e. the osmotic pressure of the suspension is

$$\Pi = \rho_c k_B T \tag{17}$$

in analogy to the ideal gas law. Like in atomic systems, the particle volume and their interaction potential has to be taken into account to describe the pressure correctly at elevated concentrations. For hard sphere colloids, an accurate description of the osmotic pressure dependence on



Fig. 2: Left: Chemical potential in a hard sphere fluid (black) and in a hard sphere fcccrystalline phase in dependence of volume fraction, as calculated from eqs. 20 and 22. The volume fraction at the intersection is $\phi_c \approx 0.54$. Right: Osmotic pressure of a hard sphere fluid (black) and a hard sphere fcc-crystalline phase (red) in dependence of volume fraction. Full lines are calculated from eqs. 18 and 21 and symbols are from computer simulations [43] where full squares represent the pressure at coexistence [44].

particle volume fraction is given by the Carnahan-Starling equation of state

$$\widetilde{\Pi} = \frac{\phi_c + \phi_c^2 + \phi_c^3 - \phi_c^4}{(1 - \phi_c)^3},$$
(18)

where the dimensionless pressure is $\widetilde{\Pi} = \Pi v_S / k_B T$ with $v_S = 4\pi R_c^3/3$ the volume of a single sphere with Radius R_c . Again, in analogy to atomic systems we expect some kind of phase transition to occur, when the particle volume fraction becomes sufficiently high. At the transition, the pressure and the chemical potential μ of the two coexisting phases have to be equal. To find the transition for the hard sphere system we have thus to have an expression for its chemical potential. This can be obtained from the Gibbs-Duhem relation, which for constant temperature reads $d\Pi = \rho_c d\mu$. With $\rho_c = \phi_c / v_S$ the chemical potential can formally be written as

$$\int d\mu = \mu = \mu^0 + v_S \int_0^{\phi_c} \frac{1}{\phi'_c} \frac{d\Pi}{d\phi'_c} d\phi'_c$$
(19)

where the ideal gas reference term $\mu^0 = k_B T ln (\Lambda^3/v_S)$ with the particles De Broglie wavelength $\Lambda = h\sqrt{2\pi m_S k_B T}$, h Planck's constant and m_S the single particle mass [42]. Calculating $d\Pi/d\phi_c$ from eq. 18 and solving the integral on the right hand side, the dimensionless chemical potential $\tilde{\mu} = \mu/k_B T$ is obtained as:

$$\widetilde{\mu} = ln \frac{\Lambda^3}{v_S} + ln\phi_c + \frac{3 - \phi_c}{(1 - \phi_c)^3} - 3.$$
(20)

Theoretical and simulation studies in the mid of the twentieth century revealed that hard spheres fluid will transform into a face centered cubic ordered system at sufficiently large volume fractions. Analytical approximations for the pressure

$$\widetilde{\Pi}_{fcc} = \frac{3\phi_c}{1 - \phi_c/\phi_{cp}} \tag{21}$$



Fig. 3: *PMMA* particles suspended in an refractive index matching solvent mixture at effective volume fractions indicated by the numbers at the bottom. The samples are illuminated obliquely from behind with white light. Homogeneous colors indicate a random distribution of particles while opacity is due to Bragg diffraction, indicating crystal structures with lattice constants in the wave length range of visible light. Image reproduced from reference [45] with permission.

and the chemical potential

$$\widetilde{\mu}_{fcc} = ln \frac{\Lambda^3}{v_S} + \frac{27}{8\phi_{cp}^3} + 3ln \left(\frac{\phi_c}{1 - \phi_c/\phi_{cp}}\right) + \frac{3}{1 - \phi_c/\phi_{cp}}$$
(22)

, of this ordered system can be found in, e. g. reference [34] together with a comprehensive derivation of the entire subject. Here $\phi_{cp} \approx 0.74$ is the volume fraction at crystalline close packing of spheres .

As shown in the left panel of Fig. 2, the two curves representing eqs. 20 and 22 intersect at $\phi_c \approx 0.54$. This value is indicating the highest volume fraction up to which the fluid state of the hard sphere system can exist, since the condition for coexisting phases is $\tilde{\mu} = \tilde{\mu}_{fcc}$. At this volume fraction, the pressure of the ordered system has its minimal, i. e. the coexistence value $\Pi_{coex} \approx 6.0$. According to eq. 18 the fluid system reaches this pressure at a volume fraction of $\phi_c \approx 0.49$ as shown in the right panel of Fig. 2, where we compare the calculated data according to eqs. 20 and 22 to results of computer simulations [44, 43].

An experimental verification of the fluid to crystal transition in a quasi hard sphere colloidal system is displayed in Fig. 3, which shows an image of samples containing increasing volume fractions of colloidal spheres in an index matching solvent mixture [45]. The samples are illuminated obliquely from behind with white light, which results in an opaque appearance due to Bragg diffraction, if there is long range order. Samples or parts of the samples showing homogeneous color have a fluid like structure. Thus, below an effective volume fraction of ~ 0.5 the sample is completely isotropic, while in the range $0.51 < \phi < 0.53$ a fluid coexists with a crystalline structure. Up to $\phi \leq 0.58$ the samples are fully crystalline while at even higher volume fractions a new isotropic state occurs, which is a colloidal glass.

It is important to note that even the apparently simple system of hard colloidal spheres undergoes phase separation at sufficiently high concentrations. However, there is no phase coexistence of a gas-like and a liquid-like phase, thus the hard sphere fluid is always super-critical.



Fig. 4: Sketch of the semi-grand potential in dependence of the colloid volume fraction in a colloid-polymer mixture. The common tangents connect points of equal slopes, identifying the the colloid volume fraction of coexisting phases. Left: gas-like/liquid-like coexistence; middle: three phase coexistence of a gas-like, a liquid-like and a solid-like state; right: coexistence of a supercritical fluid state with a crystalline state.

To induce a gas/liquid-like phase coexistence, attractive interactions between the particles are required, which we will discuss in the next section.

3.2 Depletion interaction and colloidal phase behavior

In crowded systems, particle-particle interactions will become more and more important the higher the macromolecular contents. In experimental model systems, electrostatic interactions can be reduced by adding large amounts of electrolytes and van der Waals interactions can be minimized by adjusting the solvent's refractive index to that of the suspended particles. Differently, depletion interactions will inevitably occur, due to excluded volume interactions. To predict the phase diagram of such systems, the semi-grand potentials, Ω , of the various possible states of the system have to be calculated numerically, using free volume theory. Since the chemical potential of a component is the derivative of the semi-grand potential with respect to its concentration, the slope of a Ω vs. ϕ curve is a chemical potential. This is illustrated in Fig. 4 where we sketch the semi-grand potential in dependence of the colloid volume fraction ϕ_c in a colloid polymer mixture for three different situations. To identify the colloid content of coexisting phases, conditions have to be identified by the common tangent (the thin, non vertical lines in the graphs) construction, connecting the points of the curves with equal slope, indicating identical chemical potentials of the colloid $\mu_c^T = \mu_c^{TI}$ in the coexisting phases.

According to standard thermodynamics $\Omega = -\Pi V + \mu N$. Thus, the system pressure at coexistence is related to the intercept of the common tangent. In the left panel of Fig. 5 we show a phase diagram in the $\phi_d^R - \phi_c$ plane of a polymer-colloid mixture at q = 0.6, where ϕ_d^R is the depletant polymer volume fraction in a reservoir. The plot represents a comparison of results calculated using free volume theory [36] with data from computer simulations [46]. At volume fractions of polymer and colloid, which are both below ~ 0.5 the mixture consists of a single liquid-like phase. If the polymer volume fraction is $0.5 \leq \phi_d^R \leq 1.3$ the system separates into a gas-like phase, coexisting with a liquid-like phase and at polymer volume fractions $\phi_d^R \gtrsim 1.3$ a supercritical fluid is coexisting with the crystalline phase. For colloid volume fractions $\phi_c \gtrsim 0.55$ the system is in the crystalline state, independent of the polymer content. A triple point, at which the three phases coexist is located at $\phi_c \approx 0.5$ and $\phi_d^R \approx 1.3$.

An experimental example for this behavior is displayed in the right panel of Fig. 5. The samples shown, contain polystyrene latex particles with $R_c = 67$ nm at a volume fraction $\phi_c = 0.16$ and



Fig. 5: Left: Phase diagram of a colloid polymer mixture in the $\phi_p^R - \phi_c$ -plane at q = 0.6. Full lines were calculated using free volume theory [36] and symbols are data from computer simulations (reproduced with permission from reference [46]). Right: Image of 67 nm polystyrene latex aqueous suspensions with different amounts of hydroxyethylcellulose (HEC), indicated as weight per volume percentage by the numbers at the bottom. The sample with 0.3 %w/v shows the coexistence of a solid, a liquid-like and a gas-like phase (reproduced with permission from reference [47]).

hexaethylcellulose (HEC) as the polymeric depletant at concentrations which are indicated in units of % (w/v). The sample containing 0.3% HEC, obviously has the triple point composition, as three phases are clearly visible. The top phase boundary remains horizontal, if the sample is tilted, showing that the two adjacent phases are fluid. Differently, the lower phase boundary follows the inclination of the vial, showing that the bottom phase has solid-like properties.

4 Reaction equilibria and kinetics

4.1 Equilibrium constants

Most proteins and other biopolymers in living cells assemble into compound structures or pass through transient complexes, to execute their designed function. It is obvious that attractive interactions among the constituents will favor the development of these complexes. However, also excluded volume effects my have a pronounced effect on their formation. It is difficult to measure these unspecific contributions directly, but their effect can be observed indirectly by determining equilibrium constants [48, 49, 50, 51, 52] and reaction rates [53, 54, 55, 56]. Many of these effects can be predicted quantitatively, using thermodynamic approaches. Here we will discuss the change of the equilibrium constant for the association of the two species A and B into a complex AB. In Fig. 6 we show a thermodynamic cycle, which basically consist of four steps (i) the association/dissociation reaction in dilute solution, (ii) the transfer of the product from the dilute system into a crowded environment, (iii) the association/dissociation



Fig. 6: Thermodynamic cycle demonstrating the difference between the standard reaction free energy, ΔG_{AB}^0 , in dilute solution and the free energy change, ΔG_{AB} , of the same reaction in a crowded environment. The difference is determined by the changes of free energy ΔG_X^{crowd} , connected to transferring the species $X \in A, B, AB$ from a dilute solution to a crowded suspension.

reaction in the crowded environment and (iv) the transfer of the reactants from the crowded environment back into the dilute solutions. Since the change of Gibbs free energy along this circular path has to be zero, we realize that the difference between the reaction free energy in the crowded system and the standard reaction free energy is

$$\Delta\Delta G_{AB} \equiv \Delta G_{AB} - \Delta G_{AB}^{0} = \Delta G_{AB}^{crowd} - \Delta G_{A}^{crowd} - \Delta G_{B}^{crowd}$$
(23)

where ΔG_X^{crowd} is the molar work it takes to transfer the species X, which comprises A, B and AB, from a dilute solution to the crowded environment. Approximating the reactants and the formed complex as convex spheroid bodies, we can calculate this work applying results from scaled particle theory [38, 57, 58]. The work required to accommodate a spherocylinder with a tip-to-tip length L and a cross section radius R_{cs} , which is equal to the radius of the hemispherical end-caps, into a suspension of hard spherical particles with radius R_1 at volume fraction ϕ is approximately given by

$$\frac{\Delta G_X^{crowd}}{RT} = -\ln(1-\phi) + A_1Q + A_2Q^2 + A_3Q^3 \tag{24}$$

where R is the gas constant, $Q = \phi/(1 - \phi)$ is the ratio of the volume occupied by the particles over the unoccupied volume and

$$A_{1} = R_{cs}^{3} + 3R_{cs}^{2} + 3R_{cs} + \frac{3}{2}(l+1)(R_{cs}^{2} + 2R_{cs} + 1)$$

$$A_{2} = \frac{3}{2}(2R_{cs}^{3} + 3R_{cs}^{2}) + \frac{9}{2}(l+1)(R_{cs}^{2} + R_{cs})$$

$$A_{3} = 3R_{cs}^{3} + \frac{9}{2}(l+1)R_{cs}^{2}.$$
(25)





Fig. 7: Variation of the equilibrium constant for the formation of a spherocylindrical complex from two spherical reactants in dependence of volume fraction of spherical crowders with same radius as sphere A. Different colors of the full lines refer to various aspect ratios of the product as indicated in the panel at the right.

Here *l* is the length of the right cylinder, i. e. the spherocylinder without the end-caps, expressed in units of the cross section radius $l = L/2R_{cs} - 1$, thus l = 0 represents the case of a sphere with radius R_{cs} .

Let us now consider an association reaction between two particles, one with radius $R_A = R_1$ the other with radius $R_B = rR_1$ into a spherocylindrical object with $R_{cs} = R_B$ and various length, demanding that the volume of the spherocylinder shall be equal to the sum of the volumes of the two reactant spheres which requires that $l = 2/3r^3$. The equilibrium constants of the reaction in dilute solution, K_{AB}^0 , and in the crowded environment, K_{AB} , and their ratio are given by standard chemical thermodynamics as

$$K_{AB}^{0} = \exp\left\{-\Delta G_{AB}^{0}/RT\right\}$$

$$K_{AB} = \exp\left\{-\Delta G_{AB}/RT\right\}$$

$$\frac{K_{AB}}{K_{AB}^{0}} = \exp\left\{-\Delta \Delta G_{AB}/RT\right\}$$
(26)

Result for the ratio of equilibrium constants are shown in Fig. 7. Obviously crowding can increase the tendency of species A and B to associate, when the resulting complex is rather compact, i. e. $l \rightarrow 0$. On the other hand crowding may hamper association, if the dimer becomes increasingly elongated. In this case the resulting complex excludes more volume to the crowder than the two reactant spheres. For systems in which a larger number, N, of proteins associate, similar effects are predicted, which increase drastically with N [59].

4.2 Reaction kinetics

The kinetics of association reactions are dominated mainly by two effects, the diffusion of reactants which controls the probability of two reactants meeting and the energy barrier related to the formation of intermediate reaction states. According to the two extreme cases, the reaction rate in diffusion limited cases and in transition-state controlled situations is here called



Fig. 8: Left: Effect of PEG chain length on the rate constant of the BLIB dimerization. Rate constants, normalized by the rate constant in pure buffer, are plotted versus the solvent viscosity, normalized by the viscosity of the pure buffer, for solutions containing PEG with different chain length, N, as indicated in the legend. Points represent experimental data reproduced from reference [53] and lines are guides to the eye. Right: Sample turbidity as a function of time following the polymerization reaction of HIV-1 capsid protein in suspensions with various contents c_F of the synthetic crowding agent Ficoll 70. Symbols represent data reproduced from reference [54] and the lines are best fits with the Hill-function [62].

 k_d and k_t , respectively. In general the reaction rate of the association can be approximated by a a combination of the two limiting cases [60] as

$$k_a = \frac{k_d k_t}{k_D + k_t} \tag{27}$$

In the two limiting cases, the rate constants depend in different ways on crowding. As k_d is proportional to the reactants mobility in its surrounding, increasing volume fractions of crowders are expected to cause a monotonic decrease of k_d . At the same time depletion causes an effective attraction between the reactants, favoring dimerization, which may partially counterbalance but generally not overcompensate the effect of decreasing k_d . It is argued that k_t should increase with crowding, because crowding will favor the formation of compact transition states thereby lowering the energy barrier for their formation. Thus the association reaction. As a rule of thumb, fast associations are usually diffusion limited, while slow associations are dictated by the energy barrier of the transition state [61]. Therefore, crowding is generally expected to slow down fast and to accelerate slow reactions.

Examples for both limits are shown in Fig. 8. In the left panel the effect of buffer viscosity and crowding on the rate constant for the dimerization of β -lactamase inhibitor protein (BLIB) is demonstrated. Rate constants, normalized by the rate constant in pure buffer are plotted versus the normalized viscosity of the suspending solution, which contains poly-ethylene-glycol (PEG) with different degrees of polymerization, N. In ethylene-glycol (N = 1) and in PEG-200 ($N \approx 3$) solution, the reaction rates are significantly decreased due to increasing viscosity of the solvent, which results in a corresponding reduction of the BLIP diffusion constant. However, in PEG-8000 solution ($N \approx 130$) the polymer coils are large enough to cause significant attractive depletion interaction between the proteins, favoring association and thereby counterbalancing the effect of decreasing k_d to a large extent. In the right panel the optical density of suspensions of Human Immunodeficiency Virus Type 1 Capsid Protein (HIV-1 CA) is plotted as a function of time after triggering their assembly reaction by the addition of the crowding agent Ficoll70, which is a highly branched copolymer of saccharose and epichlorohydrin. As the total content of HIV-1 CA in the sample is constant over time, the turbidity is a rough measure for the average size of the assemblies in the solution. Either the inverse of the time t_{50} , at which the optical density reaches half of its final value OD = 0.5ODF, or the rate of change of the optical density k_{OD} in the region where OD depends more or less linear on time, is a measure for the rate constant of the assembly reaction. It is thus obvious that the reaction is significantly sped up by the presence of the crowder, which is just occupying volume and increasing the system viscosity, being inert in any other respect.

5 Conclusions

In may biological systems, like blood, cytoplasm or the nuclei of eucaryotic cells there is a large volume fraction of macromolecular solutes, inevitably causing excluded volume interactions. These can drastically alter the structure and the functioning of biological macromolecules, as compared to their properties in dilute suspensions.

Here we have discussed that excluded volume effects will lead to colloidal phase transition and the coexistence of various phases at volume fractions which are typical for many biological system. It is thus speculated that phase separation, due to macromolecular crowding is the basis for micro-compartmentation in cytoplasm and nucleoplasm.

Further protein function is strongly influenced by macromolecular crowding, by e. g. fostering reactions, leading to the reduction of total excluded volume. However, the magnitude of the effect is expected to dependent drastically on the relative size and shape of the involved species. Finally we demonstrated that reaction rates are influenced by crowding effects. Generally we may expect that fast, diffusion limited reactions are slowed down by crowding while slow, transition state- controlled reactions are sped up.

References

- [1] G. B. Ralston, J. Chem. Edu. 67(10), 857 (1990).
- [2] A. P. Minton, Biophysical Journal **31**(1), 77 (1980).
- [3] A. P. Minton, Biochemistry 20, 4821 (1981).
- [4] R. J. Ellis, Trends in biochemical sciences 26(10), 597 (2001).
- [5] R. J. Ellis, Current opinion in structural biology 11(1), 114 (2001).
- [6] G. Foffi, Phys. Biol. 10, 040301 (2013).
- [7] A. P. Minton, Biophys. J. 78, 101 (2000).
- [8] A. P. Minton, Curr.Opin. Struct. Biol. 10, 34 (2000).
- [9] A. P. Minton, Biophys. J. 88, 971 (2005).
- [10] M. S. Cheung, D. Klimov, and D. Thirumalai, Proc. Natl. Acad. Sci. 102, 4753 (2005).
- [11] M. S. Cheung, Current Opin. Struct. Biol. 23, 1 (2013).
- [12] T. Hoppe and J.-M. Yuan, J. Phys. Chem. B 115, 2006 (2011).
- [13] A. P. Schlesinger, Y. Wang, X. Tadeo, O. Millet, and G. J. Pielak, Journal of the American Chemical Society 133(21), 8082 (2011).
- [14] E. Chen, A. Christiansen, Q. Wang, M. S. Cheung, D. S. Kliger, and P. Wittung-Stafshede, Biochemistry 51(49), 9836 (2012).

- [15] Y. Wang, M. Sarkar, A. E. Smith, A. S. Krois, and G. J. Pielak, Journal of the American Chemical Society 134(40), 16614 (2012).
- [16] L. A. Benton, A. E. Smith, G. B. Young, and G. J. Pielak, Biochemistry 51(49), 97739775 (2012).
- [17] H. Walter and D. E. Brooks, FEBS letters 361(2-3), 135 (1995).
- [18] R. Hancock, Biol. Cell 96, 595 (2004).
- [19] W. M. Aumiller Jr, B. W. Davis, and C. D. Keating, Int. Rev. Cell Mol. Biol 307, 109 (2013).
- [20] H.-X. Zhou, Journal of Molecular Recognition 17(5), 368 (2004).
- [21] H.-X. Zhou, G. Rivas, and A. P. Minton, Annu. Rev. Biophys. 37, 375 (2008).
- [22] M. Vis, J. Opdam, I. S. J. van t Oor, G. Soligno, R. van Roij, R. H. Tromp, and B. H. Ern, ACS Macro Letters 4(9), 965 (2015).
- [23] G. B. Benedek, Appl. Opt. 10, 459 (1971).
- [24] M. Delaye and A. Tardieu, Nature 302, 415 (1986).
- [25] G. Foffi, G. Savin, S. Bucciarelli, N. Dorsaz, G. M. Thurston, A. Stradner, and P. Schurtenberger, Proceedings of the National Academy of Sciences 111(47), 16748 (2014).
- [26] S. Asakura and F. Oosawa, J. Chem. Phys. 22, 1255 (1954).
- [27] S. Asakura and F. Oosawa, J. Polym. Sci 33(126), 183 (1958).
- [28] A. Vrij, Pure Appl. Chem. 48, 471 (1976).
- [29] S. D. Durbin and G. Feher, Annu. Rev. of Phys. Chem. 47, 171 (1996).
- [30] J. Blouwolff and F. Seth, Journal of Crystal Growth 303, 546 (2007).
- [31] E. Zhou, X. Trepat, C. Park, G. Lenormand, M. Oliver, S. Mijailovich, C. Hardin, D. Weitz, J. Butler, and J. Fredberg, Proceedings of the National Academy of Sciences 106(26), 10632 (2009).
- [32] M. A. Mourão, J. B. Hakim, and S. Schnell, Biophysical journal 107(12), 2761 (2014).
- [33] B. R. Parry, I. V. Surovtsev, M. T. Cabeen, C. S. OHern, E. R. Dufresne, and C. Jacobs-Wagner, Cell 156(1), 183 (2014).
- [34] H. N. W. Lekkerkerker and R. Tuinier, *Colloids and the Depletion Interaction*, vol. 833 of *Lecture Notes in Physics* (Springer, Dordrecht, Heidelberg, London, New York, 2011).
- [35] P. R. Lang, D. Vlassopoulos, and W. Richtering, *Polymer Science: A Comprehensive Reference* (Elsevier, Amsterdam, The Netherlands, 2012), chap. Polymer/Colloid Interactions and Soft Polymer Colloids.
- [36] H. N. W. Lekkerkerker, W. C.-K. Poon, P. N. Pusey, A. Stroobants, and P. B. Warren, EPL (Europhysics Letters) 20(6), 559 (1992).
- [37] G. J. Fleer and R. Tuinier, Advances in Colloid and Interface Science 143(1-2), 1 (2008).
- [38] H. Reiss, H. L. Frisch, and J. L. Lebowitz, The Journal of Chemical Physics 31, 369 (1959).
- [39] M. L. Broide, C. R. Berland, J. Pande, O. O. Ogun, and G. B. Benedek, Proceedings of the National Academy of Sciences 88(13), 5660 (1991).
- [40] C. Liu, N. Asherie, A. Lomakin, J. Pande, O. Ogun, and G. B. Benedek, Proceedings of the National Academy of Sciences 93(1), 377 (1996).
- [41] N. Dorsaz, G. M. Thurston, A. Stradner, P. Schurtenberger, and G. Foffi, Soft Matter 7, 1763 (2011).
- [42] D. A. Mc Quarrie, *Statistical Mechanics* (University Science Books, Sausalito, 2000).
- [43] A. Fortini, M. Dijkstra, and R. Tuinier, J. Phys.: Condens. Matter 17(50), 77837803 (2005).
- [44] W. G. Hoover and F. H. Ree, J. Chem. Phys. 49, 3609 (1968).

- [45] P. N. Pusey and W. van Megen, Nature 320, 340 (1986).
- [46] M. Dijkstra, R. van Roij, R. Roth, and A. Fortini, Phys. Rev. E 73, 041404 (2006).
- [47] M. A. Faers and P. F. Luckham, Langmuir 13(11), 2922 (1997).
- [48] H. J. Bosma, G. Voordouw, A. De Kok, and C. Veeger, FEBS Letters 120, 179 (1980).
- [49] J. Wilf and A. P. Minton, Biochimica et Biophysica Acta (BBA) Protein Structure 670(3), 316 (1981).
- [50] S. B. Zimmerman and S. O. Trach, Nucleic Acids Research 16(14A), 6309 (188).
- [51] T. Díaz-López, C. Dávila-Fajardo, F. Blaesing, M. P. Lillo, and R. Giraldo, Journal of Molecular Biology 364(5), 909 (2006).
- [52] X. Aguilar, C. F. Weise, T. Sparrman, M. Wolf-Watz, and P. Wittung-Stafshede, Biochemistry 50(14), 3034 (2011).
- [53] N. Kozer and G. Schreiber, Journal of Molecular Biology 336(3), 763 (2004).
- [54] M. del Álamo, G. Rivas, and M. G. Mateu, J. Virol. 79(22), 14271 (2005).
- [55] Y. Phillip, E. Sherman, G. Haran, and G. Schreiber, Biophysical Journal **97**(3), 875 (2009).
- [56] Y. Phillip and G. Schreiber, FEBS Letters 587(8), 1046 (2013).
- [57] M. A. Cotter, J. Chem. Phys. 66, 1098 (1977).
- [58] G. Rivas, J. A. Fernández, and A. P. Minton, Proceedings of the National Academy of Sciences 98(6), 3150 (2001).
- [59] A. P. Minton, Biopolymers **20**(10), 2093 (1981).
- [60] H. Zhou and A. Szabo, Biophysical Journal 71(5), 2440 (1996).
- [61] R. Alsallaq and H.-X. Zhou, Biophysical Journal 92(5), 1486 (2007).
- [62] A. V. Hill, J. Physiol. 40, 4 (1910).
D3 Neuronal signaling

F. Müller Cellular Biophysics Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

2	Basi	ic properties of neurons	2
	2.1	Neurons are the functional units of the nervous system	2
	2.2	Neurons generate membrane potentials	3
3	Neu	ronal signaling	6
	3.1	Local signaling depends on passive membrane properties	6
	3.2	Long-distance signaling is mediated by action potentials	9
	3.3	Neurons show complex electrophysiological features	11
	3.4	Information processing arises from interaction in neuronal networks	12

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

In every instant of our life, our brain manages incredibly complex tasks. It does this so efficiently, that we take its performance as granted and hardly ever think about the mechanisms that determine how we perceive, behave, think, and remember – or, in other words, how the brain generates our individuality. Our brain is extremely good at acquiring and processing information. While such information must be processed within milliseconds, e.g. to enable rapid action in a dangerous situation, the brain is also capable of storing information in form of memories for an entire lifetime. How do neurons manage these tasks?

In our nervous system, individual neurons process and propagate information in form of electrical voltages at their membrane. At the connection between two neurons, the synapse, the electrical signal is translated into a chemical signal by the release of small molecules, socalled neurotransmitters. (These chemical synapses make up the majority of synaptic connections in the nervous system. Other synapses work by direct electrical coupling of neurons. See Chapter C7 for details.) These neurotransmitters on the other hand, generate electrical signals in the subsequent neuron. Hence, when information is exchanged between neurons within a neuronal network, it is permanently translated from electrical into chemical signals and back. Understanding the molecular mechanisms that underly neuronal signaling is essential for the understanding of high-level function, such as cognitive processes or memory formation. On the other hand, the cellular and molecular mechanisms that provide neurons with their signaling abilities are subjected to disease processes and are targets for toxins that result in malfunction of the nervous system. Knowledge of the cellular and molecular processes in neurons is, therefore, fundamental to understanding the brain in health and disease. In this chapter, we will focus on the generation of electrical signals. For further reading we refer to textbooks on neurobiology and electrophysiology [1], [2], [3].

2 Basic properties of neurons

2.1 Neurons are the functional units of the nervous system

In the first half of the nineteenth century, the cell was recognized as the fundamental unit of all living organisms. Neurons share the repertoire of organelles present in our body cells, like the nucleus, ribosomes, endoplasmic reticulum and a plasma membrane that encloses the cell and separates it from its environment. What distinguishes neurons from all other cells is their specialization for signaling and processing of information and for intercellular communication. This specialization is apparent in the elaborate morphology of neurons, in the complex properties of their plasma membrane and in the specific patterns of connection with other neurons in form of neuronal networks (for detailed aspects of network function, see Chapter E7). Neurons come in an astonishing diversity (for examples see Fig. 1). There are probably several hundred types of neurons found in our body. In a "typical" neuron, several processes emerge from the cell body. They can be distinguished into dendrites (usually numerous processes) and an axon (one long process). The dendritic arborization can be quite large and elaborate. Morphological classification of neurons mostly relies on the architecture of the dendritic tree. Dendrites are the site of synaptic input from other neurons. This input is integrated by the cell and finally relayed to other neurons via the axon. The axonal ending is the site of synaptic output to other neurons. In our body, the length of axons varies vary between few µm and one m.

The human brain harbors some 10^{11} neurons with appr. 10^{14} synaptic connections. Although the brain constitutes only 2% or our body weight, it consumes 20% of our metabolic energy.



Fig. 1: The morphology of neurons can differ substantially. Note the different scale bars. A, drawing of Purkinje Cells from the cerebellum, stained with the Golgi Method by Santiago Ramon y Cajal, around 1900. All dendrites (top) stratify in one plane, axons project towards the bottom. B, retinal ganglion cells, filled with fluorescent dyes. The dendritic tree is radially symmetric around the soma, all axons project to the optic nerve (towards top right corner). C, retinal bipolar cells as examples of locally acting interneurons (cells are genetically modified to express a green fluorescent protein). Note the small dendritic tree (top) and the short axon with a well-established axon terminal system.

2.2 Neurons generate membrane potentials

Neurons employ the electrical voltage across their plasma membrane to encode and propagate information. This voltage can be measured by impaling the cell with a microelectrode. The membrane voltage V_m is the difference between the intracellular electrical potential V_i and the extracellular potential V_o . As V_o is defined as zero, V_m equals V_i . As soon as the microelectrode penetrates the cell, a resting membrane potential between -40 and -90 mV can be measured, depending on the type of neuron. As the interior of the cell is negatively charged, an electrophysiologist might also say that the membrane is polarized. When a neuron receives and processes information, the membrane potential changes. We call it hyperpolarization, if the membrane potential becomes more negative, and depolarization, if the membrane potential becomes more negative, and depolarization, if the membrane activation or excitation, while hyperpolarization means inhibition.

The generation of electrical potentials across the plasma membrane rests on three key features: first, the lipid bilayer of the membrane acts like an insulator, separating the two compartments on both sides of the membrane; second, there is a difference in concentration of certain ions across the membrane; third, the membrane displays a selective permeability for some of these ions. The latter two features are not observed in a "naked" lipid bilayer. Rather, they are endowed to the plasma membrane by two types of membrane spanning proteins. Transporters actively move ions into or out of the cell against their concentration gradient. Ion channels allow the flow of ions in the direction of their electrochemical gradient. The structure and function of transporters and ion channels is described in Chapter B4. Let us briefly summarize the properties of an ion channel. An ion channel is usually formed by several proteins, that span the membrane. They make up an aequous pore, that is usually closed by a gate. Upon activation of the channel, the gate opens due to a conformational

change in the protein and allows the flux of ions until the channel pore closes again. The pore selects one or few ion species against the other ions, providing the channel with a certain ionic selectivity. Our genome encodes several hundred kinds of ion channels.

As described in Chapter B4, the Na^+-K^+ -ATPase utilizes energy derived from the cleavage of ATP, the universal energy coin of the cell, to pump Na^+ ions out of the cell into the extracellular medium and K^+ ions into the cell. Due to the action of other transporters, concentration gradients are established for each physiologically relevant ion. The ion distribution in the mammalian – and, hence, in human - brain, is summed up in table 1.

	extracellular (mM)	intracellular (mM)	Nernst potential (mV)
Na ⁺	145	5 to 15	~+60
K ⁺	5	140	-85
Ca ²⁺	1 to 2	0.0001	+120
Cl	145	4 to 15	~ -65
Large anions	10	140	impermeable

Table 1. Concentrations of important ions in the extracellular and intracellularcompartments and their respective Nernst potential.

To understand how such ion distributions result in electrical potentials across the membrane, let us consider a simple hypothetical case (Fig. 2). Only two species of ions are present in this case: K^+ and Cl, and both are at equal concentrations inside and outside the cell. The membrane only harbours ion channels that are selectively permeable to K^+ . Let us further assume that these ion channels are spontaneously open most of the time and are, therefore, termed leak channels. Under these conditions, no electrical potential can be measured across the membrane, as electrical charges are identical on both sides. The number of K^+ ions that flow into the cell through the open channels is identical to the number of ions that flow in the opposite direction, hence the net flux of ions is zero. Now we increase the K^+ concentration inside the cell. This increases the chemical driving force and, therefore, will result in a K^+ efflux. As positively charged ions leave the cell, the cell interior becomes negative and an electrical potential is generated. The potential will grow until equilibrium is reached. At the electrochemical equilibrium, the two opposing forces are balanced: first, the chemical potential that supports K^+ efflux, and second, the electrical potential that impedes the efflux of K^+ ions. The equilibrium potential E_K can be calculated by the Nernst equation, where R is the gas constant, T is the absolute temperature (in Kelvin), z is the valence (i.e. the electrical charge, here +1) of the ion, and F is the Faraday constant. The determinative term in the equation is the ratio of the concentration of K^+ outside $[K^+]_0$ and inside $[K^+]_i$ the cell. The equation can be simplified for conditions at room temperature (eq. (2)).

$$E_{K} = \frac{RT}{zF} \ln \frac{[K^{+}]_{o}}{[K^{+}]_{i}}$$
(1)

$$E_{\kappa} = 58mV\log\frac{[K^+]_o}{[K^+]_i} \tag{2}$$



Fig. 2: Generation of the membrane potential in a simplified model. A, K^+ concentration is identical on both sides of the membrane, $V_m = 0$ mV, net flux through leak channels is zero. B, K^+ concentration inside is increased resulting in K^+ efflux and generation of a negative membrane potential. C, stable E_K has been achieved, net flux is zero.

In our simplified case, the Nernst equation would correctly predict the membrane potential of the cell. Interestingly, the fraction of intracellular K^+ ions that have to leave the cell to generate this potential is very small ($10^{12} K^+$ ions per cm² of membrane). It can be calculated by assuming that the plasma membrane behaves like a capacitor with a specific capacitance of $1 \mu F/cm^2$ membrane. Depending on the ratio of membrane surface to cell volume, one out of 10^4 to 10^5 intracellular K^+ ions must leave the cell to generate the membrane potential. In other words, the ionic composition of the cell does not change. Under physiological conditions, the Nernst potential for K^+ is close to -90 mV. Changing the K^+ concentrations on one side of the membrane by a factor of 10 would elicit a change in membrane potential by 58 mV (compare equation (2)).

In neurons, the situation is more complex as K^+ ions are not the only ion species that contribute to the membrane potential. This situation can be described using the Goldman equation (also termed Goldman-Hodgkin-Katz or GHK equation), in which K^+ , Na⁺, and Cl⁻ are the primary permeant ions.

$$V_{m} = \frac{RT}{F} \ln \frac{P_{K}[K^{+}]_{o} + P_{Na}[Na^{+}]_{o} + P_{Cl}[Cl^{-}]_{i}}{P_{K}[K^{+}]_{i} + P_{Na}[Na^{+}]_{i} + P_{Cl}[Cl^{-}]_{o}}$$
(3)

In this equation, P indicates the permeability of the membrane for the ion of interest. The permeability depends on the number of open ion channels that are permeable for this ion species. Let us first consider a simplified case in which all channels for Cl⁻ and Na⁺ ions are closed and, hence, P_{Na} and P_{Cl} are zero. In this case the Goldman equation simplifies to the Nernst equation for K⁺. If, on the other hand, the membrane is only permeable for Na⁺, the Goldman equation simplifies to the Nernst equation for Na⁺, yielding a membrane potential of appr. +60 mV. If both K⁺ and Na⁺ are permeable, V_m will achieve a value intermediate to the

two potentials. Under resting conditions, P_{Na} and P_{Cl} are not zero, but significantly smaller than P_K . Hence P_K dominates the membrane potential. The resting membrane potential of neurons is, therefore, often around -70 mV, i.e. close, but not identical to E_K .

3 Neuronal signaling

3.1 Local signaling depends on passive membrane properties

Let us now look at the processes, that happen when a neuron becomes activated. The simplified sensory neuron in Fig. 3 has a cell body and an axon that connects to another neuron via a synapse. The sensory compartment of the cell is specialized for the transduction of the sensory stimulus into an electrical signal. In many sensory cells, the external stimulus, e.g. an odorant, activates a cellular signaling cascade that often employs G-protein coupled receptors (for details, see Chapters B9 and F1) and eventually leads to the opening of Na⁺ permeable ion channels. There are two reasons why this leads to a strong Na⁺ influx: first, the Na⁺ concentration outside is much higher than inside the cell; second, the membrane potential favors the influx of Na⁺ ions. In other words, the electrochemical gradient for Na⁺ ions is large. The Na⁺ influx can be calculated following Ohm's law. The conductance g_{Na} of the membrane for Na⁺ is equal to $1/r_{Na}$, (the resistance of the membrane for Na⁺) and depends on the number of open channels, the electrochemical potential is given by V_m-E_{Na}.

$$I_{Na} = g_{Na}(V_m - E_{Na}) \tag{4}$$

The influx of Na⁺ ions depolarizes the cell, i.e. V_m becomes more positive. In a sensory neuron, we call this shift in membrane potential a "receptor potential". The generation of a receptor potential is one way to change V_m and the activity of a neuron and is found in all sensory cells. We can describe the depolarization using the Goldman equation: P_{Na} increases and V_m shifts towards E_{Na}. Because at the sensory compartment V_m has adopted a value different from the rest of the membrane, the local current produced by the opened channels will flow passively along the plasma membrane of the soma and the axon (with some decrement as we shall see later). Let us assume it reaches the presynaptic ending of the sensory neuron. The synapse consists of three parts: the presynaptic ending of our sensory neuron and the postsynaptic process of the following neuron, separated by the synaptic cleft (see Chapter C7). In the presynaptic ending, we find voltage-activated Ca^{2+} channels (see Chapter B4). These channels become activated, when V_m crosses a critical value called the threshold and conduct the influx of Ca²⁺ ions at the synaptic ending, again because the electrochemical gradient for Ca^{2+} ions is high. But what matters at the presynaptic site is the change in internal Ca^{2+} concentration rather than the amplitude of the ionic flux. In their function as intracellular messenger, once inside the cell the Ca^{2+} ions trigger the fusion of synaptic vesicles with the presynaptic membrane. These vesicles contain neurotransmitter molecules. During fusion of vesicle and plasma membrane, transmitter molecules are released into the synaptic cleft – the electrical signal has been translated into a chemical signal (for details, see Chapter C7). At the postsynaptic membrane, the transmitter molecules bind to ion channels. This in turn triggers a conformational change in the protein and subsequently opening of the pore (ligand-gated ion channels, see Chapter B4). If these ion channels are permeable to Na^{\dagger} , the resulting Na^{\dagger} influx will depolarize the postsynaptic cell (i.e. it will become excited). The translation of the chemical neurotransmitter signal into an electrical signal is another way to open Na⁺ channels and depolarize a neuron. We find it in all neurons. We call the depolarization of the postsynaptic cell a "synaptic potential", or in this case more specifically "excitatory postsynaptic potential" or EPSP.





In essence, Fig. 3 shows how neurons work. Opening of Na^+ permeable channels leads to depolarization of the membrane. The depolarization spreads. If the depolarization at the presynaptic membrane crosses the threshold for activation of voltage-gated Ca^{2+} channels, transmitter is released to the postsynaptic cell. The neuron "tells" the postsynaptic cell that it is active and relays its information to the postsynaptic cell. However, there are several problems that must be adressed.

First, the properties of the membrane affect the time course of electrical signaling. In principle, the membrane can be described by an equivalent circuit with two elements in parallel: a resistor with the resistance r_m and a capacitor with the capacitance c_m . The product of both parameters yields the time constant of the membrane, τ .

$$\tau = r_m c_m \tag{5}$$

If a rectangular current pulse is injected into the neuron in Fig. 4A, the change in membrane potential can be described by an exponential relationship:

$$V_t = V_{\infty} (1 - e^{-t/\tau})$$
 (6)

where V_{∞} is the steady-state value of V_m , e is the base of natural logarithms, t is the time, and τ is the time constant. After the end of the current pulse, V_m declines also exponentially :

$$V_t = V_\infty e^{-t/\tau} \tag{7}$$



Fig. 4: Passive properties of the membrane. *A*, current injection into the axon with an intracellular electrode. *B*, rise and decay in the membrane voltage changes exponentially with the time constant τ . *C*, change in membrane potential decays exponentially with distance from the site of current injection ($\lambda =$ length constant).

Second, compared to electric wires, neurons are poor conductors. The passive electrical properties of a neuron can be determined by measuring along the membrane the voltage change that is triggered by the initial depolarization. The amplitude of the voltage change decays exponentially with increasing distance from the site of the initial depolarization (Fig. 4C):

$$V_x = V_0 e^{-x/\lambda} \tag{8}$$

where V_x is the voltage at any distance x along the axon membrane, V_0 is the initial depolarization, and λ is the length constant of the axonal membrane. The length constant can be calculated as:

$$\lambda = \sqrt{ar_m/r_i} \tag{9}$$

The membrane resistance r_m (specific resistance of 1 cm² membrane, unit: Ω cm²) depends on the number of ion channels in the membrane. Current can dissipate through open channels, therefore, a low number of channels in the membrane will increase r_m and λ . The specific axial resistance r_i (Ω cm) is given for a cylindrical piece of axon of 1 cm length and crosssectional area of 1 cm². With a diameter a, r_i/a is the actual resistance along 1 cm length. Increasing the axon diameter increases λ . Here is the problem: remember that the neuron will only release transmitter to the postsynaptic neuron, if the depolarization at the presynaptic site crosses the threshold for activation of the voltage-activated Ca²⁺ channels. With λ values typically in the range from 0.1 to 1 mm, this only works for small neurons. In neurons with long axons projecting to distant areas in the brain, a passively conducted signal would die out on the way and would never reach the presynaptic site. In these cases, the signal must be boosted along the way by an active mechanism: the "action potential".

3.2 Long-distance signaling is mediated by action potentials

The action potential is a brief voltage pulse in which the electric polarization of the membrane reverses. It is based on the stereotypic consecutive activation of two kinds of voltage-gated ion channels: Na^+ channels and K^+ channels (for details see Chapter B4). To generate action potentials, the axonal membrane must be endowed with these types of channels. They are usually found in low density in the membrane of dendrites and soma, but at higher density in the axonal membrane. Action potentials are, therefore, usually initiated at the transition from cell body to the axon, the axon hillock. In mammals, the action potential lasts 1 - 2 ms or even less. Electrophysiologists say that a neuron "fires" action potentials. As described above, neurons can depolarize due to sensory or synaptic input. Once V_m crosses the threshold for voltage-gated Na⁺ channels, some of these channels open and conduct an inward Na⁺ current. This Na⁺ current causes further depolarization of the membrane, recruiting more Na⁺ channels etc. Based on this positive feedback loop V_m changes rapidly towards E_{Na} (as P_{Na} now dominates in the Goldman equation). V_m can reach values of +20 or +40 mV (overshoot). However, the depolarization leads not only to the opening of the Na^+ channels, but also subsequently to their inactivation. In the inactivated mode, the pore is plugged by a cytoplasmic loop of the channel protein and the channel does no longer conduct Na⁺. Hence, in the action potential, Na^+ influx is very transient. Note that the inactivated state differs from the closed state in that the channel cannot proceed to the open state anymore - it is locked in the inactivated state. At about the time Na⁺ channels inactivate, voltage-activated K⁺ channels open. As at depolarized V_m the electrochemical driving force for K^+ is very high, a strong efflux of K⁺ ions is initiated. As a consequence, V_m rapidly returns towards the resting value, a process called repolarization. Often, the action potential is followed by a short hyperpolarization, the after potential. As many K^+ channels have been opened during the action potential, K^+ efflux is higher than under resting conditions and V_m becomes more negative (i.e. closer to E_K). Once the resting membrane potential is reached, K⁺ channels close again and the Na⁺ channels can overcome the inactivation and proceed to the closed state. This switch only happens when V_m is close to the resting membrane potential. As this switch takes some time, a neuron that has just finished an action potential is not capable of generating another action potential for a short period of time (1 to few ms, depending on the cell). This time span is called the refractory phase. Action potentials are also called "spikes".



Fig. 5: Action potentials are formed by the consecutive activation of voltage-gated Na⁺ channels and K⁺ channels. A, resting membrane potential is dominated by K⁺ leak channels (yellow). B, depolarization by current through voltage-activated Na⁺ channels. C, repolarization by current through voltage-activated K⁺ channels (red), while Na⁺ channels are inactivated (note the green ball plugging the pore). D, after the refractory phase, voltage-gated Na⁺ channels and K⁺ channels proceed to the closed state.

As action potentials are generated by the stereotypic sequence of activation of Na⁺ and K⁺ channels, their amplitude is usually constant in a given neuron. Action potentials are either fired, if V_m crosses the threshold, or not fired, if V_m remains below the threshold (all-or-none rule). Once an action potential has been elicited at the axon hillock, it propagates along the entire axon without decrement. As the action potential recruites closed Na⁺ channels along the axonal membrane, it is always restored to its full amplitude. The speed of propagation depends on several aspects. As described above, a part of the local current generated by the activation of the Na^+ channels will spread along the axon, depolarizing the membrane in the adjacent section of the axon above the threshold and thereby eliciting an action potential at this site. As spreading depends on the length constant λ (see equation (9)), propagation speed is higher in axons with large diameter. Invertebrates like insects or molluscs, have developed axonal fibers with large diameter. The work of Alan Hodgkin and Andrew Huxley on the famous squid giant axon with a diameter of 1 mm (!) layed the foundation for our understanding of neuronal signaling and was awarded with the Nobel prize in 1963 (together with John Eccles). In vertebrates like mammals (and, hence, in man), the number of neurons in the nervous system is much higher than in invertebrates. Hence, the axon diameter must be limited due to space constraints. Here, nature has developed another means to increase propagation speed. Glial cells (supporting cells in the nervous system) ensheath the axonal fiber with stacks of membrane in form of large thin processes. This feature is called myelination. It is the myelin sheaths of axon tracts that give the white matter of the brain its appearance (as compared to the grey matter that harbors cell bodies). Myelination is functionally equivalent to increasing the thickness of the axonal membrane by as much as 100 times. Because the capacitance of a parallel-plate capacitor like the membrane is inversely proportional to the thickness of the insulation material, myelination strongly decreases c_m and, thus, the time constant τ (see equation (5)). A short time constant means that the membrane can be charged much more rapidly and, therefore, conduction is much faster in myelinated than in unmyelinated axons. However, if the entire axon were insulated by myelin, there would be no place for current flow needed to generate the action potential. Therefore, the myelin sheath is interrupted at regular intervals (every few hundred μm up to 2 mm). At these so-called nodes of Ranvier, the axonal membrane is exposed for a segment of $1 - 2 \mu m$ and Na^+ and K^+ channels are found at high density. The currents generated by these channels at each node spread passively but rapidly (due to the low c_m) within the myelinated segment until the next node is reached, where a new action potential is generated. Hence, excitation leaps from node to node, a process called saltatory. Conduction speed in unmyelinated axons ranges from 0.1 to 10 m/s. Myelinated axons can conduct at velocities of up to 120 m/s.

It is easy to conceive that loss of myelin as it occurs in diseases like multiple sclerosis, causes severe neurological problems. Signals may dissipate in unmyelinated segments and timing in information processing may be compromised due to changes in conduction velocity.

Saltatory conduction is also favorable from the standpoint of metabolic costs. Due to the unfavorable ratio of membrane surface to axonal volume, the fraction of K^+ ions that must leave the cell to generate the negative membrane potential is higher in thin axons than in large cell bodies. At high action potential rate, the Na⁺ and K⁺ gradients tend to run down in thin axons and gradients must be restored by the Na⁺-K⁺-ATPase. In myelinated axons, the ion gradient is only changed at a fraction of the axon length (i.e. at the nodes). Hence, compared to unmyelinated axons, less energy must be expended by the Na⁺-K⁺-ATPase in restoring the ion gradients.



Fig. 6: Saltatory action potential conduction along a myelinated axon. At t₁, Na⁺ influx through Na⁺ channels initiates an action potential at the left node of Ranvier. Local currents flow through the axon below the myelin sheath depolarizing the membrane. When V_m crosses the threshold at the next node at t₂, an action potential is generated there, while the membrane at the previously active node is refractory.

Please remember that the axonal membrane becomes refractory following an action potential, due to the inactivation of voltage-activated Na^+ channels. This prevents subsequent reexcitation of this segment of the membrane and, hence, backpropagation of action potentials. In other words, the refractory phase ensures the directed propagation of neuronal signals from the axon hillock towards the synaptic terminal at the end of the axon. Once the action potentials have reached the synaptic terminal, voltage-activated Ca^{2+} channels open and transmitter is released as discussed previously. While the action potential is often considered a hallmark of neurons, it should be stressed that not every neuron generates action potentials. As discussed above, passive conduction may be sufficient to trigger transmitter release in small neurons. These cells work locally in neuronal networks and are called interneurons. Often, they do not fire action potentials. They encode the stimulus intensity (or excitation) in graded potentials that passively spread throughout the cell. If action potential are involved in signal propagation, the stimulus strength must be encoded in the action potential frequency, as the action potential amplitude remains constant in a given cell. A strong stimulus triggers a burst of many action potentials, weaker stimuli elicit fewer action potentials.

3.3 Neurons show complex electrophysiological features

It is important to note that not all neurons are equal. There is no such thing as THE neuron. As much as neuronal cell types can differ in their morphology (see Fig. 1), they can differ in their electrophysiological signature. Figure 7 shows the response of three neuronal cell types to the same current injection (bottom trace). Each cell represents a particular neuronal type, that we call A, B, and C, respectively. All cells respond with depolarization to the current injection. While cell type A displays a continuous regular firing of action potentials (depicted as upward spikes at this resolution), in cell type B the timespan between successive action potentials becomes progressively longer, a process called adaptation. Cell type C responds with bursts of action potentials that are separated by long gaps of silence. Such cells behave like pacemakers. Why do different cell types respond so differently to the same stimulus? Because they differ in the repertoire of ion channels that are encoded by our genome. Ion

channels differ in their ionic selectivity, the way they are activated etc. For example, during the repetitive firing of action potentials, voltage-activated Ca^{2+} channels open and the intracellular Ca^{2+} concentration rises continuously. Some, but not all neuronal types express K^+ channels that become activated by intracellular Ca^{2+} . Such channels would tend to hyperpolarize the cell, making it more difficult to fire the next action potential. While cell type A does not express such channels, the cell types B and C do. If we take into account that a given neuron may express a dozen or more different types of ion channels, it is conceivable how the complex interaction between these channel type shapes the physiological properties of a given neuronal cell type. Hence, each neuronal type expresses a specific inventory of ion channel types necessary for its physiological function. This inventory is sometimes so characteristic that it can be used as an electrophysiological fingerprint to classify neurons [4].



Fig. 7: Different types of neurons can differ substantially in their electrophysiological signature. Upon depolarization, cell type A shows continuous firing, cell type B displays spike frequency adaptation, and cell type C fires in short bursts. Duration of current injection ca. 200 ms.

3.4 Information processing arises from interaction in neuronal networks

So far we have discussed, how neurons generate and propagate information in form of electrical potentials at their membrane. But how do neuronal networks process information? Each neuron may receive input from a large number of synapses (up to 100.000 per cell in extreme cases!). Please note that such synapses are not formed randomly between neurons, but very specifically only between certain types of neurons within a neuronal network. Each neuron "knows" to which neurons it must connect. For the most simple computation, you need two algebraic signs: + and -. In neuronal networks, this is provided by different types of synapses. A synapse can be excitatory (+) or inhibitory (-), depending on the type of neurotransmitter that is released by the presynaptic cell and on the type of postsynaptic receptor. Some transmitters, like glutamate or acetylcholine activate ion channels in the postsynaptic membrane, that conduct a Na⁺ inward current that depolarizes the postsynaptic cell. These synapses are excitatory, because they drive V_m towards a) the threshold for activation of Na⁺ channels and, hence, action potential generation, or b) the threshold for activation of Ca²⁺ channels and, hence, transmitter release. Inhibitory transmitters like glycine or γ -aminobutyric acid (GABA) open Cl⁻ conducting channels, that lead to Cl⁻ influx and, hence, hyperpolarization of the membrane (inhibitory synapse). In cells in which E_{Cl} is identical to the resting membrane potential, opening of Cl⁻ channels does not hyperpolarize the cell. However, a large number of open channels electrically shunts the membrane, resulting in a process called shunting inhibition that clamps V_m at the resting value and counteracts depolarization induced by excitatory input. Hence, inhibitory input makes it more difficult for the cell to reach the thresholds for Na^+ or Ca^{2+} channel activation.

Excitatory and inhibitory synapses and the postsynaptic potentials they induce form the basis for the computational processes that underly neuronal information processing. As synaptic potentials last for several ms, synaptic potentials can be summated even if they are generated at slightly different time points (temporal summation) or at different sites in the dendritic tree (spatial summation). The overall membrane potential of a cell is changed by the sum of its synaptic input. If the excitatory input prevails, V_m is depolarized. If this depolarization is sufficiently high to cross the thresholds discussed previously, action potentials can be elicited and transmitter can be released at the synaptic terminals of the cell. The neuron "speaks" to its neighbors and provides them with the information it has received. If the inhibitory input outweighs excitatory input, depolarization is impeded. The neuron is less active or even remains quiet.



Fig. 8: Computation in a neuronal network. The cell in the center receives excitatory input from three cells (+) and inhibitory input from one cell (-). Excitatory input from one cell triggers depolarization that remains below the threshold for the generation of action potentials. Spatial and temporal summation of three inputs increases depolarization that crosses the threshold and triggers action potentials to relay the information to the presynaptic ending. Inhibitory input leads to hyperpolarization that can counteract excitatory input.

By combining excitatory and inhibitory input in well organized neuronal networks, our brain can process information in a most sophisticated way. Important examples of directed inhibition are feedback inhibition and lateral inhibition. Feedback inhibition is a means to dampen signals and to make them more transient. For information processing, our brain favors transient signals over constant signals. Lateral inhibition is a widespread mechanism in the nervous systems. In the visual system, it forms the basis for contrast enhancement.

But synaptic connectivity is not the only key to understanding brain function. Our nervous system employs a plethora of neurotransmitters and receptors that can fine tune existing synapses and neuronal function. Whether we consider the activity of individual neurons, information processing in neuronal networks, cognitive functions or memory formation, all these processes are generated and shaped by the complex interaction of these elaborate signaling processes.

References

- [1] Principles of Neural Sciences. Eds: E.R. Kandel, J.H. Schwartz, T.M. Jessel, S.A. Siegelbaum, A.J. Hudspeth, 5th Edition, McGraw-Hill Education Ltd 2012
- [2] Neuroscience. Eds.: D. Purves, G.J. Augustine, W.C. Hall, A.S. LaMantia, L.E. White, 5th Edition, Sinauer 2012 (new edition will be available in 2018)
- [3] Ionic Channels of Excitable Membranes. B. Hille, 3rd Edition, Sinauer 2001
- [4] E. Ivanova and F. Müller, Vis Neurosci 23(2):143 (2006)

D4 Sequential Bottom-Up Assembly of Synthetic Cells

Ilia Platzman, Joachim P. Spatz

Department of Cellular Biophysics, Max Planck Institute for Medical Research, Jahnstr. 29, and the Department of Biophysical Chemistry, University of Heidelberg, 69120 Heidelberg, Germany

Contents

Introduction	•••••••••••••••••••••••••••••••••••••••	2
1	Formation, functionalization and characteriazation of microfluidicdroplets	3
	 Amphiphilic triblock- and gold-linked diblock-copolymer surfactants. 	4
	1.2 Microfluidic devices for droplets formation and their manipulation.	4
	1.3 Biofunctionalization of copolymer-stabilized droplets	5
2	Formation of droplets-stabilized GUVs (dsGUVs)	6
	2.1 Lipid concentration.	6
	2.2 Formation of dsGUV compartments by small unilamellar vesicles (SUVs) encapsulation.	6
	2.3 Formation of dsGUVs by means of pico-injection	7
	2.4 Fluorescence intensity analysis and FRAP measurements of dsGUVs.	8
3	Sequential bottom-up assembly of dsGUVs	9
-	3.1 Protein reconstitutions in dsGUVs	9
	3.2 Mobility of transmembrane proteins in dsGUVs	11
4	Approaches to release GUVs from dsGUVs 1	1
5	Summary and outlook for the	
3	future 1	13
References		4

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

Introduction

Compartmentalization, the formation of lipid-membrane compartments in which specific metabolic activity takes place, is one of the distinguishing features of eukaryotic cells.[1, 2] Synthetic biologists have concentrated on an attempt to develop cell-like compartments for biochemical reactions.[3, 4] The most common studies have focused on the development of enclosed cell-size volumes of aqueous space with organic lipid-based membranes as in giant unilamellar vesicles – GUVs, or polymer-based membranes as in polymersomes and polymer-stabilized emulsion droplets.

GUV-based compartments provide a suitable model system for mimicking the cell membrane. The lipid matrix with reconstituted proteins resembles the in vivo frame.[5] Therefore, free-standing GUVs have been utilized in various synthetic biology applications.[4, 6] However, the GUV-based compartment system has several drawbacks that are related mostly to its poor chemical and mechanical stability, and therefore limited potential for manipulation.

Due to their increased stability and lifetime, polymersomes made of amphiphilic block-copolymers in a continuous water phase are commonly used as an alternative in synthetic biology applications.[7] By adjusting the molecular properties of the block-copolymers, the thickness and properties such as the bending and stretching moduli of the membrane can be finely tuned.[8] Despite their increased mechanical and chemical stability, the encapsulation of biomolecules and further manipulation of such bio-containing polymersomes still represent big challenges due to the lack of technological means that allow efficient incorporation or loading of bio-ingredients for biochemical activity.

Compartments based on block copolymer-stabilized water-in-oil droplets have a potential to overcome the drawbacks related to technological limitations associated with "traditional" polymersomes. In the context of synthetic biology, droplet-based compartment systems possess the advantages of polymersomes and, in addition, can easily be integrated into microfluidic technologies for controlled and precise loading with biologically relevant materials.[9] Nevertheless the ability of these droplets or polymersomes to serve as optimal cell-like compartments is mainly hindered by their inability to mimic the biophysical properties of cellular membranes.[10] Combining the biophysical properties of cellular lipid membranes, the stability of copolymer-stabilized droplets and the ability to be adapted for high-throughput manipulation may be the solution.

These lecture notes describes our recently developed approach that merges lipid vesicle formation and droplet microfluidics for the generation of stable, uniformly sized and easily manipulated droplet-supported GUV (dsGUV) compartments (Fig. 1).[11] Contrary to conventional GUVs, the stability of the dsGUVs imparted by the supporting droplet interface is enhanced and their size can be precisely controlled. This novel system is poised to overcome the fundamental limitations associated with manipulation of currently employed cell-like compartments and possesses great potential for enabling efficient bottom-up assembly of minimal synthetic cells.



Fig. 1: (a) Schematic of a droplet-supported GUV compartment. The water-in-oil droplet is stabilized by triblock-copolymer surfactants. Optional copolymer surfactants are the gold-linked diblock-copolymers which can be employed to bio-functionalize the droplets' inner interface. The lipid bilayer is supported on the copolymer-stabilized oil-water interface of the droplet. (b) and (c) present chemical structures of PFPE-PEG-PFPE triblock- and PFPE-PEG-Gold diblock-copolymer surfactants, respectively [12].

1 Formation, functionalization and characteriazation of microfluidic droplets

Droplet-based microfluidics combines principles of science and technology, and enables the user to handle, process and manipulate droplets of very small volumes – down to less than a few picoliters – via microchannels. This technology permits the integration of multiple laboratory functions into one single microfabricated chip, requires minimal manual user intervention and sample consumption, and allows for enhanced data analysis speed and precision. The potential for application of this technology biological [12-14], chemical [15, 16] and medical [17] research is vast. This section will describe the preparation and characterization of copolymer-stabilized droplets required for the dsGUVs formation. The droplet dimension sets the size of the dsGUVs. Thus, droplet-based microfluidics offers (i) unprecedented control on the size and monodispersity of the resulting stable vesicles, a feature not yet available for conventional GUVs, and (ii) precise control and flexibility of the encapsulated content and the constituting membrane of the dsGUVs at high yield.

1.1 Amphiphilic triblock- and gold-linked diblock-copolymer surfactants

The most important key factor for the stability of water-in-oil emulsion droplets is the block copolymer surfactants at their interface. Nonionic fluorosurfactants made of perfluorinated polyether (PFPE) hydrophobic blocks attribute long-term stability to the droplets by preventing their coalescence, whereas polyethylene glycol (PEG) hydrophilic blocks serve as a biocompatible, inert droplet interface.[18] Optionally, to provide active sites for biochemical interactions within the droplets, a new type of surfactants that is covalently linked to gold nanoparticles (~ 5 nm in diameter) can be employed.[12] Such gold-linked

surfactants are mixed with gold-free surfactants (at mixing molar ratios ranging from 1:1000 to 1:2000) to create stable droplets functionalized with gold nanoparticles (Fig. 1).

1.2 Microfluidic devices for droplets formation and their manipulation

Simple agitation of an aqueous phase and the oil with dissolved copolymer surfactants will generate the formation of droplets. However, droplet-based microfluidic devices are required to create monodisperse droplets with the defined size. Moreover, microfluidic devices are neccesary for further droplets manipulation (e.g. injection, sorting and time-lapse analysis). All microfluidic devices used in this research are fabricated from poly(dimethylsiloxane) (PDMS) using photo- and soft-lithography methods.[19, 20] PDMS is a common material in microfluidic technology due to its low price, good biocompatibility and permeability to gasses, high transparency and low fluorescent background. [21]

Droplets are generated in a flow-focusing geometry junction, in which an aqueous phase is cut off by a surfactant-containing oil phase (Fig. 2 a, b). Following the formation, water-in-oil droplets are stabilized by accretion of block-copolymer surfactants at the water-oil interface leading to reduction of the oil/water interfacial tension[22] from 52.1 mN/m to a value of 19.5 mN/m for TRI7000 and 3.1 mN/m Tri2500.[23] The droplet diameter is mainly controlled by the channel dimensions, but can also be regulated to some extend by the variation of flow rates of the aqueous and oil phase.

To allow precise delivery of various biological components into preformed droplets, the microfluidic devices can be integrated with small and compact electrodes to apply electric fields in the microchannels. These electric fields induce destabilization (poration) of the surfactant (mono)layer and facilitate controlled injection (pico-injection) of aqueous phase into the droplets. The design of our droplet-based pico-injection unit is adapted from Abate et al.[9] A microfluidic flow control system was used to introduce droplets into the pico-injection unit. The spacing between the droplets was controlled through addition of oil with surfactants via the second oil channel as presented in Fig. 2c. Following the separation step, isolated droplets passed an electric alternating current (AC) field (frequency of 1 kHz, voltage of 250 V) generated by a signal generator and amplified by a amplifier and two electrodes made of Indalloy 19 (51 % indium, 32.5 % bismuth, 16.5 % tin). This process destabilizes the droplet interface and allows introduction of biological reagents via a pressurized injection channel as presented in Fig. 2d. The injection volume can be controlled precisely between 1 to 100 pl dependent on the applied pressure in the injection channel.



Fig. 2: Presentation of two droplet-based microfluidic operating units. (a) and (b) show the bright-field images of the flow-focusing junctions where the droplets with diameters of 40 and 100 μ m are generated, respectively. Bright-field images (c) and (d) show the pico-injection microfluidic unit for controlled introduction of biomaterials into the droplets. (c) The spacing between the droplets with different biological components is controlled through addition of oil via the second oil channel. (d) An AC field (1 kHz, 250 V) reduces the stability (poration) of the surfactant layer at the droplet surface allowing the injection of an aqueous solution of biological reagents from the pico-injection channel. The droplets' size before and after the injection can be compared. Photographs of (e) a droplet production and (f) a pico-injection device. Channel and tubing are filled with ink for better visualization. [21]

1.3 Biofunctionalization of copolymer-stabilized droplets

Two different approaches can be implemented to achieve efficient biofunctionalization of the nanostructured copolymer-stabilized droplets.[12] The first approach is based on the functionalization of the created droplets containing gold-linked copolymer surfactants (see Section 1.1) using a nitrilotriacetic acid (NTA) and hexahistidine-tagged (His6-tag) protein chemistry. The second approach involved two experimental steps: 1) synthesis of gold-linked surfactants coupled to bioactive molecules; and 2) formation of biofunctionalized droplets.

Two types of droplets loaded with His6-GFP-NTA-thiol were investigated: 1) those stabilized only by TRI7000 PFPE-PEG-PFPE (2.5 mM) surfactants and 2) droplets stabilized by a mixture of TRI7000 PFPE-PEG-PFPE (2.5 mM) and Gold-PEG-PFPE (3 μ M) surfactants. Fig. 3 shows fluorescence images of the (His6)-GFP-NTA-thiol within the gold-nanostructured droplets (Fig. 3a) and within the droplets containing no gold-linked surfactants (Fig. 3b). It can easily be observed that the fluorescent signal is confined to the surface of the

nanostructured droplets whereas it is homogeneously distributed in the entire volume of the non-nanostructured droplets.



Fig. 3: Representative fluorescence images of the (His6)-GFP-NTA-thiol within (a) goldnanostructured and (b) gold-free polymer-stabilized droplets measured 1 day after formation.[11]

FRAP analysis of the GFP-labeled gold-linked surfactants revealed similar diffusion coefficients of 0.21 ± 0.05 and $0.20 \pm 0.05 \ \mu m2/s$ when mixed with TRI7000 and TRI2500, respectively. These values are slightly lower than the values obtained from the "traditional" water-in-water polymersomes.[10] The lower values of water-in-oil droplets can be explained by the fact that the oil, FC-40, viscosity is 3.8 times higher than the viscosity of water. Oil viscosity might be also a potential reason for the similar diffusion coefficients, independent of the surfactant molecular weight.

2 Formation of droplets-stabilized GUVs (dsGUVs)

In this section, we describe the approach that merges lipid bilayer vesicle formation and droplet-based microfluidics for the generation of stable and easily manipulated dsGUV compartments (Fig. 1). To achieve the formation of dsGUV compartments, aqueous solution containing liposomes was encapsulated within droplets. To create a lipid bilayer within the droplets, we adapted the key factors that are necessary for an efficient formation of planar supported lipid bilayers (e.g. ion type and its concentration). In the following, we will describe the essential theoretical and experimental steps that are necessary for the development and analysis of the novel compartment system.

2.1 Lipid concentration

Droplet-based microfluidics allows for a high-throughput generation of monodisperse droplets (diameter difference < 1 %) with precise volume and surface area.[24] Therefore, it is possible to estimate the necessary amount of lipid required for formation of continuous lipid bilayer at the droplet surface. Because droplets are spherical, simple arithmetic gives for the necessary lipid concentration needed to cover the droplet by a bilayer.

2.2 Formation of dsGUV compartments by small unilamellar vesicles (SUVs) encapsulation

To create the dsGUV compartments, SUVs in MilliQ water are encapsulated into copolymerstabilized water-in-oil droplets by means of droplet-based microfluidics. However, no transfer of the encapsulated sUVs to the droplet interface in the form of lipid bilayer was observed, if no additional bivalent ions were added (Fig. 4). Therefore, to create dsGUV compartments, SUV solution containing 10 mM MgCl₂ was used as aqueous phase for the droplet creation. Box 30-1c shows that the fluorescence intensity was localized at the droplet interface, in comparison to the homogeneous distribution of the fluorescence signal seen in Fig. 4a where no Mg2+ ions were applied. From studies on the planar supported lipid bilayers formation, Mg2+ ions are known to be the most efficient mediators of lipid vesicle rupturing due to promotion of adhesion to the substrate.[25, 26]



Fig. 4: Formation process of dsGUVS by means of droplet-based microfluidic technology. In a first step, lipids in the form of (a) SUVs or (b) GUVs were encapsulated into water-in-oil droplets. Mg^{2+} ions (10 mM) were introduced into the droplets during droplet creation, pathway 1, or via pico-injection microfluidic technology, pathway 2, in order to transfer the encapsulated SUVs or GUVs into the form of supported lipid bilayer at the droplet interface (c). Fluorescence signal in the droplets is due to ATTO 488-labeled DOPE, which is part (1%) of the lipids mixture consisting of DOPC:DOPE:DOPS 8:1:1. [11]

2.3 Formation of dsGUVs by means of pico-injection

Microfluidic pico-injection technology (see Section 1.2) can be used as an alternative approach to create dsGUVs. Two experimental steps are required: the first is the formation of copolymer-stabilized droplets containing SUVs or GUVs (Box 30-1); the second step is the injection of an ionic aqueous solution of $MgCl^2$ into these droplets to a final ionic concentration of 10 mM. The pico-injection approach allows the analyzation of the dsGUV formation process, especially when free-standing GUVs are used for supported lipid bilayer formation.

Following the pico-injection step, time-lapse microscopy can be used to observe the dynamics of the dsGUVs formation. In case of SUVs-containing droplets, the fusion process of the SUVs to the droplet interface is observed immediately and lasts no longer than a minute. In contrast to SUVs, fusion dynamics of GUVs to the droplet interface is significantly slower – in the range of half an hour. This can be attributed to different diffusion coefficients of SUV and GUV species.

General notes: The easiest way to generate dsGUVs is by the "all-in-one" approach (Fig 4), i.e. encapsulation of SUVs solution containing 10mM of Mg2+ in to the copolymer

stabilized droplets. "Step-by-step" – pico-injection approach (Fig. 4) should be applied in case an analysis of dsGUVs formation is required.

2.4 Fluorescence intensity analysis and FRAP measurements of dsGUVs

To compare the membrane fluorescence intensity of encapsulated GUVs to that of dsGUVs as shown in Fig. 5, GUVs (egg PC:egg PG, 9:1, including 0.5 % ATTO 488-labeled DOPE) were encapsulated into the droplets and dsGUVs were produced from the same lipid composition using the "all-in-one" approach (Fig. 4). Both types of droplets were evaluated with identical settings of the confocal microscope. At least twenty intensity profiles were extracted for each droplet type.



Fig. 5: (*a*) and (*c*) Phase-contrast and fluorescence images of the encapsulated GUVs and dsGUVs (egg PC:egg PG, 9:1, 0.5 % ATTO 488 DOPE), respectively. (*b*) and (*d*) Fluorescence intensity profiles along the indicated lines as presented in (*a*) and (*c*), respectively. [11]

From both, microscope images as well as evaluated data, a light blur close to the droplet interface is observed. This is caused by refraction and diffraction at the water-oil interface due to a slight difference in the refractive indices of water (1.333) and FC-40 oil (1.290). This effect causes a widening of the intensity profile (Fig. 5b) and a reduction of the fluorescence intensity amplitude of the GUV part close to the droplet interface. Therefore, to compare fluorescence intensities, a Gaussian function with a background correction was fitted to the intensity peak profiles using a nonlinear least-square fit (Matlab 2015 SP1). The fitting revealed similar integrated intensity values of 42 ± 8 and 44 ± 4 a.u. × µm for dsGUVs and encapsulated GUVs, respectively. These findings suggest that freely suspended GUVs and dsGUVs consist of the same lipid bilayer conformation.

FRAP measurements were performed to compare the lipid diffusion coefficients in encapsulated GUVs and dsGUVs. Please note that full recovery of the bleaching spot was observed in all measurements. Table 1 presents the summary of diffusion coefficients of encapsulated GUVs and dsGUVs, consisting of various lipid composition and fluorophore types. The data shows a weak slowdown of diffusion in the dsGUV membrane. This outcome can be related to the fact that supported lipid membranes are subject to pertubation from the copolymer shell of the droplet, whose mobility is an order of magnitude lower (see Section 1.3). It is known that bivalent ions bind to phosphatidylcholine membranes[27] and lead to a decrease in the self-diffusion of lipids in the membrane.[28] Furthermore, FRAP and fluorescence correlation spectroscopy (FCS) measurements performed in other studies with

similar lipid compositions revealed diffusion coefficients in the same range as well as a similar tendency to lower values in the case of supported lipid membranes.[25, 29]

Diffusion coefficient Lipid Composition	Encapsulated GUV [µm2/s]	dsGUV [µm2/s]
DOPC:DOPE:DOPS (8:1:1) + 1 % ATTO488-labeled DOPE	3.52 ± 0.26	3.31 ± 0.77
Egg PC:Egg PG (1:1) + 1 % ATTO488-labeled DOPE	3.96 ± 0.51	2.88 ± 0.06
DOPC + 1 % Rhodamine B-labeled DOPE	4.42 ± 0.65	4.11 ± 0.59

Table 1. Summary of the diffusion coefficients obtained by FRAP measurements.

3 Sequential bottom-up assembly of dsGUVs

In order to dissect complex cellular sensory machinery by means of an automated dropletbased microfluidic approach, dsGUVs have to be adapted to allow the functional bottom-up assembly of various sub-cellular functional units. Thus, we focused on developing highthroughput strategies for incorporating transmembrane proteins, such as integrin and ATP synthase, and for immobilizing proteins to the bio-functionalized dsGUVs.

3.1 Protein reconstitutions in dsGUVs

Two experimental steps are required for proteoliposomes fusion (i.e., liposomes containing TAMRA-labeled $\alpha_{IIIb}\beta_3$ integrin or ATTO 488-labeled F_0F_1 -ATP synthase) with the preformed dsGUVs consisting of DOPC:DOPE:DOPS (8:1:1), including 1 % ATTO488-labeled DOPE or 1 % Rhodamine B-labeled DOPE, respectively. The first step is the creation of dsGUVs with the diameter of 40 µm using a solution of 800 µM liposomes and integrin activation buffer or ATP Synthase working buffer following the all-in-one approach. The second step is injection of proteoliposome solution into these droplets by means of pico-injection technology as schematically shown in Fig. 6. Colocalization of proteoliposomes with the dsGUVs.



Fig. 6: (a) Schematic of the process for proteoliposome fusion into dsGUVs via highthroughput pico-injection microfluidics. (b) and (d) are representative fluorescence images of the dsGUVs (DOPC:DOPE:DOPS, 8:1:1) 10 minutes after pico-injection, containing 1 % ATTO 488- or Rhodamine B-labelled DOPE, respectively. (c) and (e) are representative

fluorescence images of TAMRA-labelled $\alpha_{IIb}\beta_3$ integrin and ATTO 488-labeled F_0F_1ATP synthase incorporated in the dsGUVs presented in (b) and (d), respectively [11]

The same pico-injection approach was applied to obtain a bottom-up reconstitution of the actin cytoskeleton or microtubules within the dsGUVs. This was done in the following two-step process: (i) the formation of dsGUVs (90% DOPC, 9% DOPS and 1% RhB DOPE) in the presence of an actin or tubulin polymerization buffer; (ii) the pico-injection of G-actin (10 μ M final concentration, including 1% Alexa 488-labeled G-actin) or tubulin (10 μ M final concentration, including 10% ATTO 488-labeled tubulin) solution into the dsGUVs. Moreover, we compared this pico-injection approach to the pre-mixed (one-step) approach, in which G-actin or tubulin proteins were mixed with SUVs prior to droplet formation.

The successful reconstitution of the actin filaments and microtubules within the dsGUVs could be obtained by sequential pico-injection only (Fig. 7 lower row). When using the one-step premixed approach, vesicle fusion to the droplets' inner surface was suppressed. In the case of microtubules none and in the case of F-actin only partial fusion was observed (Fig. 7 upper row). The inability to form dsGUVs in the presence of microtubules is related to the fact that SUVs are subject to perturbations stemming from the amphiphilic nature of tubulin. The sequential pico-injection approach enables the prior formation of stable dsGUVs with appropriate ionic conditions followed by the injection of proteins without perturbing the lipid bilayer of the dsGUVs.



Fig. 7: Representative bright field images of cross-junction (top), pico-injection (bottom) microfluidic devices and fluorescence images of microtubules (10% ATTO 488-labeled tubulin, right panel) or the actin cytoskeleton (1% Alexa 488-labeled actin, left panel) in droplets containing RhB-labeled DOPE lipids, as obtained by either the pre-mixed (top) or the pico-injection approach (bottom). Scale bar 20 μ m. [11]

3.2 Mobility of transmembrane proteins in dsGUVs

FRAP measurements were performed to investigate the mobility of transmembrane proteins reconstituted in to the dsGUVs. Table 2 presents the summary of diffusion coefficients of

lipids D_{Lip} in protein-decorated dsGUVs and of the corresponding incorporated proteins D_{Pro} . In all cases, the measured D_{Lip} values were lower in comparison to the diffusion coefficients of the dsGUVs containing no proteins (see Table 1). Lower diffusion coefficient values can be attributed to the fact that the lipid lateral diffusion is a subject to steric and charge-related perturbations from the incorporated proteins.[30] As can be observed from Table 2, the FRAP measurements indicated similar diffusion coefficient values of $D_{Pro} \approx 0.7 \ \mu m^2/s$ for integrin, independently on whether they were introduced as pure proteins or as proteoliposomes. These values are in good agreement with previously published studies on integrin $\alpha_{IIb}\beta_3$ lateral mobility in planar supported lipid bilayers or in the cellular membranes as obtained by FRAP [31, 32] and FCS measurements [33], respectively.

To test the functionality of the incorporated integrin proteins, nanostructured dropletscontaining RGD-linked surfactants (see Section 1.3) were used to provide binding sites for integrin adhesion. In this case (see Table 2) the diffusion coefficient of integrin dropped significantly ($D_{Pro} = 0.13 \ \mu m^2/s$) to values that represent the mobility of the surfactant layer. This observation might indicate a successful establishment of a linkage between the transmembrane integrin and RGD-peptides on the copolymer droplet interface. It also reveals that at least some of the integrin are oriented correctly, i.e. that the extracellular part points towards the copolymer-stabilized droplet inner interface.

Lipid Composition	Protein	D _{Lip} [µm²/s]	D _{Pro} [µm²/s]
DOPC:DGS-NTA (9:1)	GFP		1.22±0.03
DOPC:DOPE:DOPS (8:1:1) + 1 % Rhodamine B-labeled DOPE	F ₀ F ₁ ATPsynthase	2.80±0.39	1.15±0.76
DOPC:DOPE:DOPS (8:1:1) + 1 % ATTO 488-labeled DOPE	Integrin (Proteoliposome)		0.67±0.10
DOPC:DOPE:DOPS (8:1:1) + 1 % ATTO 488-labeled DOPE	Integrin (pure)	2.27±0.22	0.70±0.06
DOPC:DOPE:DOPS (8:1:1) + ATTO 488-labeled DOPE	Integrin (Proteoliposome) + RGD	2.27±0.16	0.13±0.03

Table 2. Summary of the diffusion coefficients of lipids and proteins reconstituted in proteindecorated dsGUVs obtained by FRAP measurements.

4 Approaches to release GUVs from dsGUVs

The polymer-based surfactant shell and oil phase provid great stability to the dsGUV and, therefore allows the sequential loading of the compartment with biomolecules. However, it greatly restricts the possibility to study the behavior of the GUV-based model system in a physiological environment. Towards this end, we developed approaches to recover/extract protocells from the dsGUVs. To note, in the context of synthetic biology, protocells are synthetic, biomolecules-containing lipid-based compartments.

When designing the methods for protocell release, one has to take in to account the fragile nature of cell-size lipid vesicles due to their mechanical and chemical instabilities. Therefore, we decided to use copolymer-based destabilizing surfactants for the gentle GUV release. These demulsifiers partially replace their stable counterparts at the droplet interface and by doing so reduce the energy barrier to allow droplet coalescence. In the following paragraphs two developed bulk and microfluidic methods for GUVs release will be described.

Microfluidic device for GUVs release (Fig. 8) was designed to allow monitoring of the GUVs release process under controlled high-throughput conditions. Observation of the

release process allowed us optimization of the parameters necessary to get a high release yield. Moreover, findings obtained from microfluidic release setup were in turn used to improve bulk release conditions (Fig. 9).

In the microfluidic device for GUVs release a flow control system was used to control the pressure in aqueous and oil inlet channels. To minimize stress on the droplets (e.g. shear forces) pressures levels were set below 20 mbar with minor corrections for individual setups and experimental conditions. Moreover, the channel heights were designed to exceed the droplets diameter. Preformed dsGUVs were reinjected into the release chip and separated at the T-junction (Fig. 8b) by an additional oil flow containing 20 vol% perfluoro-1-octanol demulsifier. Following the T junction the stabilizing surfactants are replaced by the demulsifier from the spacing oil. The total flow was adjusted to allow efficient time (300 ms) to replace the stabilizing surfactants by the demulsifier prior to reaching the release unit where dsGUV encounter the aqueous phase in a wide perpendicular channel. To minimize the mechanical impact on the droplets at the oil/water junction, passive trapping structures within the microfluidic channels (i.e. rows of pillars separated by slits) have been designed and used to decelerate the droplets before coming in contact with the aqueous phase (Fig. 8c). On both sides the trapping unit is connected to the adjacent outlet oil channels. To allow only the oil phase to flow to the outlet channels, the width of these slits was designed to be smaller than the representative droplets dimensions. The droplet decelerates as it approaches the oil-water interface. Provided a sufficient concentration of destabilizing surfactant, the residual surfactant layer peels of the droplet at contact with the water phase and its content (i.e GUV) is released into the aqueous phase.



Fig. 8: Schematics of the microfluidic device for release of GUVs from dsGUVs. (a) In sake of clarity of presentation, the droplet injection channel is marked in green, the droplet separating channel (introduction of demulsifier-containing oil) and the outlet oil channels are marked in orange, the aqueous wide channel in blue. (b) Bright field microscopy image of the spacing T-junction. The displacement of each phase without droplets is indicated in a colored overlay – orange for demulsifier-containing oil surfactant and green for oil from the reinjection channel. (c) Bright field microscopy image of the passive trapping structures (i.e.

rows of pillars separated by distances smaller than the representative droplets dimensions) and the recovery area. [11]

Using destabilizing surfactants, bulk deemulsification is probably the most intuitive approach to release the content of droplets. Efficient release can be achieved by applying the destabilizing surfactant Perfluoro-1-octanol. The steps of our protocol to achieve bulk release are sketched in Fig. 9. For the bulk release approach 100 μ l of formed dsGUVs are collected in an Eppendorf tube. Due to the density differences between the FC-40 oil and water, the dsGUVs form a dense layer at the top of the tube. To provide an aqueous phase for release, 100 μ l of buffer was placed as a one large drop in the center of the droplet layer. To reduce osmotic pressure effect, it is preferable that the buffer ionic content will be identical to the buffer content within the dsGUVs. Following the addition of buffer, a 20 vol% of deemulsifier FC-40 oil solution was gently dripped on top of the buffer drop. After applying the complete volume of deemulsifier, the tube was tilted to increase the interface area and slowly rotated about its longitudinal axis. In that conditions the emulsion breakage takes less than five minutes. Residual oil drops in the aqueous phase were centrifuged down by briefly spinning at low speed with a table-top centrifuge. The aqueous solution containing GUV can be carefully removed with a pipette and immediately used for the analysis.



Fig. 9: Schematic of experimental steps for bulk release of GUVs from the dsGUVs. (a) Droplets are stored in an Eppendorf tube. (b) To provide an aqueous phase for release, 100 μ l of the buffer are added to the droplet layer. To minimize the effects that are related to osmotic pressure the buffer ionic content has to be similar to the content in dsGUVs. (c) 100 μ l of the demulsifier solution is added dropwise on the buffer drop. (d) The emulsion breaks and the released GUVs are going to aqueous phase.[21]

5 Summary and outlook for the future

In these lecture notes, we explored the capacity of droplet-based microfluidics for sequentional high-throughput bottom-up assembly of synthetic cells. To illustrate the necessity of our novel compartment system, we addressed the drawbacks of the currently available protocell systems and described technological limitations related to their manipulation. Thereafter, we summarized in detail our recently developed approach that merges lipid vesicle formation and droplet-based microfluidics for the generation of stable dsGUV compartments. By applying several droplet-based microfluidic functional techniques, including droplet generation and pico-injection, the versatility and robustness of the dsGUV compartment system for high-throughput manipulations was presented. The combination of these various technologies allows the sequential assembly of synthetic cells with a high

complexity. Moreover, by presenting methods to recover the constructed synthetic cells from the dsGUV, we show that this method allows to study their interaction with a physiological environment.

It is our expectation that the pliable biophysical properties of the dsGUV compartments and their integration to microfluidic technology might provide a system with superior properties for the assembly of a wide range of cellular functional units. A potential example of application is adhesion-associated complexes and cytoskeleton filament organization as schematically illustrated in Fig. 10.



Fig. 10: Schematic of bio-inspired "minimal synthetic cells". To provide bioactivity in terms of integrin adhesion, RGD peptides are immobilized on the gold nanoparticles via a thiol linker. G-actin and other proteins can be subsequently introduced into the dsGUVs by means of pico-injection technology. After the release, the synthetic cells and can be employed as cell-sized compartments within which interactions between different adhesion-associated proteins are systematically analyzed via a high-throughput screening platform.

References

[1] C.M. Agapakis, P.M. Boyle, P.A. Silver, Nat Chem Biol 8(6) (2012) 527-35.

[2] Y. Diekmann, J.B. Pereira-Leal, Biochem J 449(2) (2013) 319-31.

[3] M. Li, X. Huang, T.Y.D. Tang, S. Mann, Current Opinion in Chemical Biology 22 (2014) 1-11.

[4] Y. Elani, R.V. Law, O. Ces, Nature Communications 5 (2014) 5.

[5] S. Ramadurai, A. Holt, V. Krasnikov, G. van den Bogaart, J.A. Killian, B. Poolman, Journal of the American Chemical Society 131(35) (2009) 12650-12656.

[6] P. Carrara, P. Stano, P.L. Luisi, Chembiochem 13(10) (2012) 1497-1502.

[7] D.E. Discher, A. Eisenberg, Science 297(5583) (2002) 967-973.

[8] H. Bermudez, A.K. Brannan, D.A. Hammer, F.S. Bates, D.E. Discher, Macromolecules 35(21) (2002) 8203-8208.

[9] A.R. Abate, T. Hung, P. Mary, J.J. Agresti, D.A. Weitz, PNAS 107(45) (2010) 19163-19166.

[10] F. Itel, M. Chami, A. Najer, S. Lörcher, D. Wu, I.A. Dinu, W. Meier, Macromolecules 47(21) (2014) 7588-7596.

[11] M. Weiss, J.P. Frohnmayer, L.T. Benk, B. Haller, J.-W. Janiesch, T. Heitkamp, M. Börsch, R.B. Lira, R. Dimova, R. Lipowsky, E. Bodenschatz, J.-C. Baret, T. Vidakovic-Koch, K. Sundmacher, I. Platzman, J.P. Spatz, Nature Materials (2017) Accepted.

[12] I. Platzman, J.-W. Janiesch, J.P. Spatz, Journal of the American Chemical Society 135(9) (2013) 3339-42.

- [13] T.M. Pearce, J.C. Williams, Lab on a Chip 7(1) (2007) 30-40.
- [14] A.C. Rowat, J.C. Bird, J.J. Agresti, O.J. Rando, D.A. Weitz, PNAS 106(43) (2009) 18149-18154.
- [15] A.J. deMello, R.C.R. Wootton, Nature Chemistry 1 (2009) 28-29.
- [16] L.-H. Hung, A.P. Lee, Journal of Medical and Biological Engineering 27(1) (2007) 1-6.
- [17] Y. Huang, B. Agrawal, D.D. Sun, J.S. Kuo, J.C. Williams, Biomicrofluidics 5(1) (2011) 013412.
- [18] C. Holtze, A.C. Rowat, J.J. Agresti, J.B. Hutchison, F.E. Angile, C.H.J. Schmitz, S. Koster, H. Duan, K.J. Humphry, R.A. Scanga, J.S. Johnson, D. Pisignano, D.A. Weitz, Lab on a Chip 8(10) (2008) 1632-1639.
- [19] D.C. Duffy, J.C. McDonald, O.J.A. Schueller, G.M. Whitesides, Analytical Chemistry 70(23) (1998) 4974-4984.
- [20] Y. Xia, G.M. Whitesides, Annual Review of Materials Science 28(1) (1998) 153-184.
- [21] J.P. Frohnmayer, Heidelberg University, Heidelberg, 2017, pp. 1 Online-Ressource (238 Seiten).
- [22] L. Mazutis, J.C. Baret, A.D. Griffiths, Lab on a Chip 9(18) (2009) 2665-2672.
- [23] J.-W. Janiesch, Heidelberg, 2015, pp. 1 Online-Ressource (204 Seiten).
- [24] T. Thorsen, R.W. Roberts, F.H. Arnold, S.R. Quake, Physical Review Letters 86(18) (2001) 4163-4166.
- [25] T. Bhatia, P. Husen, J.H. Ipsen, L.A. Bagatolli, A.C. Simonsen, Biochimica Et Biophysica Acta-Biomembranes 1838(10) (2014) 2503-2510.
- [26] B. Seantier, B. Kasemo, Langmuir 25(10) (2009) 5767-5772.
- [27] C.G. Sinn, M. Antonietti, R. Dimova, Colloids and Surfaces a-Physicochemical and Engineering Aspects 282 (2006) 410-419.
- [28] R.A. Böckmann, H. Grubmüller, Angewandte Chemie International Edition 43(8) (2004) 1021-1024.
- [29] R. Machan, M. Hof, Biochim Biophys Acta 1798(7) (2010) 1377-91.
- [30] S. May, D. Harries, A. Ben-Shaul, Biophysical Journal 79(4) (2000) 1747-1760.
- [31] E.M. Erb, K. Tangemann, B. Bohrmann, B. Muller, J. Engel, Biochemistry 36(24) (1997) 7395-7402.
- [32] S. Goennenwein, M. Tanaka, B. Hu, L. Moroder, E. Sackmann, Biophysical Journal 85(1) (2003) 646-655.
- [33] J.B. Edel, M. Wu, B. Baird, H.G. Craighead, Biophysical Journal 88(6) (2005) L43-L45.

D 5 Theory of biological force sensing

B. Sabass

Theoretical Soft Matter and Biophysics Institute of Complex Systems -2 Forschungszentrum Jülich GmbH

Contents

1	Intro	oduction	2		
2	Mec	hanosensitive membrane channels	2		
	2.1	Tension-sensing membrane channels	2		
	2.2	Channels sensing force through attached tethers	4		
	2.3	A topical issue: Piezo channels	6		
3	Stretch-sensing structural molecules				
	3.1	Talin responds to differential levels of stretch	6		
	3.2	Opening and closing of stretched molecular folds	8		
4	Force measurement with a two-state sensor				
	4.1	A generic model for a molecular force sensor	9		
	4.2	Precision of force measurement with a two-state sensor	9		
5	Sequ	iential, threshold-based force sensing	10		
	5.1	The CUSUM test for change detection	10		
	5.2	Can optimal sequential tests be realized with membrane channels?	12		
A	Mechanical interaction of a channel with the surrounding membrane				
	A.1	Calculation of the membrane deformation	14		
	A.2	Free energy of the membrane around a tethered channel	16		
B	Revi	ew of the stochastic two-state process	17		

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Mechanotransduction – the cellular transduction of mechanical information into chemical signals – occurs on the level of specialized molecules. During the last decades, a large number of mechanotransduction molecules have been discovered, leading to the general perception that molecular mechanosensors exist in almost all types of cellular organisms, from bacteria to mammalian cells and plant cells. These molecular sensors can not only detect forces but also sense mechanical properties of the environment, including viscosity, pressure, and elastic deformability. For instance, most cells have means to detect tension in their outer membrane, which enables them to maintain their mechanical integrity. Since the molecular underpinnings of biological force sensing are slowly unraveling, there is an obvious need for theory to understand generic mechanisms. This chapter will introduce paradigmatic molecular mechanosensors along with established concepts for theoretical modeling. A few original contributions regarding sensitivity and optimal detection of force changes are also presented.

Measuring minute quantities with molecular sensors may require systems that are quite different from those employed for measurement in our macroscopic world. For instance, we usually measure macroscopic mechanical forces by recording the extension of a spring as shown in Fig. 1. This type of analog measurement is predicated on the availability of a gauged reference scale. Alternatively, one may be more interested in knowing if the force exceeds a given threshold. An example is a bistable light switch. In the nanoscopic world of biological molecules, the energy scales of the sensor and the signal are not far above the thermal energy scale k_BT . Thus, any system is subject to considerable fluctuations, which makes the two exemplary macroscopic measurement approaches hard to realize. Noise affects the gauged reference scale of an analog device as well as the definite state of a binary sensor. The question is then, which strategies are employed by nature to sense and interpret stochastic signals robustly and quickly.

2 Mechanosensitive membrane channels

Mechanosensitive channels are a class of membrane channels that open or close upon mechanical stimulation. When open, they allow passage of ions or solvent through the membrane, thus transducing the mechanical stimulus into a chemical signal. Membrane channels for mechanotransduction are expressed in almost all cells, including prokaryotes and mammalian cells. In fact, our perception of the world ultimately relies on these channels since they are responsible for our sense of touch or pain and enable hearing. Mechanosensitive channels are also required for maintaining tissue integrity, for blood pressure regulation, and osmoregulation [25]. Although the first studies of mechanosensitive channels are more than 30 years old [14], the various biophysical mechanisms of mechanotransduction by membrane channels have only become a very active research field during the last decade [9, 29, 3].

2.1 Tension-sensing membrane channels

One of the best-studied mechanism for cellular mechanotransduction is "tension sensing". The salient feature of the tension sensing mechanism is that membrane channels react to increased tension with shape changes and an ensuing pore opening, see Fig. 2a). Tension-sensing channels were largely studied with bacterial model organisms. Bacteria can grow in environments with various concentrations of salts and sugars, leading to considerable variations in osmotic


Fig. 1: Force measurement in our macroscopic world and in molecular systems of biological cells. a) Macroscopic device for analog force measurement consisting of a spring and a ruler. b) A light switch is a binary force sensor that flips if sufficiently strong forces are applied for a sufficiently long time. c) Adherent cells sense and control tensional forces ranging from about 10^{-12} N to about 10^{-9} N. d) Cellular membrane channels detect minute forces resulting from poking, shear flow, or osmotic pressure on the membrane.



Fig. 2: *a) Sketch of a tension-sensing membrane channel. b) Opening changes the area that the channel occupies. c) Measured opening of MscL and fit to Eq. (1). Data points taken from Ref. [30].*

pressure. It is thought that bacteria have evolved membrane channels that open if the tension in the membrane exceeds a critical level, thereby avoiding deformation and bursting of their outer membrane under osmotic pressure. Channels that could fulfill this function are the mechanosensitive channel with small conductance (MscS) and large conductance (MscL), which open at a tension around 10 mN/m, approaching the lytic tension of bacterial membranes. A popular physical mechanisms of tension-sensing relies on a radial expansion of the channel molecule [29, 23]. If a channel opens, it changes the area that it occupies in the membrane. Given a finite tension γ of the membrane, the free energy change associated with a positive area change ΔA is $\Delta \mathcal{F} = -\gamma \Delta A$. Denote the state of the channel by x with x = 0 being the closed state and x = 1 being the open state. Assigning internal energies of ϵ_0 and ϵ_1 to the closed and open states, the overall energy is given by $\mathcal{F}_x(\gamma) = (1 - x)\epsilon_0 + x\epsilon_1 - x\gamma\Delta A$. Assuming equilibrium, the probability of having an open channel is

$$p_1(\gamma) = \frac{e^{-\frac{\mathcal{F}_1(\gamma)}{k_{\rm B}\mathrm{T}}}}{\sum_{x=0,1} e^{-\frac{\mathcal{F}_x(\gamma)}{k_{\rm B}\mathrm{T}}}} = \frac{1}{e^{\frac{\epsilon_1 - \epsilon_0 - \gamma\Delta A}{k_{\rm B}\mathrm{T}}} + 1}.$$
(1)

The resulting sigmoidal dependence of the open probability on membrane tension has been measured for various channels. The corresponding data for MscL is shown in Fig. 2c). Such measurements allow to estimate the energy difference between open and closed channel. Having a radius on the order of $R_0 = 2 \text{ nm}$ in the closed state, the channel area can change up to $\Delta A = \pi (R_1 - R_0) \sim 20 \text{ nm}^2$. Using a typical opening membrane tension of $\gamma = 10 \text{ mN/m}$ we find $\epsilon_1 - \epsilon_0 \sim \gamma \Delta A = 2 \times 10^{-19} \text{ Nm} \simeq 48 \text{ k}_{\text{B}}\text{T}$. This high energetic barrier prevents thermal fluctuations from strongly affecting the opening state. Hence, tension-sensing channels with large area change are molecular sensors that selectively respond to large stimuli that would be critical for the integrity of the cell.

2.2 Channels sensing force through attached tethers

A second mechanism for mechanosensing is based on the idea that forces can be transmitted directly to the channels via attached tethers (Fig. 3a). Tethered channels occur in various biological systems, for instance in nociceptors (harm-sensing neurons) [12] and osmosensory neurons [24]. Notably, it has also been hypothesized that tethered channels are necessary for a)

tether

force F



Fig. 3: *a)* Sketch of a tethered membrane channel that opens when force is applied to the tether. *b)* Opening can proceed through a purely internal deformation (a "trap door" mechanism), but often changes the channel shape, leading to a radial or conical deformation.

h(r)

U)

hearing [15]. On the molecular level, one of the best-studied mechanosensitive ion channels is TRPN, a member of the Transient Receptor Potential channel family. It is responsible for touch sensation and hearing in Drosophila [10, 35]. TRPN has a remarkably long N-terminal module with 29 ankyrin repeats that tethers the channel to intracellular structures. Applying force to the ankyrin tether opens the channel. Interestingly, it has been shown that the ankyrin tether from TRPN can be fused to a voltage-gated potassium channel that is usually mechanoinsensitive, and then renders this channel mechanosensitive [35]. This finding raises the question whether sensing force at tethers could rely on simple micromechanical principles that are somewhat independent of the detailed molecular channel structure.

The physical mechanisms involved in opening of tethered channels are a subject of current research. Certainly, opening in response to force could be based on a rearrangement of the molecule only, which may be thought of as "trap door mechanism", see Fig. 3b). However, it is likely that the energetics of the membrane-channel interaction also play a role for opening [27]. In particular, conical deformations of the channel shape affect the membrane bending energy. To estimate the energetic effect of conical deformations, we consider the system depicted in Fig. 3a), where a channel connected to a tether bears a force F on the order of 10 pN. Now assume that the channel opening produces a tilt in the channel walls, changing its shape from a cylinder to a cone by an angle $\Delta \alpha$. The tilted boundary induces a bending on the membrane, which slightly changes the vertical position of the channel. For almost planar membranes, the height change is $\Delta h \approx \Delta \alpha R_0 \left[\Gamma_e + \log(R_0 \sqrt{\gamma/\kappa_b}/2)\right]$, where Γ_e is the Euler-Mascheroni constant and κ_b is the membrane bending modulus. A calculation of Δh is presented in Appendix A. Since a vertical force is applied to the channel, the height change Δh corresponds to a change in free energy given by

$$\Delta \mathcal{F}_{\Delta \alpha} \sim -\Delta h F = -\Delta \alpha F R_0 \left[\Gamma_e + \log(R_0 \sqrt{\gamma/\kappa_b}/2) \right].$$
⁽²⁾

The conical deformation $\Delta \alpha$ can become energetically favorable when a force F is applied to the tether. The corresponding parameter values have been measured, e.g., for the channel TREK-1 [20] and are $\Delta \alpha \sim 0.38$ rad and $R_0 \sim 2.5$ nm. With a bending modulus of $\kappa_b = 25 \,\mathrm{k_BT}$ and a typical membrane tension for eukaryotic cells of $\gamma = 10^{-3} \,\mathrm{k_BT/nm^2} \simeq 4.1 \times 10^{-3} \,\mathrm{mN/m}$ we find $|\Delta \mathcal{F}_{\Delta \alpha}| \approx 10 \,\mathrm{k_BT}$. Thus, the gain in elastic energy through conical deformation is quite large and actually exceeds the internal energetic barrier $\sim [4-5] \,\mathrm{k_BT}$ that resists deformation.

2.3 A topical issue: Piezo channels

A discussion of mechanosensitive membrane channels would not be complete without mentioning the molecules Piezo1 and Piezo2 [7]. During the last few years, it has become clear that Piezo proteins play critical roles in various mechanotransduction processes, including the sensing of mechanical harm (nociception), gentle touch [11], vascular functions, volume regulation of red blood cells [5], and may even play a role for the mechanical perception of oneself (proprioception) [32]. Mutations of the genes for Piezo in humans are also linked to hereditary diseases, see for example Ref. [34]. The physical mechanism underlying mechanosensing via Piezo channels in a physiological setting is a current subject of debate. They are more sensitive to mechanical stimulation than the bacterial "security valve" MscL, since the energy difference between closed and open state is only $\sim 9.7 \, k_B T$ [8]. Experiments clearly demonstrate that the channels can be opened by increasing membrane tension above around $5.1 \,\mathrm{mN/m}$. However, the molecules also sense cell poking and shear flow. Moreover, tethering of the membrane to the cytoskeleton or to the extracellular matrix strongly affects the response of Piezo channels. An intriguing feature of Piezo channels is that they inactivate themselves on the timescale of 100 ms after application of a constant force. This behavior may allow the channels to selectively respond to changes in the applied forces but reduces their sensitivity in the high-frequency regime [17].

3 Stretch-sensing structural molecules

3.1 Talin responds to differential levels of stretch

Eukaryotic cells have means to detect and regulate mechanical stress in the intracellular cytoskeleton [13, 28, 6]. Key cytoskeletal proteins are filamentous actin, myosin motors that produce contractile forces in the cytoskeleton, and molecules connecting the cytoskeleton to trans-membrane integrin complexes, see Fig. 4a). One molecule that plays an essential role is talin. In recent years, evidence has emerged that talin acts as a mechanosensor, responding to applied physiological forces that are generated by the actomyosin complex to strengthen the adhesion sites connecting the cell and the extracellular matrix.

Talin comprises of a head domain and 13 rod domains that are connected by an unstructured linker chain, see Fig. 4b). The molecule is up to around 100 nm long when stretched and various binding sites for actin, vinculin, and integrins are distributed along its length. Some of the binding sites are cryptically buried inside the protein structure and usually not accessible to the binding partners. When integrated into the cytoskeleton, talin constitutes a force-bearing linkage between the extracellular matrix (ECM) and the actomyosin contractile machinery by binding to integrins via the F3 domain in the N-terminal head and to actin via a number of actin-binding sites located along the talin rod [2]. In this arrangement, the domains R1 to R12 may all experience tensile forces that can lead to conformational changes making further binding sites that only become exposed under increasing tensile stress [33]. Thus, the binding affinity of talin to other molecules is differentially regulated by force. Since talin has multiple binding sites that open at increasing forces, its function is somewhat reminiscent of a classical force-measurement device employing a spring and a ruler Fig. 1a).



Fig. 4: *a)* Many eukaryotic cells possess dedicated adhesion structures that adapt their size and composition to the mechanical load that they are exposed to. b) The molecule talin connects transmembrane integrins with the intracellular actin cytoskeleton. c) If tensile forces between 5 pN and 25 pN are applied to talin, the molecule unfolds its rod domains one after the other and thereby allows binding of an increasing number of vinculin molecules. d) sketch of a hypothetical energy landscape separating the folded and unfolded state of a single fold in absence of force. e) Measured folding and unfolding rates of the R3 domain of Talin. The straight lines are fits with Eqns. (5a, 5b). Data points were taken from Ref. [33].

3.2 Opening and closing of stretched molecular folds

To calculate how the unfolding of stretched molecules depends on the applied force, we consider a linear molecule with only one folded site that can unfold in a reversible manner. Folding and unfolding in the absence of forces can be pictured as a transition between two minima in the free energy landscape determined by the molecular conformation, see Fig. 4d). The process can be idealized as a diffusive motion of the system state along a one-dimensional reaction coordinate. Calculation of the transition rates across an energetic barrier is the famous Kramers problem [16]. For energy barriers that are much larger than the thermal energy scale, one finds for the transition from one minimum (state j) to the other minimum (state i) that

$$k_{ij} \approx c_0 e^{-\frac{(\mathcal{G}_{TS} - \mathcal{G}_j)}{k_{\rm BT}}},\tag{3}$$

where \mathcal{G}_{TS} is the maximum of free energy at the transition state. The exponential dependence of a reaction rate on an energy divided by k_BT is also known as Arrhenius law. To employ this framework for folding and unfolding of a chain-like molecule, we make the highly simplified assumption that the system behaves as an elastic spring with elastic constant d and rest length L_0 in the folded state. If the molecule unfolds, the rest length of the chain increases. The rest length is $L_0 + \delta L_1$ in the transition state and $L_0 + \delta L_1 + \delta L_2$ in the unfolded state. Hence, we can assign the following energies to the folded state (0), the transition state (TS), and the unfolded state (1)

$$\mathcal{G}_0 = \mu_0 + \frac{d}{2}(y - L_0)^2,$$
(4a)

$$\mathcal{G}_{TS} = \mu_{TS} + \frac{d}{2}(y - (L_0 + \delta L_1))^2,$$
 (4b)

$$\mathcal{G}_1 = \mu_1 + \frac{d}{2}(y - (L_0 + \delta L_1 + \delta L_2))^2,$$
 (4c)

where y represents the extension of the molecule. Employing Eq. (3) and the state-dependent force $F = d(y - L_0)$ or $F = d(y - (L_0 + \delta L_1 + \delta L_2))$ we find

$$k_{10} = c_0 e^{\frac{\mu_0 - \mu_{TS} - d\delta L_1^2/2 + \delta L_1 \, d(y - L_0)}{k_{\rm BT}}} = \hat{k}_{10} e^{\frac{\delta L_1}{k_{\rm BT}}F},\tag{5a}$$

$$k_{01} = c_0 e^{\frac{\mu_1 - \mu_{TS} - d\delta L_2^2 / 2 - \delta L_2 \, d(y - (L_0 + \delta L_1 + \delta L_2))}{k_{\rm BT}}} = \hat{k}_{01} e^{-\frac{\delta L_2}{k_{\rm BT}}F},\tag{5b}$$

where \hat{k}_{10} and \hat{k}_{01} are constants. Thus, the transition rates from one state into the other depend exponentially on the force F. The prefactors $\delta L_{1,2}/(k_BT)$ in the exponent determine how differently F affects forward and reverse rates and depend on molecular details. Note that the rates naturally satisfy a local detailed balance constraint

$$\frac{k_{01}}{k_{10}} = e^{\frac{\mathcal{G}_1 - \mathcal{G}_0}{k_{\rm B} {\rm T}}},\tag{6}$$

and therefore are thermodynamically admissible. In reality, biological molecules rarely behave like ideal elastic springs and unfolding is usually a complex process involving multiple transition states and pathways. Nevertheless, Eqns. (5a, 5b) are often a good approximation, as can be seen for example by comparison with the folding rates of the R3 domain of talin, see Fig. 4e). The measured data points for this comparison were taken from Ref. [33].

4 Force measurement with a two-state sensor

4.1 A generic model for a molecular force sensor

In this section, we discuss a simple model that epitomizes the theory of molecular sensors. We idealize a force sensor, for instance a membrane channel, as a two-state system. The system state is described by the binary variable $x \in \{0, 1\}$. The rate constant for a transition $0 \rightarrow 1$ is denoted by k_{10} and the rate constant for $1 \rightarrow 0$ is denoted by k_{01} . We assume Arrhenius-type approximations for the rates

$$k_{10} = \hat{k}_{10} e^{rF}, \qquad \qquad k_{01} = \hat{k}_{01} e^{-rF}, \qquad (7)$$

where \hat{k}_{10} , \hat{k}_{01} , and r are constants. Thus, a positive force F increases the probability to be in x = 1 by increasing the transition rate into this state and by also decreasing the transition rate out of this state. To describe a sequence of transitions, we employ a two-state Markov process also known as random telegraph process. $P(x, t|x_0, t_0)$ is the probability to be in state x at time t given a state x_0 at time t_0 . The probabilities obey

$$1 = P(1, t|x_0, t_0) + P(0, t|x_0, t_0)$$
(8)

and the Master equation reads

$$\partial_t P(1,t|x_0,t_0) = -k_{01} P(1,t|x_0,t_0) + k_{10} (1 - P(1,t|x_0,t_0)).$$
(9)

The steady state expectation values and variances are given by

$$\langle x \rangle_{ss} = \frac{k_{10}}{k_{10} + k_{01}},\tag{10}$$

$$\sigma_{ss}^2 = \langle x^2 \rangle_{ss} - \langle x \rangle_{ss}^2 = \langle x \rangle_{ss} - \langle x \rangle_{ss}^2.$$
(11)

A more detailed discussion of the two-state model is presented in Appendix B.

4.2 Precision of force measurement with a two-state sensor

In this subsection, we study how the two-state sensor can be used to measure the value of the force F. The probabilities to be in either of the state are a unique function of F, as illustrated for example in Fig. 2c). Thus, cells could in principle determine the magnitude of a constant force by recording the statistics of the sensor state. If the sensor state x is monitored for a long time T, the average is $\bar{x} = \frac{1}{T} \int_0^T x \, dt \approx \langle x \rangle_{ss}$. Then, an estimate for F follows from inverting $\langle x \rangle_{ss} = 1/(1 + \hat{k}_{01}e^{-2rF}/\hat{k}_{10})$, which gives

$$F \approx \log[\hat{k}_{01}\bar{x}/(\hat{k}_{10} - \hat{k}_{10}\bar{x})]/(2r).$$
(12)

The error of this estimate depends on the length of the measurement time T. Since F is constant, all the error must come from the stochastic fluctuations in the sensor state, characterized by the variance of \bar{x} . This variance is to be calculated as $\bar{\sigma}_{ss}^2 \equiv \langle \bar{x}\bar{x} \rangle_{ss} - \langle \bar{x} \rangle_{ss}^2$. Using the correlation function (46) derived in the appendix together with Eq. (11) we find

$$\bar{\sigma}_{ss}^2 = \frac{1}{T^2} \int_0^T \int_0^T \langle x(t)x(t') \rangle_{ss} \, \mathrm{d}t \mathrm{d}t' - \langle \bar{x} \rangle_{ss}^2 \approx \frac{2}{T(k_{01} + k_{10})} \sigma_{ss}^2, \tag{13}$$

where we assumed that the measurement time is much longer than the characteristic timescale of the sensor $T \gg 1/(k_{01} + k_{10})$. The meaning of Eq. (13) is that time averaging reduces the variance in the measurement of \bar{x} to leading order by a factor of $(k_{01} + k_{10})^{-1}/T$. Knowing the variance of the estimated state \bar{x} , we can now proceed to calculate the corresponding uncertainty in force measurement δF . A Taylor expansion for small δF yields

$$\left(\frac{\partial \bar{x}}{\partial F}\right)^2 (\delta F)^2 = \sigma_F^2. \tag{14}$$

When solving for $(\delta F)^2$ and inserting Eq. (13), we find a remarkably simple result

$$(\delta F)^2 = \frac{1}{2r^2T} \left(\frac{1}{\hat{k}_{10}e^{rF}} + \frac{1}{\hat{k}_{01}e^{-rF}} \right).$$
(15)

The relation implies that the uncertainty is dominated by the smaller of the two rates. Good measurements are only possible if the timescale of both state transitions is much shorter than the measurement time T. Eq. (15) is almost identical to a well-known formula for the precision of concentration sensing in biological systems, the so-called Berg-Purcell limit [4, 1, 31]. However, in the context of force sensing, Eq. (15) poses a rather stringent constraint on the range of forces that can be measured since the rates k_{ij} depend exponentially on F. The signal-to-noise ratio is $F^2/(\delta F)^2 \sim F^2 \hat{k}_{ij}T r^2 \exp(-r|F|)$. If the force-sensing module is, for example, a molecular folding site that is similar to those in talin, we expect $r \sim 5 \text{ nm/k}_{B}T \sim 1/\text{pN}$. The signal to noise ratio shows that extending the working range of such a force sensor over a range of [0 - 10] pN requires taking extensive statistics with $\hat{k}_{ij}T \gtrsim 10^3$. Note that this result only holds if both rates are force-dependent. Having one force-independent rate may be somewhat advantageous.

In summary, however, measuring the analog value of forces with a two-state sensor can be challenging. Also, one may argue that biological cells often do not need to know the precise magnitude of a mechanical stimulus. Rather, it is important to quickly detect forces, stresses, or tensions if they exceed a certain threshold. The focus of the next section will be optimal on-line detection of such events.

5 Sequential, threshold-based force sensing

One of the purposes of biological force sensing is to help conserve mechanical integrity of cells and tissue. To be able to respond to external forces appropriately, cells must be able to identify those situations where the mechanical load is too high. Typically, one can characterize destructive amounts of mechanical load by a threshold. Then, the challenge for cells is to react as quickly as possible if this threshold is surpassed. It is thus interesting to ask if there are strategies that are optimal for the detection of such events.

5.1 The CUSUM test for change detection

Consider a stochastic two-state sensor with a time-dependent state x that is affected by the signal F. The sensor could be represented, for instance, by the model introduced in Sec. 4.1. For simplicity, we assume that F changes during the observation time in a step-wise fashion from a value F_0 to a different value F_1 as shown in Fig. 5a). We will now introduce a procedure to detect this change quickly on the fly, without producing many false detections.

 g_k

Denote the probability to observe a sequence of states $[x_1, \ldots, x_k]$ by $\prod_{i=1}^k P_F(x_i)$. To quantify the relative probability of having $F = F_1$ versus $F = F_0$, we consider the logarithm of the probability ratios

$$\mathcal{L}_{k} = \log\left(\frac{\prod_{i=1}^{k} P_{F_{1}}(x_{i})}{\prod_{i=1}^{k} P_{F_{0}}(x_{i})}\right) = \sum_{i=1}^{k} \log\left(\frac{P_{F_{1}}(x_{i})}{P_{F_{0}}(x_{i})}\right) = \sum_{i=1}^{k} \Delta \mathcal{L}_{i}.$$
(16)

As long as $F \sim F_0$, the log likelihood decreases since the denominator is larger than the numerator. After the change, when $F \sim F_1$, the log likelihood increases. Hence, we can determine a change point by locating the minimum in \mathcal{L} at which the function switches from a decreasing trend to an increasing trend. However, to avoid detection errors, it is advisable to wait a little bit after \mathcal{L} has reached its minimum. By sampling \mathcal{L} for a bit following a minimum, we ensure that we are not detecting random fluctuations but instead pick up the real trend. To constantly test if the log likelihood has a positive trend, we iteratively calculate a decision function g_k using the likelihood increments $\Delta \mathcal{L}_i = \log \left(\frac{P_{F_1}(x_i)}{P_{F_0}(x_i)}\right)$ as

$$g_k = g_{k-1} + \Delta \mathcal{L}_k \qquad \text{if } g_{k-1} + \Delta \mathcal{L}_k \ge 0, \qquad (17a)$$

$$= 0 \qquad \qquad \text{if } g_{k-1} + \Delta \mathcal{L}_k < 0. \tag{17b}$$

This procedure is continued until the decision function g exceeds a fixed value h. Then the decision is made that F has changed from F_0 to F_1 . If we assign a time t_k to every measurement x_k , the time at which the decision is made can be formally expressed as

$$\tau_d \equiv \min(t_k | g_k > h). \tag{18}$$

A graphical illustration of the procedure is given in Fig. 5a)-c). This method for detecting changes was first suggested by E.S. Page [22] around 50 years ago. Since it is based on an evaluation of a cumulative sum, it is commonly referred to as CUSUM test. It can be shown that the CUSUM test is optimal in the following sense: if a detection threshold h is chosen such that false detections occur with a mean period that is larger than a constant γ , then the CUSUM test has the smallest worst mean delay for detection of a real change. The proof of this optimality statement is quite technical and comes in different variants, for asymptotic optimality when $\gamma \to \infty$ [19], in a Bayesian framework [26], and for continuous times and different stochastic dynamics determining $\Delta \mathcal{L}_k$ [21]. Here, we will content ourselves with an estimation of how quick the response of the CUSUM test is.

Assuming that the force change occurs at a time τ_c , the delay between τ_c and the decision time $\tau_d \geq \tau_c$ is determined by the requirement of taking a random number of m samples. To estimate m, we consider the sequence of log likelihood values just after the change occurred $\Delta \mathcal{L}_1, \Delta \mathcal{L}_2 \dots \Delta \mathcal{L}_m$. For simplicity, we assume that the $\Delta \mathcal{L}_i$ all have the same first moment $\langle \Delta \mathcal{L}_i \rangle_{F_1} = \langle \Delta \mathcal{L} \rangle_{F_1}$ where $\langle \dots \rangle_{F_1}$ denotes the expectation value with respect to the distribution determined by F_1 . We can write

$$\langle \sum_{i=1}^{m} \Delta \mathcal{L}_i \rangle_{F_1} = \langle \sum_{i=1}^{\infty} \Theta(m-i) \Delta \mathcal{L}_i \rangle_{F_1} = \sum_{i=1}^{\infty} \langle \Theta(m-i) \Delta \mathcal{L}_i \rangle_{F_1}$$

$$= \sum_{i=1}^{\infty} \langle \Delta \mathcal{L}_i \rangle_{F_1} \langle \Theta(m-i) \rangle_{F_1} = \langle \Delta \mathcal{L} \rangle_{F_1} \langle m \rangle_{F_1},$$
(19)



Fig. 5: The classical CUSUM test discriminates between two values of a continuous signal. a) It is assumed that the signal F jumps between two values F_0 and F_1 . Measurements are conducted at discrete times. b) The logarithmic likelihood of F_1 versus F_0 has a decreasing trend for $F \sim F_0$ and has a increasing trend for $F \sim F_1$. c) The decision function g_k records increases in the log likelihood. If $g_k \ge h$, the decision is made that F_1 is the new signal value.

where we employed the unit step function $\Theta(y) = 0$ for y < 0 and $\Theta(y) = 1$ for $y \ge 0$. At the decision time τ_d , the procedure of the CUSUM test requires $\langle \sum_{i=1}^m \Delta \mathcal{L}_i \rangle_{F_1} \approx \langle g_k \rangle_{F_1} \approx h$. It follows from Eq. (19) that the expected number of measurements producing a delay between force change and sensor response is approximated by

$$\langle m \rangle_{F_1} \approx \frac{h}{\langle \Delta \mathcal{L} \rangle_{F_1}} = \frac{h}{\langle \log\left(\frac{P_{F_1}(x)}{P_{F_0}(x)}\right) \rangle_{F_1}}.$$
 (20)

The term in the denominator of Eq. (20) is called Kullback-Leibler divergence and is a measure for how different the two distributions P_{F_1} and P_{F_0} are. The Kullback-Leibler divergence only becomes zero if the two distributions are equal, leading to a divergent mean detection delay.

5.2 Can optimal sequential tests be realized with membrane channels?

To study a concrete example of how the CUSUM test can be important in biology, we consider a force-sensing membrane channel that is described by the two-state model introduced in Sec. 4.1. In our two-state model, the waiting times in both states are exponentially distributed. The probability to observe a sequence of states is thus given by $P_F(\{\tau_1, \tau_2, \ldots\}) = k_{ij}e^{-k_{ij}\tau_1}dt_1 k_{ji}e^{-k_{ji}\tau_2}d\tau_2 \ldots$ with $ij \in \{01, 10\}$, see Appendix B. The probability to observe a waiting time τ_n during which the system remains in one state $x \in \{0, 1\}$ can be conveniently written as

$$P_F(\tau_n) = e^{-(xk_{01}(F) + (1-x)k_{10}(F))\tau_n} k_{01}^{j_0^{10}}(F) k_{10}^{j_1^{10}}(F) \,\mathrm{d}\tau.$$
(21)

Here, the indicator functions j_n^{01} and j_n^{10} are only non-zero if the state changes at the end of τ_n . If the system state changes as $0 \to 1$, we set $j_n^{01} = 1$ and if the change is $0 \to 1$ we set $j_n^{10} = 1$.



Fig. 6: Results for detection of a sudden force change from 0 to F > 0 at time τ_c . The parameters of the CUSUM test are $F_0 = 0$, $F_1 = 4$, and h = 100. We set r = 0.1, $k_{01} = k_{10} = 1$. a) Simulation results for the waiting times from the proper CUSUM test, Eqns. (17a, 17b), compare well with the analytical approximation given in Eq. (26). b) The approximate CUSUM test realized with concentrations of signaling molecules c(t), Eqns. (24a, 24b), responds to force changes with almost the same delay as the CUSUM test.

The increment of the logarithmic likelihood is thus given by

$$\Delta \mathcal{L}(\tau_n) = \log\left(\frac{P_{F_1}(\tau_n)}{P_{F_0}(\tau_n)}\right) = (k_{01}(F_0) - k_{01}(F_1))x(t)\tau_n + j_n^{01}\log\left(\frac{k_{01}(F_1)}{k_{01}(F_0)}\right) + (k_{10}(F_0) - k_{10}(F_1))(1 - x(t))\tau_n + j_n^{10}\log\left(\frac{k_{10}(F_1)}{k_{10}(F_0)}\right).$$
(22)

Next, we assume that the force is initially small and set $F_0 = 0$. Using the expressions (7) for the force-dependent transition rates yields

$$\Delta \mathcal{L}(\tau_n) = [\hat{k}_{01}(1 - e^{-rF_1}) + \hat{k}_{10}(e^{rF_1} - 1)]x(t)\tau_n - \hat{k}_{10}(e^{rF_1} - 1)\tau_n + (j_n^{01} - j_n^{10})rF_1 \sim [\hat{k}_{01}(1 - e^{-rF_1}) + \hat{k}_{10}(e^{rF_1} - 1)]x(t)\tau_n - \hat{k}_{10}(e^{rF_1} - 1)\tau_n.$$
(23)

In the second line we have neglected the terms $\sim \pm rF_1$ since these represent short kicks of alternating sign and therefore do not produce a continuous trend in the likelihood function.

If cells are to make use of the CUSUM test to detect a force $F_1 > 0$, the detection procedure must be implemented biochemically, which includes a repeated evaluation of Eq. (23). To see how this occurs naturally in membrane channels, we consider the dynamics of signaling molecules, e.g., calcium ions, that pass through the channel. We denote the intracellular concentration of the signaling molecule by c(t) and assume that the molecule is available in excess outside of the cell. We then have $\frac{dc(t)}{dt} = a x(t) - b(c(t))$ where a is the rate at which molecules can pass through the open channel and b(c(t)) denotes the rate at which the signaling molecules are removed. Typically, the molecules are being pumped actively out of the cytosol. In the case of calcium ions, an established expression for the rate of pumping is $b(c) = \hat{b}_1 c^2 / (\hat{b}_2 + c^2)$ with two constants $\hat{b}_{1,2}$ [18]. Assuming that the ion pumps operate always at maximum speed, we take $\hat{b}_2 \ll c^2$ and $b(c) \approx \hat{b}_1$. Then, the concentration changes in each time step Δt as

$$c(t + \Delta t) \approx c(t) + ax(t)\Delta t - \hat{b}_1\Delta t \qquad \text{if } c(t) + ax(t)\Delta t - \hat{b}_1\Delta t \ge 0, \tag{24a}$$

$$c(t + \Delta t) = 0 \qquad \qquad \text{if } c(t) + ax(t)\Delta t - b_1\Delta t < 0. \tag{24b}$$

The concentration increment in Eq. (24a) has the same form as the likelihood increment given in Eq. (23), except that the waiting times τ_n are replaced by infinitesimal time steps Δt . The concentration c(t) also obeys the same dynamics as the decision function g_k for the CUSUM test, Eqns. (17a,17b). Therefore, we suggest that c(t) can act as a continuous approximation for g_k to decide whether the force has exceeded a prescribed threshold. The reaction parameters aand \hat{b}_1 determine the force $F_1 > 0$ and comparison with Eq. (23) yields

$$a/\hat{b}_1 = [\hat{k}_{01}(1 - e^{-rF_1}) + \hat{k}_{10}(e^{rF_1} - 1)]/[\hat{k}_{10}(e^{rF_1} - 1)].$$
(25)

The threshold in c(t) leading to a detection of a force can be realized, for example, by a chemical reaction that depends non-linearly on c(t) to produce a step-like response if $c(t) \ge \tilde{h}$.

To estimate the delay until detection we consider the waiting time probability of consecutive closed and open states $P_F(\tau_0, \tau_1) = k_{10}(F)e^{-k_{10}(F)\tau_0} d\tau_0 k_{01}(F)e^{-k_{01}(F)\tau_1} d\tau_1$. Using the same reasoning as for Eq. (20), we obtain the average number of consecutive pairs of state changes at force F as

$$\langle m_{\rm oc} \rangle_F \approx \frac{h}{\int \int_0^\infty \log\left(\frac{P_{F_1}(\tau_0,\tau_1)}{P_{F_0}(\tau_0,\tau_1)}\right) P_F(\tau_0,\tau_1) \,\mathrm{d}\tau_0 \mathrm{d}\tau_1} = \frac{h}{2[\cosh\left(rF\right) - \cosh\left(rF_1 - rF\right)]}.$$
 (26)

With the average number of state-change pairs given, an approximate delay time results as $\langle (\tau_0 + \tau_1) \rangle_F \langle m_{\rm oc} \rangle_F = (1/k_{10}(F) + 1/k_{01}(F)) \langle m_{\rm oc} \rangle_F.$

In Fig. 6a), results from a simulation of the CUSUM test for detection of a suddenly applied force F are presented. The likelihood is calculated with the constant parameter $F_1 = 4 \times 0.1/r$, while the true magnitude of the force F is varied. It can be seen the sensor only responds for $F > F_1/2$. Moreover, the sensor responds to forces that are above the threshold $F > F_1$ on average faster than for $F = F_1$. The analytical approximation for the delay time agrees well with the simulation results. Figure 6b) shows a comparison of the proper CUSUM test based on evaluation of g_k , Eqns. (17a, 17b), with the approximate biological realization based on the concentration of signaling molecules c(t), Eqns. (24a, 24b). Clearly, the performance of the biological realization is very similar to the performance of the full CUSUM test. Therefore, we can surmise that biology employs a signal integration strategy like the CUSUM test if optimal detection of forces is required.

Finally, we mention that membrane-channel systems combined with a threshold-like response to a critical concentration of signaling molecules is a common motive in cell biology. Examples include depolarization of nerve cells and local intracellular calcium responses. It is tempting to conclude that signal integration principles akin to those idealized by the CUSUM test are a generic feature of such systems.

Appendices

A Mechanical interaction of a channel with the surrounding membrane

A.1 Calculation of the membrane deformation

We assume that the membrane-channel system is in mechanical equilibrium and consider an almost planar lipid membrane. The height h of the membrane above a reference plane is to be

a unique function of a two-dimensional position vector **r** lying in the reference plane. For the position vector, we employ a cylindrical coordinate system with radius *r* and angular coordinate φ , see Fig. 3. The nabla operator ∇ operates in the two-dimensional reference plane and, since we assume small gradients, the functional determinant is approximated as $\sqrt{1 + (\nabla h)^2} \approx [1 + \frac{1}{2}(\nabla h)^2]$. Within this framework, a fluid membrane can be described with the following energy

$$\mathcal{H}_h = \int \frac{\kappa_{\rm b}}{2} (\nabla^2 h)^2 + \gamma [1 + \frac{1}{2} (\nabla h)^2] \mathrm{d}^2 r, \qquad (27)$$

where the surface integration extends over the entire membrane. Here, κ_b is the bending constant of the membrane and γ is the membrane tension. The first term in Eq. (27) is an energy penalty resulting from non-zero mean curvature while the second term penalizes changes of the area. A variation of the energy yields

$$\delta \mathcal{H}_{h} = \int \kappa_{\rm b} (\nabla^{2} h) (\delta \nabla^{2} h) + \gamma (\nabla h) (\delta \nabla h) d^{2} r = \int \nabla^{2} (\kappa_{\rm b} \nabla^{2} h - \gamma h) \delta h d^{2} r + \int [\kappa_{\rm b} \nabla^{2} h (\delta \nabla h) - \nabla (\kappa_{\rm b} \nabla^{2} h - \gamma h) (\delta h)] \, \mathbf{s} ds,$$
(28)

where we employed partial integration and the divergence theorem with s denoting a normal vector pointing outwards from the membrane area on the contour path s. Since \mathcal{H}_h is minimal in mechanical equilibrium, the equation determining h follows from the first line of Eq. (28) as

$$\nabla^2 (\nabla^2 - \xi^2) h = 0, \tag{29}$$

with $\xi^2 \equiv \gamma/\kappa_b$. We next assume that the membrane forms a radially symmetric annulus around a circular channel protein with radius R. The membrane extends far out to a radius $L \gg R$. At the outer contour of the membrane, we fix the membrane height and slope as

$$h(L) = 0, \tag{30a}$$

$$\partial_r h(r)|_{r=L} = 0. \tag{30b}$$

For the inner contour surrounding the channel we assume the boundary conditions sketched in Fig. 3a) with a height h(R) and a contact angle α given by

$$h(R) = h(0), \tag{31a}$$

$$\partial_r h(r)|_{r=R} = \alpha. \tag{31b}$$

If a constant vertical force F is applied via the tether at r = 0, we need to employ a free energy that takes the constant force acting on the membrane contour into account

$$\mathcal{F} = \mathcal{H}_h - \int_0^{2\pi} \frac{Fh(R)}{2\pi R} R \mathrm{d}\varphi.$$
(32)

The variation $\delta \mathcal{F}$ yields the expression (28) minus $1/(2\pi) \int_0^{2\pi} F \delta h(R) d\varphi$. When calculating $\delta \mathcal{F}$, the first term in the second line of Eq. (28) is irrelevant since the contact angle at the channel is fixed by Eq. (31b). A vanishing $\delta \mathcal{F}$ in mechanical equilibrium requires

$$F = -\kappa 2\pi \partial_r \left(\xi^2 h - \nabla^2 h\right) r|_{r=R},\tag{33}$$

as well as equation (29) for the interior of the membrane. These equations determine h completely. The solution fulfilling the above differential equations and boundary conditions is

$$h = \frac{F \log(L/r)}{2\gamma\pi} - \frac{(F + 2R\alpha\gamma\pi)K_0(r\xi)}{2R\gamma\pi\xi K_1(R\xi)},$$
(34)

where K_n are the Bessel K functions of n-th order and we have dropped all terms that decay exponentially with $\xi L \ll 1$. Usually, the membrane tension is weak enough to guarantee that the lengthscale ξ^{-1} set by the tension and bending constant is much larger than the nanometerscale that is characteristic for membrane channels. Hence, we assume $\xi r \ll 1$ and $\xi R \ll 1$ to obtain

$$h \approx \frac{F}{2\gamma\pi} \log(L/r) + \frac{(F + 2\alpha\gamma\pi R)}{2\gamma\pi} \left[\Gamma_e + \log\left(R\xi/2\right)\right],\tag{35}$$

where Γ_e is the Euler-Mascheroni constant. If F is held constant and α changes, the resulting height change is $\Delta h = \Delta \alpha R \left[\Gamma_e + \log(R\xi/2)\right]$, which is the expression used in the main text.

A.2 Free energy of the membrane around a tethered channel

For a full analysis of how the force F changes the free energy of the membrane-channel system, we need to calculate the expression in (32) explicitly. On employing the the identity $(\nabla h)^2 = \nabla \cdot (h\nabla h) - h\nabla^2 h$ along with Eq. (29), the deformation energy becomes

$$\mathcal{H}_{h} = -\frac{\gamma}{2} \int h \partial_{r} h \, r \mathrm{d}\varphi_{b}|_{r=R} + \frac{\gamma}{2} \int h \partial_{r} h \, r \mathrm{d}\varphi_{b}|_{r=L} - \frac{\gamma}{2} \int h^{H} \nabla^{2} h^{S} \mathrm{d}^{2} r, \qquad (36)$$

where h^H and h^S are the parts of h that fulfill $\nabla^2 h^H = 0$ and $(\nabla^2 + \xi^2)h^S = 0$. Again, we ignore terms that decay exponentially with $\xi L \gg 1$ and obtain up to an L-dependent constant

$$\mathcal{H}_{h} = -\gamma \pi R^{2} - \frac{F^{2} K_{0}(R\xi)}{4R\gamma \pi \xi K_{1}(R\xi)} + \frac{F^{2} \log(L/R)}{4\gamma \pi} + \frac{(2R\gamma \pi \alpha)^{2} K_{0}(R\xi)}{4R\gamma \pi \xi K_{1}(R\xi)}.$$
 (37)

The work related to application of force F is given by

$$Fh(0) = Fh(R) = \frac{F^2 \log(L/R)}{2\gamma\pi} - \frac{(F^2 + 2R\gamma\pi\alpha F)K_0(R\xi)}{2R\gamma\pi\xi K_1(R\xi)}.$$
(38)

Adding the last two equations, we obtain for the overall free energy

$$\mathcal{F} = -\gamma \pi R^2 + \frac{(F + 2\pi\kappa_b \xi^2 R\alpha)^2}{4\pi\kappa_b \xi^2} \frac{K_0(R\xi)}{R\xi K_1(R\xi)} - \frac{F^2 \log(L/R)}{4\pi\gamma}.$$
(39)

We can again expand this result for $\xi R \ll 1$ to obtain

$$\mathcal{F} \approx -\gamma \pi R^2 - \alpha R F \left[\Gamma_e + \log\left(R\xi/2\right)\right] - \frac{F^2}{4\pi\gamma} \left[\Gamma_e + \log\left(L\xi/2\right)\right]. \tag{40}$$

B Review of the stochastic two-state process

We consider a binary state variable $x \in \{0, 1\}$ and denote by $P(x, t|x_0, t_0)$ the probability to be in state x at time t given a state x_0 at time t_0 . The probabilities obey

$$\mathbf{l} = P(1, t | x_0, t_0) + P(0, t | x_0, t_0)$$
(41)

and the probabilities evolve according to

$$\partial_t P(1,t|x_0,t_0) = -k_{01} P(1,t|x_0,t_0) + k_{10} (1 - P(1,t|x_0,t_0)).$$
(42)

The solution is for $P(x, t|x_0, t_0)$ is

$$P(x,t|x_0,t_0) = \frac{k_{10}\delta_{x,1} + k_{01}\delta_{x,0}}{k_{10} + k_{01}} + (\delta_{x,1} - \delta_{x,0})\frac{e^{-(k_{10}+k_{01})(t-t_0)}}{k_{10} + k_{01}}(k_{01}\delta_{x_0,1} - k_{10}\delta_{x_0,0}), \quad (43)$$

The steady-state expectation values are thus given by

$$\langle x \rangle_{ss} = \lim_{t \to \infty} \sum_{x=1,0} x P(x,t|x_0,0) = \frac{k_{10}}{k_{10} + k_{01}},$$
(44)

$$\sigma_x^2 = \langle x^2 \rangle_{ss} - \langle x \rangle_{ss}^2 = \langle x \rangle_{ss} - \langle x \rangle_{ss}^2.$$
(45)

The two-time correlations for $t \ge t'$ in steady state are given by

$$\langle x(t)x(t')\rangle_{ss} = \lim_{t''\to-\infty} \sum_{x,x',x''=1,0} xP(x,t|x',t')x'P(x',t'|x'',t'') = P(1,t|1,t')\langle x\rangle_{ss}$$

$$= \frac{k_{10}^2 + k_{10}k_{01}e^{-(k_{10}+k_{01})(t-t')}}{(k_{10}+k_{01})^2} = \langle x\rangle_{ss}^2 + \frac{k_{10}k_{01}e^{-(k_{10}+k_{01})(t-t')}}{(k_{10}+k_{01})^2}.$$

$$(46)$$

Next, we aim to calculate the likelihood of a given sequence of states and start by considering a single state change. Assuming that the state is initially given by $j \in \{0, 1\}$ with x = j at $t = t_0$, we are looking for the survival probability $G(\tau|j, t_0)$ that quantifies how likely it is that the system remains in the same state for a time τ . The evolution equation for the two-state process dictates that $G(\tau|j, t_0)$ obeys

$$\partial_{\tau}G(\tau|j,t_0) = -k_{ij}G(\tau|j,t_0). \tag{47}$$

Note that this equation also holds if the rates are time-dependent. Assuming constant rates, the differential equation yields $G(\tau|j, t_0) = e^{-k_{ij}\tau}$. We are now interested in the probability $p(\tau|j, t_0) d\tau$ that the state survives until τ and then changes in the infinitesimal time interval $[t + \tau, t + \tau + d\tau)$. Using the expression Eq. (47) for the rate of occurrence of the change, we have

$$p(\tau|j, t_0) \,\mathrm{d}\tau = -\partial_\tau G(\tau|j, t_0) \,\mathrm{d}\tau = k_{ij} e^{-k_{ij}\tau} \,\mathrm{d}\tau.$$
(48)

Next, let us consider a given sequence of changes with waiting times $\{\tau_1, \tau_2, \tau_3...\}$ starting at x = j. Using the probability density $p(\tau|j, t_0)$ derived above, the probability of finding this given sequence is

$$P(\{\tau_1, \tau_2, \tau_3 \dots\}) = k_{ij} e^{-k_{ij}\tau_1} d\tau_1 k_{ji} e^{-k_{ji}\tau_2} d\tau_2 k_{ij} e^{-k_{ij}\tau_3} d\tau_3 \dots$$
(49)

References

- G. Aquino, N. S. Wingreen, and R. G. Endres. Know the single-receptor sensing limit? think again. J. Stat. Phys., 162(5):1353, 2016.
- [2] K. Austen, P. Ringer, A. Mehlich, A. Chrostek-Grashoff, C. Kluger, C. Klingner, B. Sabass, R. Zent, M. Rief, and C. Grashoff. Extracellular rigidity sensing by talin isoform-specific mechanical linkages. *Nat. Cell Biol.*, 17(12):1597, 2015.
- [3] N. Bavi, Y. A. Nikolaev, O. Bavi, P. Ridone, A. D. Martinac, Y. Nakayama, C. D. Cox, and B. Martinac. Principles of mechanosensing at the membrane interface. In *The Biophysics* of Cell Membranes, page 85. Springer, 2017.
- [4] H. C. Berg and E. M. Purcell. Physics of chemoreception. *Biophys. J.*, 20(2):193, 1977.
- [5] S. M. Cahalan, V. Lukacs, S. S. Ranade, S. Chien, M. Bandell, and A. Patapoutian. Piezo1 links mechanical forces to red blood cell volume. *Elife*, 4:e07370, 2015.
- [6] B. Chen, B. Ji, and H. Gao. Modeling active mechanosensing in cell-matrix interactions. *Annu. Rev. Biophys.*, 44:1, 2015.
- [7] B. Coste, J. Mathur, M. Schmidt, T. J. Earley, S. Ranade, M. J. Petrus, A. E. Dubin, and A. Patapoutian. Piezo1 and piezo2 are essential components of distinct mechanically activated cation channels. *Science*, 330(6000):55, 2010.
- [8] C. D. Cox, C. Bae, L. Ziegler, S. Hartley, V. Nikolova-Krstevski, P. R. Rohde, C.-A. Ng, F. Sachs, P. A. Gottlieb, and B. Martinac. Removal of the mechanoprotective influence of the cytoskeleton reveals piezo1 is gated by bilayer tension. *Nat. Commun.*, 7(10366), 2016.
- [9] N. Dan and S. A. Safran. Effect of lipid characteristics on the structure of transmembrane proteins. *Biophys. J.*, 75(3):1410, 1998.
- [10] T. Effertz, B. Nadrowski, D. Piepenbrock, J. T. Albert, and M. C. Göpfert. Direct gating and mechanical integrity of drosophila auditory transducers require trpn1. *Nat. Neurosci.*, 15(9):1198, 2012.
- [11] A. Faucherre, J. Nargeot, M. E. Mangoni, and C. Jopling. piezo2b regulates vertebrate light touch response. J. Neurosci., 33(43):17089, 2013.
- [12] S. L. Geffeney, J. G. Cueva, D. A. Glauser, J. C. Doll, T. H.-C. Lee, M. Montoya, S. Karania, A. M. Garakani, B. L. Pruitt, and M. B. Goodman. Deg/enac but not trp channels are the major mechanoelectrical transduction channels in a c. elegans nociceptor. *Neuron*, 71(5):845, 2011.
- [13] B. Geiger, J. P. Spatz, and A. D. Bershadsky. Environmental sensing through focal adhesions. *Nat. Rev. Mol. Cell Biol.*, 10(1):21, 2009.
- [14] F. Guharay and F. Sachs. Stretch-activated single ion channel currents in tissue-cultured embryonic chick skeletal muscle. J. Physiol., 352(1):685, 1984.

- [15] J. Howard and S. Bechstedt. Hypothesis: a helix of ankyrin repeats of the nompc-trp ion channel is the gating spring of mechanoreceptors. *Curr. Biol.*, 14(6):R224, 2004.
- [16] H. A. Kramers. Brownian motion in a field of force and the diffusion model of chemical reactions. *Physica*, 7(4):284, 1940.
- [17] A. H. Lewis, A. F. Cui, M. F. McDonald, and J. Grandl. Transduction of repetitive mechanical stimuli by piezo1 and piezo2 ion channels. *Cell Rept.*, 19(12):2572, 2017.
- [18] Y.-X. Li and J. Rinzel. Equations for insp3 receptor-mediated [ca2+] oscillations derived from a detailed kinetic model: a hodgkin-huxley like formalism. J. Theo. Biol., 166(4):461, 1994.
- [19] G. Lorden. Procedures for reacting to a change in distribution. Ann. Math. Stat., 42(6):1897, 1971.
- [20] G. Maksaev, A. Milac, A. Anishkin, H. R. Guy, and S. Sukharev. Analyses of gating thermodynamics and effects of deletions in the mechanosensitive channel trek-1: comparisons with structural models. *Channels*, 5(1):34, 2011.
- [21] G. V. Moustakides. Optimality of the cusum procedure in continuous time. *Ann. Stat.*, 32(1):302, 2004.
- [22] E. S. Page. Continuous inspection schemes. *Biometrika*, 41(1-2):100, 1954.
- [23] R. Phillips, T. Ursell, P. Wiggins, and P. Sens. Emerging roles for lipids in shaping membrane-protein function. *Nature*, 459(7245):379, 2009.
- [24] M. Prager-Khoutorsky, A. Khoutorsky, and C. W. Bourque. Unique interweaved microtubule scaffold mediates osmosensory transduction via physical interaction with trpv1. *Neuron*, 83(4):866, 2014.
- [25] S. S. Ranade, R. Syeda, and A. Patapoutian. Mechanically activated ion channels. *Neuron*, 87(6):1162, 2015.
- [26] Y. Ritov. Decision theoretic optimality of the cusum procedure. *Ann. Stat.*, 18(3):1464, 1990.
- [27] B. Sabass and H. A. Stone. Role of the membrane for mechanosensing by tethered channels. *Phys. Rev. Lett.*, 116(25):258101, 2016.
- [28] U. S. Schwarz and S. A. Safran. Physics of adherent cells. Rev. Mod. Phys., 85(3):1327, 2013.
- [29] S. Sukharev and D. P. Corey. Mechanosensitive channels: multiplicity of families and gating paradigms. Sci. STKE, 2004(219):re4, 2004.
- [30] S. I. Sukharev, W. J. Sigurdson, C. Kung, and F. Sachs. Energetic and spatial parameters for gating of the bacterial large conductance mechanosensitive channel, mscl. J. Gen. *Physiol.*, 113(4):525, 1999.
- [31] P. R. ten Wolde, N. B. Becker, T. E. Ouldridge, and A. Mugler. Fundamental limits to cellular sensing. J. Stat. Phys., 162(5):1395, 2016.

- [32] S.-H. Woo, V. Lukacs, J. C. De Nooij, D. Zaytseva, C. R. Criddle, A. Francisco, T. M. Jessell, K. A. Wilkinson, and A. Patapoutian. Piezo2 is the principal mechanotransduction channel for proprioception. *Nat. Neurosci.*, 18(12):1756, 2015.
- [33] M. Yao, B. T. Goult, B. Klapholz, X. Hu, C. P. Toseland, Y. Guo, P. Cong, M. P. Sheetz, and J. Yan. The mechanical response of talin. *Nat. Commun.*, 7(11966), 2016.
- [34] R. Zarychanski, V. P. Schulz, B. L. Houston, Y. Maksimova, D. S. Houston, B. Smith, J. Rinehart, and P. G. Gallagher. Mutations in the mechanotransduction protein piezo1 are associated with hereditary xerocytosis. *Blood*, 120(9):1908, 2012.
- [35] W. Zhang, L. E. Cheng, M. Kittelmann, J. Li, M. Petkovic, T. Cheng, P. Jin, Z. Guo, M. C. Göpfert, L. Y. Jan, et al. Ankyrin repeats convey force to gate the nompc mechanotransduction channel. *Cell*, 162(6):1391, 2015.

D6 Mechanobiology of Animal Cells

R. Merkel Biomechanics Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Introduction		
2	Reaction of Cells to Substrate Stiffness		
	2.1	Phenomena	5
	2.2	Some underlying molecular processes	6
3	Possible Strain-sensing Mechanisms		
Ref	erence	es	9

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Within the organism all cells are constantly exposed to mechanical signals like strain or stiffness of the immediate environment, cf. Fig. 1. While this has been known for centuries, the role of mechanical signals for cell behavior and function has been mostly ignored for three reasons. First, within the body many biochemical signals are acting in parallel to mechanical signals and it is almost impossible to disentangle their influences. Second, because of the numerous uncertainties in the interpretation of experiments on living organisms most experiments on living cells are performed on cells in artificial environment, typically cells grown on either plastic or glass in an artificial medium containing essential chemical factors. Here it is extremely difficult to apply calibrated mechanical signals of physiological relevance and to quantify their impact. Third, the last reason is the structure of teaching. We are either taught in depth about mechanics (as students of mechanical engineering or at some places also of physics) or about cells (as students of biology or medicine) but rarely in both fields.



Fig. 1: Two examples for mechanical signals in the organism. Left: An anatomical drawing of the thorax region (taken from [1]). All tissues shown are mechanically highly active. Local strains in the heart reach 30% at frequencies of about 1 Hz, lung alveoli extend their inner area by about 100% during a deep breath and pressures during coughing can be extremely high. Right: A cut through the brain (taken from [2]). Here tissues of widely varying stiffness are seen. For example, modules of elasticity range from about 10 GPa (cortical bone), via ~10 MPa (cartilage, e.g. in ear and nose), 400 - 800 kPa (skin), 17 - 600 kPa (blood vessels), down to 200 - 1500 Pa for brain tissue. Note, reported strain and stiffness values vary widely and depend strongly on measurement techniques and preparation conditions.

It was only in the last decades that these shortcomings were rectified. A first hallmark experiment was done by Harris and coworkers who cultivated endothelial cells on thin membranes of heat crosslinked silicone and correctly interpreted the observed wrinkles as indicating strong contractile forces [3]. Unfortunately wrinkling is an inherently non-linear process which substantially complicated force retrieval. This problem was circumvented only almost 20 years later by Dembo and coworkers [4,5] who cultivated cells on layers of soft polyacrylamide with embedded marker beads to visualize cell-force induced substrate deformations. This technique, called traction force microscopy, still involves formidable numerical problems but is by now well established and used by many groups [6-11]. These works spurred experiments on cells grown on soft or microstructured substrates where most substantial changes of morphology, forces and cytoskeletal structures were observed. Especially notable observations are stiffness dependent phenotypes and differentiation of cells [12-16], attractiveness of sharp edges to cells - a phenomenon called contact guidance - [17,18], and cell reorientation upon external strain [19-25].

Over the years the study of the interplay of living animal cells with mechanical signals has evolved from mostly describing phenomena to the analysis of molecular mechanisms. By now it is a field of science in its own right and is often referred to as "mechanobiology". Please note: it is impossible to do justice to such a dynamic and fascinating field of science in a short article. The best the author can hope for is that this contribution might transport some of the fascination of the field and provides some starting points for reading.

2 Reaction of Cells to Substrate Stiffness

By now substrate stiffness is a well-established mechanobiological signal. For experiments with relevance to the in vivo situation substrate stiffness should be in the physiological range of 100 Pa (brain and glandular tissue) to some 100 kPa (e.g. tendon, extracellular matrix within bone). Biocompatible substrates in this stiffness range are most often produced from polyacrylamide or hyaluronic acid hydrogels [4,5,6,15,26]. Similar stiffnesses can be reached by silicone (PDMS) elastomers [7,8]. Both material classes exhibit viscoelastic response with silicones displaying much higher loss modules than typical hydrogels. Moreover, all substrates must be coated with proteins of the extracellular matrix (e.g. fibronectin, collagens or laminins) to be cell adhesive. Both aspects, protein coating and viscous loss, have been observed to influence cell behavior and should be considered for the interpretation of experimental results.

The easily visible reactions of animal cells to substrate stiffness include systematic variation of spreading area, polarization of the cytoskeleton (especially the actin cytoskeleton) and cellular adhesions plaques. Beyond these morphological changes, substantial variations of cytoskeletal tension and cell forces have been reported. The aforementioned cell reactions are reversible, that is, if the same cells are released from their substrate and brought onto another one of different stiffness, they will follow this change. However, there are also irreversible changes of cells as response to substrate stiffness. Reversible cell changes are called adaptation, irreversible ones differentiation. For example, upon experiencing substrate stiffness above about 10 kPa fibroblasts of the heart differentiate into myofibroblasts (see Fig. 2). This type of cellular reaction is an irreversible change of cell type. Remarkably, it has been shown that substrate stiffness plays its role in steering the differentiation of stem cells towards specific lineages [14].



Fig. 2: Fluorescence micrographs of fibroblasts prepared from hearts of rat embryos. Cells were cultivated on fibronectin coated silicone elastomer of 1 kPa stiffness (left) and 16 kPa stiffness (right). Cells were fixed and immunostained for α smooth muscle actin (α -SMA, green, top) and vinculin (red, bottom). Vinculin is a molecular marker for cell adhesions. Upon substrate stiffening cardiac fibroblasts differentiate into myofibroblasts. In the course of this process, α -SMA expression in massively increased. Moreover, cell size increases enormously and the actin cytoskeleton is reorganized to produce strong forces. Healthy heart tissue exhibits a stiffness below the threshold for differentiation (about 10 kPa) while scar tissue is stiffer. Thus this differentiation processes is initiated upon wounding, e.g., in heart strokes and acts to contract wounds [13].

Because cells of the connective tissue (i.e. fibroblasts) also secrete matrix proteins and rearrange the extracellular matrix they actively change the mechanical properties of their immediate environment. Therefore the aforementioned mechanobiological reactions to substrate stiffness give rise to a feedback cycle of substantial importance: Cells shape the micromechanics of their environment to which they react again by changes of cell behavior. This feedback enables adaptation of cellular and as a consequence also tissue mechanics to external signals like mechanical loads or scar formation.

3 Reaction of Cells to Substrate Strain

As shown in Fig. 1 almost every part of the human body is constantly exposed to mechanical strain. Some structures within the body are clearly pointing towards mechanical adaptation. Examples are structures in bones that follow mechanical strain within the bone or the alignment of endothelial cells along the direction of fluid flow in blood vessels. However, because in living organisms each tissue at each time point experiences a plethora of mechanical, electrical and chemical signals, it is almost impossible to disentangle their effects

and to identify cause-and-effect relations. Therefore one has to resort to in vitro experiments where cells are cultivated in devices designed for the application of mechanical signals while all other conditions remain as constant and physiological as possible.

3.1 Phenomena

In one of the first experiments on cell straining, fibroblasts were embedded in a hydrogel formed from collagen [27]. The cross-linked collagen material was cyclically stretched and the cell response observed. The authors report a substantial reorganization of the collagen network accompanied by a markedly peaked distribution of the orientations of the elongated cells. Here and in similar experiments these cell aligned along the stretch direction. However, in the end it remains unclear if cell reorientation is causing collagen fiber alignment or vice versa.

Therefore (and for reasons of experimental simplicity) an alternative experiment design was developed. Here cells are plated on lamellae of synthetic elastomer; the most frequently used material being silicone elastomer. On one hand these artificial substrates cannot be altered by cell activity and, on the other hand, local substrate strain can be calibrated quite accurately. This experimental design is mimicking physiological situations for cell types that grow on basal laminas like endothelial or epithelial cells. Here again cyclic stretching was necessary to induce clear-cut cell reactions. In contrast to the aforementioned collagen embedded cells, here cells reorient away from the stretch direction, i.e., here cells evade stretching by rotating away from strain direction [19-25]. Strain amplitude and repeat frequency strongly influence the final degree of reorientation and the kinetics of the build-up of orientational order. Obviously, all responses are strongly dependent on cell type and many details of the culture conditions. After cyclic straining cell orientations randomize again [28]. Concomitant with cell reorientation cytoskeletal structures, especially actin bundles, appear strengthened and better oriented (see Fig. 3).



Fig. 3: Fibroblasts from human umbilical cord cultivated on silicone lamellas. Cells were fixed and immunostained for filamentous actin. Left, control cells. Right, cells after 16 h cyclic strain (32% amplitude and 9 mHz repeat frequency). Stretch direction was horizontal. Note realingment and strengthening of fibers. This cell type orients towards the direction of zero strain which due to the Poisson ratio of the substrate material is not 90° to stretch but close to 70°. For experimental details see [24].

3.2 Some underlying molecular processes

In life cell microscopy on cells during stretch many processes are observed simultaneously. For example, some actin fiber bundles dissolve, others seem to split or rotate. Thus massive reorganization processes of these structures are occurring. Similar reorganizations happen in adhesion sites [29,30].

Please note that all such structures are multiprotein complexes where one set of structural proteins, like actin that forms force-bearing fibers, is interaction with another set of accessory proteins that modify the properties of "their" target structure. Several processes of this type are known by now; here we will just focus on one specific structure.

Upon stretch, even after only one short strain pulse, the accessory proteins zyxin and VASP bind to actin fiber bundles. These molecules bind localized along the fibers. These foci of protein accumulation appear as bright spots in life cell fluorescence microscopy (see Fig. 4). As only fibers oriented in stretch direction exhibit these foci, they are stretch dependent. Presumably zyxin and VASP bind to local nicks in the actin fibers where they initiate the polymerization of actin. This assumption is supported by the observation that these foci are transient and dissolve again once the actin fiber is repaired.

Taken together mechanical strain initiates a plethora of molecular processes involving not only binding of accessory proteins but also protein modification (like phosphorylation [31,32,33]) and protein degradation [34]. By which mechanisms these processes could be initiated will be discussed in the next section.



Fig. 4: Human umbilical vein endothelial cells (HUVECs) were grown on elastomeric substrates and stretched by a single strain pulse. A) Live cell microscopy on cells transfected with GFP-zyxin. B) and C) Stretched and GFP-zyxin transfected cells were fixed and additionally stained for VASP and actin (B) or vinculin and actin (C). Alexa633 phalloidin selectively labels actin fibers. Note that zyxin and VASP colocalize along stretched actin fibers whereas vinculin does not even though it is also an actin accessory protein.



Fig. 5: Schematic gating mechanism of MscL. Top: The molecule (red cylinder) is embedded in the bacterial membrane (blue) that is tensed (green arrows, membrane tension τ). An open channel (right) will occupy a larger area in the membrane than a closed one (left). This influences the free energy of the chemical reaction (lower row). Both closed (subscript c) and open (subscript o) state are local minima of the free energy. Without external membrane tension the closed state is exhibiting the lowest free energy and forms the ground state (blue line). An external membrane tension results in an additive term $-\tau A$ that enters the free energy (red straight line). At high enough membrane tension the open state of the channel will be the one of lowest free energy. That is, the chemical equilibrium between both states is steered by mechanical tension [40].

4 Possible Strain-sensing Mechanisms

As was described above mechanical stretching of cells causes very many molecular processes within cells. Thus there must be molecules that act as strain or stress sensors. At present three different mechanisms are discussed: Tension-induced opening of ion channels in the membrane [35,36], force-induced conformational changes of cytoskeletal elements [37], and force-induced unfolding of protein domains [38,39]. In the first case inflowing ions (e.g. Ca²⁺) act as secondary messengers to influence many intracellular processes. Force-dependent conformational changes of e.g. actin filaments enable load-dependent binding of accessory proteins that may either strengthen loaded structures or destabilize unloaded ones [37]. In the last case, unfolding of proteins, mechanical load liberates molecular sites that were formerly buried deep in their structure and therefore inaccessible. These sites may be either binding sites for accessory proteins, target sites for kinases (i.e. enzymes that covalently modify the protein by adding phosphate groups) or may even act as kinases themselves.

Here we will shortly discuss how mechanical force influences the kinetics of such processes. This has been most clearly demonstrated for the case of MscL, a bacterial channel protein that opens during osmotic shocks and acts figuratively spoken as "molecular safety vent" [40]. The mechanism is depicted in Fig. 5. This specific mechanism of mechanical steering of chemical equilibria has been originally proposed by Bell for the breakage of specific bonds [41] and later applied to experiments on mechanical breakage of single bonds [42,43] and on mechanical unfolding of proteins [44].

References

- H. Gray, Anatomy of the Human Body (Lea and Febiger, Philadelphia and New York, 1918)
- [2] W. Braune and C. Schmiedel, Topographisch-Anatomischer Atlas (Verlag von Veit & Comp., Leipzig, 1872)
- [3] A. P. Harris, P. Wild, and D. Stopak, Science 208, 177 (1980)
- [4] M. Dembo, T. Oliver, A. Ishihara, and K. Jacobson, Biophys. J. 70, 2008 (1996)
- [5] M. Dembo, and Y.-L. Wang, Biophys. J. 76, 2307 (1999)
- [6] J. Butler, I. Tolic-Norrelykke, B. Fabry, and J. Fredberg, Am J Physiol Cell Physiol 282, C595 (1991)
- [7] N. Balaban, U. Schwarz, D. Riveline, P. Goichberg, G. Tzur, I. Sabanay, D. Mahalu, S. Safran, A. Bershadsky, L. Addadi, and B. Geiger, Nat Cell Biol 3, 466 (2001)
- [8] C. Cesa, N. Kirchge
 ßner, D. Mayer, U. Schwarz, B. Hoffmann, and R. Merkel, Rev Sci Instrum 78, 034301 (2007)
- [9] R. Merkel, N. Kirchgeßner, C. M. Cesa, and B. Hoffmann, Biophys J 93, 3314 (2007)
- [10] B. Sabass, M. L. Gardel, C. M. Waterman, and U. S. Schwarz, Biophys J 94, 207 (2008)
- [11] X. Trepat, L. Deng, S. S. An, D. Navajas, D. J. Tschumperlin, W. T. Gerthoffer, J. P. Butler, and J. J. Fredberg, Nature 447, 592 (2007)
- [12] D. E. Disher, P. A. Janmey, and Y. L. Wang, Science 310, 1139 (2005)
- [13] J. M. Goffin, J. Pittet, G. Csucs, J. W. Lussi, J.-J. Meister, and B. Hinz, J Cell Biol 172, 259 (2006)
- [14] A. J. Engler, S. Sen, H. L. Sweeney, and D. E. Disher, Cell 126, 677 (2006)
- [15] F. Rehfeldt, A. E. X. Brown, M. Raab, S. Cai, A. L. Zajac, A. Zemel, D. E. Disher, Integr. Biol. 4, 422 (2012)
- [16] A. Schellenberg, S. Joussen, K. Moser, N. Hampe, N. Hersch, H. Hemeda, J. Schnitker, B. Denecke, Q. Lin, N. Pallua, M. Zenke, R. Merkel, B. Hoffmann, and W. Wagner, Biomaterials 35, 6351 (2014)
- [17] A. Curtis, and M. Riehle, Phys. Med. Biol. 46, R47 (2001)
- [18] M. Thery, J. Cell Sci. 123, 4201 (2010)
- [19] R. C. Buck, Exp. Cell Res. 127, 479 (1980)
- [20] P. C. Dartsch, and E. Betz, Basic Res. Cardiol. 84, 268 (1989)
- [21] T. Iba, and B. E. Sumpio, Mircrovasc. Res. 42, 245 (1991)
- [22] K. Hayakawa, N. Sato, and T. Obinata, Exp. Cell Res. 268, 104 (2001)
- [23] A. M. Goldyn, B. A. Roioja, J. A. Spatz, C. Ballestrem, and R. Kemkemer, J. Cell Sci. 122, 3644 (2009)
- [24] U. Faust, N. Hampe, W. Rubner, N. Kirchgessner, S. Safran, B. Hoffmann, and R. Merkel, PLoS One 6, 0028963 (2011)

- [25] C. Sears, and R. Kaunas, J. Biomech. 49, 1347 (2016)
- [26] T. Yeung, P. C. Georges, L. A. Flanagan, B. Marg, M. Ortiz, M. Funaki, N. Zahir, W. Ming, V. Weaver, and P. A. Janmey, Cell Mot. Cytoskel. 60, 24 (2005)
- [27] M. Eastwood, V. C. Mudera, D. A. McGrouther, and R. A. Brown, Cell Mot. Cytoskelet. 40, 13 (1998)
- [28] S. Jungbauer, H. Gao, J. P. Spatz, and R. Kemkemer, Biophys. J. 95, 3470 (2008)
- [29] D. Riveline, E. Zamir, N. Q. Balaban, U. S. Schwarz, T. Ishizaki, S. Narumiya, Z. Kam, B. Geiger, A. D. Bershadsky, J. Cell Biol. 153, 1175 (2001)
- [30] A. Carisey, R. Tsang, A. M. Greiner, N. Nijenhuis, N. Heath, N. Nazgiewicz, R. Kemkemer, B. Derby, J. Spatz, and C. Ballestrem, Curr. Biol. 23, 271 (2013)
- [31] Y. Sawada, M. Tamada, B. J. Dubin-Thaler, O. Cherniavskaya, R. Sakai, S. Tanaka, and M. Sheetz, Cell 127, 1015 (2006)
- [32] D. M. Donato, L. M. Ryzhova, L. M. Meenderink, I. Kaverina, and S. K. Hanks, J. Biol. Chem. 285, 20769 (2010)
- [33] V. Niediek, S. Born, N. Hampe, N. Kirchgessner, R. Merkel, and B. Hoffmann, Eur. J. Cell Biol. 91, 118 (2012)
- [34] A. Ulbricht, F. J. Eppler, V. E. Tapia, P. F. M. van der Ven, N. Hampe, N. Hersch, P. Vakeel, D. Stadel, A. Haas, P. Saftig, C. Behrends, D. O. Fürst, R. Volkmer, B. Hoffmann, W. Kolanus, and J. Höhfeld, Curr. Biol. 23, 1 (2013)
- [35] K. Hayakawa, H. Tatsumi, and M. Sokabe, J. Cell Sci. 121, 496 (2008)
- [36] A. J. Kuipers, J. Middelbeek, and F. N. van Leeuwen, Eur. J. Cell Biol. 91, 834 (2012)
- [37] K. Hayakawa, H. Tatsumi, and M. Sokabe, J. Cell Biol. 195, 721 (2011)
- [38] M. Rief, M. Gautel, F. Oesterhelt, J. M. Fernandez, and H. E. Gaub, Science 276, 1109 (1997)
- [39] X. Zhuang, and M. Rief, Curr. Op. Struct. Biol. 13, 88 (2003)
- [40] S. I. Sukharev, W. J. Sigurdson, C. Kung, and F. Sachs, J. Gen. Physiol. 113, 525 (1999)
- [41] G. I. Bell, Science 200, 618 (1978)
- [42] R. Alon, D. A. Hammer, T. A. Springer, Nature 374, 539 (1995)
- [43] R. Merkel, P. Nassoy, A. Leung, K. Ritchie, and E. Evans, Nature 397, 50 (1999)
- [44] M. Rief, M. Gautel, A. Schemmel, and H. E. Gaub, Biophys. J. 75, 3008 (1998)

D7 Cell and Tissue Mechanics

R. E. Leube Institute of Molecular and Cellular Anatomy Uniklinik RWTH

Contents

Introduction2				
1	The	mechanics of single cell migration	2	
2	Epi	thelia: Ramparts against mechanical stress	6	
	2.1	Single-layered polarized intestinal epithelium	7	
	2.2	Multilayered epidermal epithelium	8	
3	Heart muscle: The challenges of permanent contraction cycles1			
4	Summary and conclusions			
Ref	erence	·s	13	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

Introduction

The goal of this lecture is to illustrate, why an understanding of the force balance acting on cells and tissues is important for elucidating mechanisms determing physiology and pathology. The three selected examples highlight different aspects of this theme. Each is of medical relevance and focuses on the contribution of the cytoskeleton and specialized adhesion sites.

First, we will elaborate on cell migration as a cardinal cellular trait that is an intricate feature of developmental processes, immune responses, and tumor cell metastasis. We will concentrate on single cells, that are currently best characterized, and describe an iterative cycle involving the force-generating machinery, which is dynamically linked to the extracellular substrate and involves pulling, pushing and friction. Second, we will describe the keratin-desmosome scaffold of epithelial cells. We will zoom in (*i*) on the subapical network in intestinal cells that is conserved from *C. elegans* to human, and (*ii*) on the 3D network in epidermal cells with its diversified stratum-, context-, and function-dependent molecular composition. We will present evidence, that these systems fullfil important architectural and mechanical functions in their native tissue contexts. Consequently, mutations of the encoded genes result in reduced mechanical resilience leading to pronounced tissue perturbations such as cytoplasmic invaginations and blister formation. Third, we will portray a rare type of cardiomyopathy that has been linked to disruption of cytoskeleton-junction coupling. We will show why it is necessary to understand the mechanical dysfunction for elucidating the still unknown pathomechanism of the disease.

Taken together, we conclude and propose that genetics and chemical microenvironment alone are not sufficient to understand cell and tissue function and that the mechanical microenvironment is just as important in determining cellular function, fate and dysfunction.

1 The mechanics of single cell migration

The aim of this part is to describe single cell migration as a paradigm highlighting the effect of forces on single cells. Several of the cellular components and molecular processes have been studied at much detail, since they are accessible to *in vitro* examination in cultured cells.

Migration is a fundamental property of cells. It allows directed movement in response to changes in the microenvironment supporting growth, differentiation and regeneration and, at the same time, avoiding adverse conditions. The cues driving migratory behaviour include chemical signals such as nutrients, growth factors, and toxins or physical signals such as radiation as well as mechanical (gravitational) and electromagnetic forces. Thus, one can distinguish between chemotaxis caused by chemoattractans or morphogens, haptotaxis due to varying substrate concentrations, mechanotaxis because of cell contact breakdown, electrotaxis induced by electric fields and durotaxis caused by differences in substrate stiffness. These cues induce movement of cells in relation to their environment. This requires mechanical forces involving a controllable force-generating machinery that must be linked not only to counterbalances within the cell but also to counterbalances in the outside environment. The force-generating machinery corresponds to the acto-myosin system, whose contractile activity is regulated by inducible signalling pathways. The cytoplasmic counterbalance is provided by the cytoskeleton, which is composed of three major filament systems, each with unique biomechanical, structural, chemical and functional properties

(Table 1). The extracellular counterbalance corresponds to the extracellular matrix. Its biomechanical properties vary considerably through different admixtures of its three major components, i.e. the filamentous component consisting of collagen and elastic fibres, the organic non-filamentous component containing gel-forming glycosaminoglycans, and the anorganic component including calciumphosphate, which contributes to extracellular matrix stiffness. The intra- and extracellular counterbalances are coupled through specific cell-extracellular contacts (Table 2). In addition, the cytoplasmic counterbalances are coupled across cell borders through specialized cell-cell contacts (Table 2).

Major cytoplasmic filament systems in human								
	Microtubules	Intermediate Filaments	Actin Filaments					
Diameter	24 nm	~10 nm	7 nm					
Molecular components	9 α-tubulins 10 β-tubulins	> 70 intermediate filament poly- peptides [epithelia: 28 type I keratins (K9-40) and 26 type II keratins (K1-8, K71- 86)]	6 actins [cytoplasmic: ACTB (cell cortex), ACTG1 (stress fibres); muscle- specific: ACTA1 (skeletal muscle), ACTC1 (cardiac muscle), ACTA2 / ACTG2 (smooth muscle)]					
Enzymatic activity	GTPase	None	ATPase					
Building principles	Heterodimers -> 13 proto-filaments	Keratins: heterodimer -> tetramer -> unit length filament -> 10 nm filament -> non-nolar	Double helix with 14 subunits per turn					
De novo formation	Nucleation (microtubule organizing center) -> elongation	Nucleation (spontaneous?) -> elongation (unit length filament	Nucleation -> elongation					
Assembly / Disassembly	Dynamic instability	Lateral subunit exchange - excision / insertion / annealing - turnover cycle	Treadmilling / retrograde flow - branching (Arp2/3) - plus-end growth (formins)					
Mechanics	Persistence length: ~1 mm Young's modulus: 1 000-1 500 MPa -> bear high compressive forces -> generate pushing and pulling forces	Persistence length: ~ 1 μm Young's modulus: 6-300 MPa -> elastic deformation at low mechanical load -> plastic deformation at high mechanical load	Persistence length: ~10 μm Young's modulus: 1 800-2 500 MPa -> semiflexible + myosin: biological active springs					
Filament localization	- Radial (centrosomal) - Non-centrosomal	- Pancytoplasmic - Perinuclear - Subplasmalemmal	 Filopodial Lamellipodial/lamellar Dorsal arcs Ventral / dorsal stress fibres Dense bodies (smooth muscle) Sarcomeric (striated muscle) 					
Filament organization	Single filaments - filament bundles - 3D network							
Superstructures	- Cilium, flagellum - Centrosome/basal body	 Keratohyalin granule Cornified envelope 	- Microvillus - Lamellipodium / filopodium					
Functions of identified associated proteins	Nucleating Polymerization-inhibiting Solubilizing / destabilizing Stabilizing / co-filamentous Gel-forming / cross-linking Bundling Severing (+) end capping / (-) end capping ATPase (motor protein)	Polymerization-inhibiting Gel-forming / cross-linking Membrane-attaching	Nucleating Polymerization-inhibiting Solubilizing / destabilizing Stabilizing / co-filamentous Gel-forming / cross-linking Bundling Severing (+) end capping / (-) end capping Membrane-attaching ATPase (motor protein)					
Motor proteins	Kinesins, dyneins	None	Myosins					
Posttranslational modifications	Phosphorylation, acetylation, sumoy Arginylation, detyrosination, polyglycylation, poly- glutamylation, palmitoylation	lation, ubiquitination, oxidation, glycosyla Transamidation	ation, methylation, ADP ribosylation Arginylation, nitrosylation					
Depolymerizing / stabilizing drugs	Colchicine, nocodazole, vinblastin, demecolcine / taxol	None	Cytochalasins, latrunculin / jasplakinolide, phalloidin					

Table 1

Composition of major load-bearing junctions in epithelia								
	Major transmembrane proteins	Major linker proteins	Associated cytoskeletal filaments					
Desmosome (Macula adhaerens)	Desmosomal cadherins (desmogleins and desmocollins)	Plakoglobin, plakophilins and desmoplakins	Keratin intermediate filaments					
Adherens junction	"Classical" cadherins, nectins	Catenins (p120, β , α), vinculin, afadin	Actin filaments					
Hemidesmosome	Integrins ($\alpha 6/\beta 4$), BP180	BP230, plectin	Keratin intermediate filaments					
Focal adhesion	Integrins (multiple α/β isoforms)	Talin, vinculin, kindlin	Actin filaments					

Table 2

Crawling is the predominant form of cell translocation in animals. Different types of crawling can be distinguished including amoeboid (fast, weak adhesion, poorly developed cytoskeleton) and mesenchymal/epithelial movement (slow, strong adhesion, highly developed cytoskeleton) of single cells and collective migration of cell groups arranged as chains or sheets [1, 2]. Polarization initiates cell migration defining front and rear of cells through partitioning of regulatory molecules leading to the formation of a tail at the back (uropod) and a leading edge in the front (lamellipod). The subsequent canonical mode of mesenchymal/epithelial single cell migration on a flat 2D surface consists of repetitive cycles encompassing (*i*) protrusion, (*ii*) adhesion, (*iii*) cell body translocation, and (*iv*) retraction [3-5]; Fig. 1). Each step involves changes in the local force equilibrium.

- (i) Protrusion is defined as leading edge extension and starts with filopodial exploration. Filopodia are thin finger-like protrusions (diameter: ~50 nm; length: several μm) of the plasma membrane at the leading edge each containing 15-30 bundled actin filaments [6, 7]. Polymerization/depolymerization of actin filaments at the tip of filopodia determine extension/retraction of these structures which is coordinated by formins. The spike-like filopodia are typically associated with the lamellipodium, the outermost part of the very flat sheet-like part of the cell's leading edge, which is referred to as the lamellum. It is characterized by membrane ruffling that is driven by a branched actin network. The ARP2/3 complex facilitates formation of network branches and fixes them at a 70° angle [8]. Net filament assembly at the leading edge and net filament disassembly behind the leading edge together support protrusion of the substratum. Forces are generated through attachment of actin stress fibres to focal adhesions by a clutch in a ratchet mechanism [9]. Focal adhesions, in turn are coupled to the extracellular matrix (step ii).
- (ii) Adhesion is primarily mediated through focal adhesions, which link the actin cytoskeleton to the extracellular matrix. New focal adhesion sites are constantly generated at the leading edge. They go through a process of maturation that is reflected by compositional and structural alterations. Nascent focal adhesions are formed in the transition zone between the lamellipodium and the lamellum forming small focal complexes that become larger and elongated focal adhesions under the influence of forces imposed by associated actin stress fibres and extracellular matrix components [10]. The mechanical coupling, which is subject to regulation, is referred to as the mechancial clutch [9]. Focal adhesions subsequently mature into fibrillar adhesions [11].
- (iii) Cell body translocation is facilitated by contraction at the rear end through activation of the actomyosin system involving myosin II activity, which propels the cell body forward [3-5].
(iv) Retraction is coupled to rear end release via de-adhesion [3-5]. The necessary disassembly of extracellular matrix contacts, however, is not complete. As a consequence long retraction fibres are formed, which are eventually ripped off leaving behind tracks of remnant membrane fragments, which remain attached to the extracellular matrix and contain integrin adhesion molecules.



Fig. 1: Schematic representation of the major steps in single cell migration on a flat 2D surface.

The above processes are coupled to and coordinated by **mechanosensing** and **mechanotransduction** [12-14]. Mechanosensors can transform mechanical signals into biochemical information. The most relevant mechanism for cell migration inolves force-dependent conformational alterations, which open binding domains and enzymatically active domains. The proteins are typically characterized by a modular structure. Forces impose sequential unfolding of tertiary structures within each module. Mechanosensors are integral

parts of the cytoskeleton-extracellular matrix scaffold. Thus, forces facilitate, e.g.: integrin binding of the extracellular matrix protein fibronectin, focal complex maturation by activation of the focal adhesion protein talin, induction of signaling cascades by the focal adhesion protein p130Cas, actin stress fibre formation by the focal adhesion protein zyxin, and cortical actin stabilization through filamin.

Our current knowledge in cell migration is, for the most part, restricted to the actin cytoskeleton and its associated focal adhesions. The precise role of the other cytoskeletal systems is much less understood. Microtubules serve an important function in polarized transport processes needed for membrane and cytoplasmic extension [15]. Intermediate filaments, on the other hand, have been shown to affect cell migration in different ways depending on molecular composition and context [16-18]. Similarly, the role of hemidesmosomes, which anchor intermediate filaments to the extracellular matrix (Table 2), and their interplay with focal adhesions still remains to be elucidated [19].

Cell migration is important for many physiologically occurring processes and pathologies. During embryogenesis, organogenesis and regeneration undifferentiated precursor cells are directed to distant locations. During wound healing and immune responses environmental signals inform cells to move. Migration of the wrong cell type to the wrong place is encountered in pathology with catastrophic effects on tissue homeostasis occurring in autoimmune diseases and during metastasis in carcinogenesis.

2 Epithelia: ramparts against mechanical stress

The aim of this part is to describe the keratin intermediate filament-desmosome system as a tissue- and cell-type specific scaffold providing epithelial resilience and supporting epithelial barrier function. Epithelia are exposed interfaces between the outside environment and the body and are therefore subjected to mechanical stress. At the same time, epithelia serve as large surfaces for the bidirectional exchange of molecules between the external and internal milieu by resorption and secretion. Depending on the specific mechanical and functional requirements epithelia encompass not only different cell types but are also distinguished by different tissue architectures ranging from polarized single cell layers to multilayered assemblies (Figs. 2, 3).

Keratin intermediate filaments are major components of the epithelial cytoskeleton. They are attached to desmosomal cell-cell adhesion sites and hemidesmosomal cell-extracellular matrix adhesions. Together they fullfil specific biomechanical tasks within the different epithelia. This is reflected by cell- and epithelial tissue type-specific patterns of keratin and desmosomal protein isoform expression and arrangements [20, 21]. By paradigmatically describing the situation in the single-layered intestinal epithelium and the multilayered epidermal epithelium we will highlight some of the features of the keratin-desmosome scaffold in functional tissue contexts.



Fig. 2: Scheme depicting organizational aspects of the cytoskeleton and cell adhesions in the one-layered polarized intestinal epithelium. The top part emphasizes the subapical enrichment of ntermediate filaments that are anchored to cell-cell junctions. The bottom part presents features of all three major cytoskeletal filament systems (intermediate filaments, actin filaments, microtubules) that are anchored to junctions through plaque proteins, crosslinked to each other by cytolinkers, enriched at the apical domain together with the polarity complex, delivered to specific regions by motor proteins, and co-distributed through nucleation sites (γ -tubulin ring complex). Further details in [22].

2.1 Single-layered polarized intestinal epithelium

Intestinal epithelia are prototypic single-layered polarized epithelia, which face the body interior at their basal surface and the body exterior at their apical suface. The function of the intestinal epithelium is to facilitate regulated exchange of molecules between the outside, i.e. the intestinal lumen, and the inside of the body while maintaining an intact and resilient barrier between both compartments. Intestinal cells produce a distinct subset of keratin polypeptides including the type II keratins K7 and K8 and the type I keratins K18, K19 and K20 (cf. [22]). K8 and K19 are synthesized throughout the epithelium of the small and large intestine. But the other keratins show a more restricted distribution with K18 and small amounts of K7 primarily in the undifferentiated crypt compartment and K20 in the villus. Differential distribution of desmosomal proteins along the crypt-villus axis has also been reported for desmosomal proteins [23]. It will be interesting to find out how the differing molecular composition relates to differences in local cellular specialization and properties, whereby cells differentiate and move from the crypt to the villar tip where they are constantly shed into the intestinal lumen.

A remarkable architectural feature of the intestinal keratin cytoskeleton is its concentration below the adluminal plasma membrane separating the apical organelle-free terminal web region containing the microvillar rootlets of bundled actin filaments from the rest of the cytoplasm (Fig. 2). This arrangement is conserved from human all the way to the nematode C. *elegans*, in which six intestine-specific intermediate filament polypeptides form a very dense subapical network that is anchored to the C. elegans apical junction, a multicomponent cellcell adhesion complex [22, 24]. Local perturbations of the network or collapse of the network towards the junction leads to luminal widening and cytoplasmic invaginations of the apical plasma membrane [25, 26]. These phenotypes tend to aggravate with increasing age. The most likely explanation is reduced mechanical stability of the intermediate filament-rich structure, which succumbs to chronic wear and tear. This hypothesis has been tested on vital intestine. The intestine can be easily prepared from living worms by a single cut, which results in its extrusion. Using a dual micropipette assay the dissected intestine is then accessible for stress-strain analyses. This revealed a remarkable resilience of the intestine along the longitudinal axis preventing ruptures at applied forces of up to 0.37 μ N [27]. It further showed that increased strain is observed at high forces of mutant intestines (own unpublished results).

2.2 Multilayered epidermal epithelium

The epidermis consists of multiple layers with stratum-specific differentiation features (Fig. 3). The basal layer consists of cylindrical cells that are attached to the basement membrane through hemidesmosomes. The cells divide to maintain and replenish not only the basal but also all suprabasal compartments. The first suprabasal compartment is referred to as the spinous layer. It is no longer in contact with the basement membrane and is characterized by increasingly flat cells that are circumferentially surrounded by abundant desmosomes. The next layer is the stratum granulosum, which contains prominent cytoplasmic granules. The top layers are formed upon programmed cell death (apoptosis) coupled to the formation of the cross-linked envelope, which serves as a protective barrier. The various differentiation states are reflected by different keratin expression profiles [21]. Keratins K5 and K14 are obligatory components of basal keratinocytes. In addition, K15 is often detectable with predilection for stem cells, most notably in the hair follicle bulge. On the other hand, K1 and K10 are usually found in suprabasal cells. They are complemented by K2 and K24 in the upper spinous and

granular layer [28] and by keratin 78 in the basal and first suprabasal layer [29]. Of note K9 is only produced in the mechanically most challenged palmoplantar epidermis. Other keratins such as K6, K16 and K17 are absent in the normal interfollicular epidermis of hairy skin but are switched on in certain situations, notably upon wounding, which requires movement of keratinocytes towards the lesioned area and keratinocyte proliferation. It is generally assumed that the different keratins endow cells with specific biomechanical properties. Thus, the different keratin pairs exhibit different biochemical properties (polymer stability, end domain composition) and differ in their propensity to bundle and to associate with desmosomes (e.g., [30, 31]).

Desmosomes are highly abundant in the epidermis occupying the major part of the surface of suprabasal cells and increasing in size towards the uppermost layers, where they form large corneodesmosomes. Compositionally, desmosomes also differ in the various epidermal layers thereby further fine-tuning the properties of the epidermal keratin-desmosome system [20, 32]. Taken together, the cytoplasmic keratin network in conjunction with its desmosomal anchorage sites form a transcellular network that is believed to confer epidermal resilience and plasticity. We recently proposed that this system is based on an ordered arrangement of the keratin network consisting of perinuclear filaments that are linked through radial spokes to desmosomes, which are, in turn, connected to each other by suplasmalemmal keratin filaments (the rim; [33]).

The most convincing evidence for a mechanical function of the keratin-desmosome scaffold stems from human blistering diseases that have been linked to desmosomal proteins and intermediate filament proteins [20, 31, 34, 35]. As an example, we will describe features of epidermolvsis bullosa simplex, which is an autosomal dominant disease that has been linked to single point mutations in the genes encoding K5 and K14 [31, 36-38]. Patients develop blisters in mechanically challenged regions. The blisters are restricted to the basal cell layer, where cytolysis occurs. This is linked to the formation of prominent keratin aggregates. The most common assumption is that the keratin point mutations encode keratin polypeptides acting in a dominant-negative fashion on keratin filament network assembly. In accordance, transfection of these mutant keratin has been shown by many laboratories to induce abundant granule formation and loss of keratin filaments. It is further assumed that dysfunctional keratin networks lead to the increased sensitivity of basal keratinocytes to mechanical trauma resulting in cytolysis and blister formation. While the simplicity of this argument and the copious evidence are overwhelming and highly seductive, the pathogenesis may be more complicated. Thus, (i) some of the most severe mutations do not prevent keratin filament formation in vitro, (ii) epidermolysis bullosa simplex patients develop normal-appearing keratin-desmosome scaffolds in non-traumatized regions, and (iii) patients with mutations in suprabasal keratins do not primarily develop blisters but present with hyperkeratosis. As an alternative mechanism, we have recently shown that phophomimetic mutation of a single residue is sufficient to prevent keratin filament network formation [39]. Therefore, altered mechanosensing and mechanotransduction likely contribute to the full penetrance of the epidermolysis bullosa simplex phenotype.



Fig. 3: Scheme depicting organizational aspects of the keratin cytoskeleton and its desmosomal/hemidesmosomal anchorage sites in the multilayered epidermal epithelium at steady state and upon increased pressure under normal circumstances (middle lower panel) or in the presence of keratin mutations causing epidermolysis bullosa simplex (right lower panel). Lack of keratin-mediated desmosomal connectivity results in cytolysis and subsequent blister formation [further details in [33]].

3 Heart muscle: The challenges of permanent contraction cycles

As a last example, we highlight the importance of cell adhesion and the cytoskeleton in the heart. The heart is unique by performing incessant cycles of contraction and relaxation throughout life at regular intervals. Thus, an 80 year-old heart beating at an average rate of 60 times per minute has gone through $\sim 2.5 \times 10^9$ precisely timed and coordinated contraction/relaxation cycles. This is made possible by a highly ordered arrangement of the contractile apparatus that is functionally and mechanically coupled between adjacent cardiomyocytes (Fig. 4). Specialized cell-cell contacts, i.e. the intercalated discs, are masterpieces of efficiency to realize the necessary electromechanical coupling. The intercalated disc connects neighbouring contractile and excitable cardiomyocytes by a complex arrangement of different cell-cell junctions, which are tightly interwoven within this superstructure. Stripe-like *fasciae adhaerentes* serve as anchorage sites for the contractile acto-myosin apparatus. Dot-like desmosomes facilitate attachment of desmin intermediate filaments, which enwrap the contractile system. Gap junctions (*nexus*) form transcellular channels and mediate electrical coupling.

Arrhythmogenic cardiomyopathy (AC) is a rare cardiomyopathy that has been linked to mutations of all known desmosomal proteins that are synthesized in the heart [40, 41]. Arrhythmias are the first symptoms which can lead to sudden death, occurring occasionally during endurance sports. In most cases, however, foci with necrotic cardiomyocytes appear during the acute disease stage that are replaced by fibrofatty tissue. During the chronic disease phase dilative cardiomyopathy is most prominent. The pathology is often restricted to the right ventricle but may also affect the left ventricle. Heart failure may eventually necessitate heart transplantation. Different disease mechanisms have been discussed, none of which has been conclusively affirmed by the scientific community [40, 41].

An idea favored by several researchers is that a dysfunction in adhesion triggers the disease [42-44]. In support, a reduced number of desmosomes and widened intercellular gaps have been observed in tissue samples of AC patients and murine AC models which were interpreted as indications of reduced adhesion strength [42, 43, 45]. Given that desmin mutations also induce an AC phenotype (c.f. [41]), it will be important to examine the biomechanical function of the entire cardiac desmin intermediate filament-desmosome scaffold. It may be a safeguarding system, which keeps the ordered sarcomeric structure and cellular organelles such as mitochondria in place because of its intrinsic elasticity and recoil activity. Obviously, the desmosomal protein mutations do not interfere directly with cardiac contractility and function, since homozygous mutant mice carrying the disease-causing mutations still develop normal-appearing hearts. As mechanical load increases, however, the disease develops postnatally. Thus, focal cardiomyocyte necrosis is observed during adolescent growth. The necrosis elicits an aseptic inflammatory reaction which eventually leads to scar formation and thereby alters extracellular matrix stiffness. This pathology may be the consequence of acute and localized excessive mechanical stress leading to disruption of cardiomyocyte connectivity. Later on, cardiac dilation is the most prominent feature that is coupled to a reactive hypertrophic response [46]. Major difficulties in unravelling the underlying molecular mechanisms are the lack of appropriate *in vitro* model systems and the relative inaccessibility of the actively moving cardiac muscle in vivo. Elucidation of molecular architecture of the cytoskeleton and its associated anchorage sites at cell-cell contacts (intercalated disc) and cell-extracellular matrix contacts (costamere) at high resolution in 3D, biomechanical probing of cardiac tissue slices, and functional *in vivo* imaging, however, may help to improve our understanding of the altered mechanics in the AC heart.



Fig. 4: Scheme of the intercalated disc as an integrator of mechanical forces and signal transducer between adjacent cardiomyocytes.

4 Summary and conclusions

The three examples were selected to highlight the importance of cellular mechanical stability in single cells and, more importantly, in differentiated tissue. We wanted to emphasize that the cytoskeleton and its associated adhesion sites are crucial prerequisites for cellular mobility on the one hand and maintenance of tissue integrity on the other hand. Our knowledge of the force distribution at the cellular and subcellular level is still rather limited but it is crucial to elucidate the contribution of specific molecules. A fascinating challenge is to understand how forces are sensed and translated into cellular responses. The advent of highly sensitive devices to measure and locally modify the force equilibrium herald exciting discoveries in the growing field of mechanobiology. These discoveries will help to elucidate pathomechanisms in numerous diseases and open novel avenues for the treatment of various pathologies.

Acknowledgement. I would like to thank Adam Breitscheidel for preparing the figures and members of my laboratory for fruitful discussions.

References

- Friedl P, Gilmour D. Collective cell migration in morphogenesis, regeneration and cancer. Nat Rev Mol Cell Biol. 2009;10:445-57.
- [2] Friedl P, Locker J, Sahai E, Segall JE. Classifying collective cancer cell invasion. Nat Cell Biol. 2012;14:777-83.
- [3] Lauffenburger DA, Horwitz AF. Cell migration: a physically integrated molecular process. Cell. 1996;84:359-69.
- [4] Sheetz MP, Felsenfeld DP, Galbraith CG. Cell migration: regulation of force on extracellular-matrix-integrin complexes. Trends Cell Biol. 1998;8:51-4.
- [5] Le Clainche C, Carlier MF. Regulation of actin assembly associated with protrusion and adhesion in cell migration. Physiol Rev. 2008;88:489-513.
- [6] Svitkina TM, Bulanova EA, Chaga OY, Vignjevic DM, Kojima S, Vasiliev JM, et al. Mechanism of filopodia initiation by reorganization of a dendritic network. J Cell Biol. 2003;160:409-21.
- [7] Wood W, Martin P. Structures in focus--filopodia. Int J Biochem Cell Biol. 2002;34:726-30.
- [8] Pollard TD. Regulation of actin filament assembly by Arp2/3 complex and formins. Annu Rev Biophys Biomol Struct. 2007;36:451-77.
- [9] Aratyn-Schaus Y, Gardel ML. Biophysics. Clutch dynamics. Science. 2008;322:1646-7.
- [10] Schwarz US, Gardel ML. United we stand: integrating the actin cytoskeleton and cellmatrix adhesions in cellular mechanotransduction. J Cell Sci. 2012;125:3051-60.
- [11] Gardel ML, Schneider IC, Aratyn-Schaus Y, Waterman CM. Mechanical integration of actin and adhesion dynamics in cell migration. Annu Rev Cell Dev Biol. 2010;26:315-33.
- [12] Ohashi K, Fujiwara S, Mizuno K. Roles of the cytoskeleton, cell adhesion and rho signalling in mechanosensing and mechanotransduction. J Biochem. 2017;161:245-54.
- [13] Hatzfeld M, Keil R, Magin TM. Desmosomes and Intermediate Filaments: Their Consequences for Tissue Mechanics. Cold Spring Harb Perspect Biol. 2017;9.
- [14] Hoffman BD, Grashoff C, Schwartz MA. Dynamic molecular processes mediate cellular mechanotransduction. Nature. 2011;475:316-23.
- [15] Etienne-Manneville S. Microtubules in cell migration. Annu Rev Cell Dev Biol. 2013;29:471-99.
- [16] Chung BM, Rotty JD, Coulombe PA. Networking galore: intermediate filaments and cell migration. Curr Opin Cell Biol. 2013;25:600-12.
- [17] Sanghvi-Shah R, Weber GF. Intermediate Filaments at the Junction of Mechanotransduction, Migration, and Development. Front Cell Dev Biol. 2017;5:81.
- [18] Cheng F, Eriksson JE. Intermediate Filaments and the Regulation of Cell Motility during Regeneration and Wound Healing. Cold Spring Harb Perspect Biol. 2017;9.

- [19] Tsuruta D, Hashimoto T, Hamill KJ, Jones JC. Hemidesmosomes and focal contact proteins: functions and cross-talk in keratinocytes, bullous diseases and wound healing. J Dermatol Sci. 2011;62:1-7.
- [20] Holthofer B, Windoffer R, Troyanovsky S, Leube RE. Structure and function of desmosomes. Int Rev Cytol. 2007;264:65-163.
- [21] Moll R, Divo M, Langbein L. The human keratins: biology and pathology. Histochem Cell Biol. 2008;129:705-33.
- [22] Coch RA, Leube RE. Intermediate Filaments and Polarization in the Intestinal Epithelium. Cells. 2016;5.
- [23] Patey N, Scoazec JY, Cuenod-Jabri B, Canioni D, Kedinger M, Goulet O, et al. Distribution of cell adhesion molecules in infants with intestinal epithelial dysplasia (tufting enteropathy). Gastroenterology. 1997;113:833-43.
- [24] Carberry K, Wiesenfahrt T, Windoffer R, Bossinger O, Leube RE. Intermediate filaments in Caenorhabditis elegans. Cell Motil Cytoskeleton. 2009;66:852-64.
- [25] Carberry K, Wiesenfahrt T, Geisler F, Stocker S, Gerhardus H, Uberbach D, et al. The novel intestinal filament organizer IFO-1 contributes to epithelial integrity in concert with ERM-1 and DLG-1. Development. 2012;139:1851-62.
- [26] Geisler F, Gerhardus H, Carberry K, Davis W, Jorgensen E, Richardson C, et al. A novel function for the MAP kinase SMA-5 in intestinal tube stability. Mol Biol Cell. 2016;27:3855-68.
- [27] Jahnel O, Hoffmann B, Merkel R, Bossinger O, Leube RE. Mechanical Probing of the Intermediate Filament-Rich Caenorhabditis Elegans Intestine. Methods Enzymol. 2016;568:681-706.
- [28] Min M, Chen XB, Wang P, Landeck L, Chen JQ, Li W, et al. Role of keratin 24 in human epidermal keratinocytes. PLoS One. 2017;12:e0174626.
- [29] Langbein L, Eckhart L, Fischer H, Rogers MA, Praetzel-Wunder S, Parry DA, et al. Localisation of keratin K78 in the basal layer and first suprabasal layers of stratified epithelia completes expression catalogue of type II keratins and provides new insights into sequential keratin expression. Cell Tissue Res. 2016;363:735-50.
- [30] Loschke F, Homberg M, Magin TM. Keratin Isotypes Control Desmosome Stability and Dynamics through PKCalpha. J Invest Dermatol. 2016;136:202-13.
- [31] Homberg M, Magin TM. Beyond expectations: novel insights into epidermal keratin function and regulation. Int Rev Cell Mol Biol. 2014;311:265-306.
- [32] Rubsam M, Broussard JA, Wickstrom SA, Nekrasova O, Green KJ, Niessen CM. Adherens Junctions and Desmosomes Coordinate Mechanics and Signaling to Orchestrate Tissue Morphogenesis and Function: An Evolutionary Perspective. Cold Spring Harb Perspect Biol. 2017.
- [33] Quinlan RA, Schwarz N, Windoffer R, Richardson C, Hawkins T, Broussard JA, et al. A rim-and-spoke hypothesis to explain the biomechanical roles for cytoplasmic intermediate filament networks. J Cell Sci. 2017;130:3437-45.
- [34] Szeverenyi I, Cassidy AJ, Chung CW, Lee BT, Common JE, Ogg SC, et al. The Human Intermediate Filament Database: comprehensive information on a gene family involved in many human diseases. Hum Mutat. 2008;29:351-60.
- [35] Waschke J. The desmosome and pemphigus. Histochem Cell Biol. 2008;130:21-54.

- [36] Chamcheu JC, Siddiqui IA, Syed DN, Adhami VM, Liovic M, Mukhtar H. Keratin gene mutations in disorders of human skin and its appendages. Arch Biochem Biophys. 2011;508:123-37.
- [37] Sprecher E. Epidermolysis bullosa simplex. Dermatol Clin. 2010;28:23-32.
- [38] Coulombe PA. The Molecular Revolution in Cutaneous Biology: Keratin Genes and their Associated Disease: Diversity, Opportunities, and Challenges. J Invest Dermatol. 2017;137:e67-e71.
- [39] Sawant M, Schwarz N, Windoffer R, Magin TM, Krieger J, Mucke N, et al. Threonine 150 phosphorylation of keratin 5 is linked to EBS and regulates filament assembly and cell viability. J Invest Dermatol. 2017.
- [40] Agullo-Pascual E, Cerrone M, Delmar M. Arrhythmogenic cardiomyopathy and Brugada syndrome: diseases of the connexome. FEBS Lett. 2014;588:1322-30.
- [41] Hoorntje ET, Te Rijdt WP, James CA, Pilichou K, Basso C, Judge DP, et al. Arrhythmogenic cardiomyopathy: pathology, genetics, and concepts in pathogenesis. Cardiovasc Res. 2017;113:1521-31.
- [42] Kant S, Holthofer B, Magin TM, Krusche CA, Leube RE. Desmoglein 2-Dependent Arrhythmogenic Cardiomyopathy Is Caused by a Loss of Adhesive Function. Circ Cardiovasc Genet. 2015;8:553-63.
- [43] Basso C, Czarnowska E, Della Barbera M, Bauce B, Beffagna G, Wlodarska EK, et al. Ultrastructural evidence of intercalated disc remodelling in arrhythmogenic right ventricular cardiomyopathy: an electron microscopy investigation on endomyocardial biopsies. Eur Heart J. 2006;27:1847-54.
- [44] Schlipp A, Schinner C, Spindler V, Vielmuth F, Gehmlich K, Syrris P, et al. Desmoglein-2 interaction is crucial for cardiomyocyte cohesion and function. Cardiovasc Res. 2014;104:245-57.
- [45] Kant S, Krull P, Eisner S, Leube RE, Krusche CA. Histological and ultrastructural abnormalities in murine desmoglein 2-mutant hearts. Cell Tissue Res. 2012;348:249-59.
- [46] Gercek M, Gercek M, Kant S, Simsekyilmaz S, Kassner A, Milting H, et al. Cardiomyocyte Hypertrophy in Arrhythmogenic Cardiomyopathy. Am J Pathol. 2017;187:752-66.

D 9 Microelectrode devices for *in vitro* and *in vivo* recordings of cellular activity

A. Offenhäusser, S. Weidlich, D. Kireev,

K. Srikantharajah, V. Rincón Montes

Bioelectronics

Institute of Complex Systems

Forschungszentrum Jülich GmbH

Contents

1	Intr	Introduction	
2	Interfacing electronic devices with neuronal cells		2
	2.1	Working principles	3
	2.2	Planar microelectrode arrays	4
3	Adding dimensionality – from planar to 2D ⁺ devices		5
	3.1	3D Nanoelectrodes	6
	3.2	Nanocavities	8
4	New materials		10
	4.1	Graphene-based devices	10
	4.2	Printed electronics	15
5	Towards 3D <i>in vivo</i> applications		17
	5.1	Flexible and penetrating MEAs	18
	5.2	Insertion tools	19
	5.3	Applications	20
Ref	erence	25	21

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

The human brain is an organ of vast complexity and despite ever increasing effort, scientists are still far from understanding how this network of several billion neurons accomplishes functional interaction and signaling or to find reliable methods to assist this intricate machinery in places where correct signaling fails. Due to its complexity, in vivo studies are limited to macro-scale investigations such as electroencephalography (EEG), magnetic resonance imaging (MRI), and positron emission tomography (PET). While these methods allow for the ascertainment of organ function and localization of its failures, it is impossible to characterize single cell function and the mechanisms of network communication from these studies. Therefore, our group focuses on the investigation of small-scale, two-dimensional neuronal networks (up to 1000 cells/mm²). their interaction, and possible means to influence their behavior. It is our goal to fabricate chipbased sensors that enable an efficient cell-chip coupling towards precise recording of cellular signals for *in vitro* studies and use this knowledge for *in vivo* applications. We want to facilitate a better understanding of the cell-cell communication and provide the ability to stimulate network communication by the transmission of electrical signals to the constituting neurons. Within this framework, we have developed a variety of microelectrode array (MEA) designs that enable non-invasive, parallel, multi-site recording of action potentials from primary neurons and the cardiomyocyte-like cell line HL-1. We have modified our standard planar 64 electrode MEA design with different geometries ranging from nanometer-sized cavities that allow for cellular protrusion into the sensor to mushroom-shaped 3D electrodes. Furthermore, we investigate various field-effect transistor (FET) designs with gate materials ranging from silicon nanowires to graphene. In order to realize the development of our devices, we employ both traditional cleanroom technology as well as modern inkjet printing procedures. Ultimately, the combination of our different approaches yields novel devices for the study of network development and communication for in vitro as well as in vivo applications.

This Chapter is meant to give an overview over the commonly used approaches in the area of microelectrode arrays for the study of cellular communication.

2 Interfacing electronic devices with neuronal cells

Silicon-based microstructures are continuously gaining importance in fundamental neuroscience and biomedical research. Precise and long-lasting neuro-electronic hybrid systems are at the center of research and development in this field. For extracellular signal recordings from electrically active cells in culture, two main concepts have been developed in the past: microelectrode arrays (MEAs) (see Fig. 1a) with metalized contacts on silicon or glass substrates have been used to monitor cardiac impulse propagation from dissociated embryonic myocytes [1]–[3], dissociated invertebrate neurons [4], [5], mammalian neurons [6], and spinal cord neurons [7]. Alternatively, arrays of field-effect transistors (FETs) (see Fig. 1b) are used for extracellular recordings having either non-metalized transistor gates with cells growing directly on the gate oxide [8]–[10] or metalized gates. The latter are in direct contact with the electrolyte [11] or they are electrical activity of single cells and networks of neurons can be observed over an extended period of time. Meanwhile, both concepts are growing together by designing MEAs inside a CMOS process with on-chip amplification and filtering [15], [16].



Fig. 1: (a) Substrate-embedded microelectrode: the metal electrode (red) is exposed to the electrolyte while the feedlines are covered with an isolation layer (green-blue) (b) Open-gate field-effect transistor for the recording of extracellular signals.

2.1 Working principles

The interaction of a neuronal cell with an electronic device is schematically depicted in Fig. 2. Sufficient electrical coupling between the cell and the (gate) electrode for extracellular signal recording is achieved only when a cell or a part of a cell is located directly on top of the (gate) electrode. Electrical signals recorded by these devices show lower signals and a higher noise level (owing to a weaker coupling to the (gate) electrode) compared to intracellular electrodes or patch pipettes (see Fig. 2, right).



Fig. 2: Left: Schematic of a neuron on an electronic device: intracellular (upper orange electrode) and extracellular (lower grey electrode) signals can be recorded. Right: Action potential of a neuron (approx. 100 mV) recorded by an intracellular electrode (upper trace) and an extracellular electrode (lower trace).

For a quantitative understanding of the extracellular signals recorded by electronic devices, it is necessary to consider the experimental environment in detail. A schematic picture of a typical experimental situation is depicted in Fig. 3. Here, the neuroelectronic hybrid is formed by the neuron, the cleft between neuron and the sensor surface, and the electronic device. Outside the neuron and inside the cleft, there is extracellular electrolyte solution. By electrical excitation, the ion channels in the cell's membrane open and ions can flow across the cell membrane. While in the upper part of the cell (free membrane) these ions just enter the surrounding electrolyte bath directly, the situation is different at the attached membrane. Here, the ions have to pass the cleft before entering/leaving the bath. The cleft acts as a resistance called seal resistance R_{seal} [9], [17]. The magnitude of R_{seal} is usually in the order of several 100 k Ω up to M Ω ,

corresponding to a typical cleft thickness of 40 to 150 nm [18]–[20]. The voltage V_J , which determines the potential at the (gate) electrode, is mainly determined by the seal resistance R_{seal} and the current that flows across it.

Fig. 3: Schematic representation of the neuroelectronic hybrid (Point contact model). The cell membrane is divided into free (FM) and junctional membrane (JM) with the respective values of membrane area (A_{FM}, A_{JM}) , membrane capacitance (C_{FM} , C_{IM}), and resistance (R_{FM} , R_{IM}). C_G and R_G are the capacitance and the resistance of the (gate) electrode, respectively. The seal resistor R_{seal} represents the electrical properties of the cleft between the membrane and the sensor surface. In case of patchclamp experiments, the intracellular voltage V_M can be measured directly.



2.2 Planar microelectrode arrays

Planar MEAs were first reported by Thomas *et al.* in 1972 [1]. The aim of their research was to establish an electrophysiology technique that enables a parallel, multi-site, non-invasive recording from electrically active tissues and cell cultures over periods of days to weeks. This goal was not achievable with the contemporary electrophysiological methods, such as sharp- or patch electrodes, due to their invasive nature and resulting loss of cell viability. Thomas's MEA design consisted of 30 electrodes arranged in two lines 50 μ m apart, a distance of 100 μ m between the electrodes within each line, and an area of 50 μ m² per electrode (Fig. 4).



Fig. 4: First MEA as developed by Thomas et al. in 1972. Their design consisted of 30 electrodes of an area of $50 \mu m^2$ per electrode, arranged in two lines $50 \mu m$ apart. After deposition of platinum black, these devices were employed to measure action potentials from embryonic chick cardiomyocytes. Image from [1].

They deposited platinum black on each electrode and succeeded in recording action potentials of up to 2.5 mV in amplitude from confluent, contracting layers of embryonic chick cardiomyocytes. Since then, various designs and materials have been employed for the fabrication of novel devices for improved signal quality and stability. While our group investigates a variety of different design, our most common layout exhibits 64 electrodes arranged in an 8x8 grid as shown in Fig. 5.



Fig. 5: Schematic representation of our 64 electrode MEA layout. In the center of the MEA, the electrodes are arranged in an 8x8 grid. An electrode opening in the range of 6 to $24 \mu m$ enables the interaction with the biological system while the remaining feedline is buried underneath a passivation layer. Bondpads in the periphery enable the connection to the measurement electronics. Image from [21].

While many improvements of planar MEAs have been developed over the years, ranging from increasing the number and density of electrodes, introducing high surface areas via platinum black, or chemical modification to increase cell adhesion and thus decrease the cleft between cell and electrode, planar MEAs still suffer from limited signal quality. Furthermore, it remains impossible to record subthreshold potentials using this method [22]. Due to these reasons, a variety of more complex designs based on the idea of extracellular, non-invasive MEAs have been developed. The following Sections will provide information on such possible improvements.

3 Adding dimensionality – from planar to **2D**⁺ devices

While microelectrode arrays offer a range of advantages such as their non-invasive nature and thus the ability to interact with cellular networks over extended times, multi-site measurement capabilities, and excellent temporal resolution, they also exhibit the drawback of recording strongly attenuated signals. A possible approach for the improvement of the recording capabilities of MEA-based devices is to increase the cell-electrode contact via the introduction of additional dimensionality, either through protruding structures that enable phagocytosis-like

events (3D Nanoelectrodes) or via cavities allowing for cellular protrusion into the sensor (Nanocavities). Fig. 6 depicts the different designs employed in our group. Despite their threedimensional nature, these nanoscale strategies are termed $2D^+$ in this Section to enable a distinction to the microscale penetrating 3D electrodes discussed in Section 5.

Different $2D^+$ strategies. **Fig. 6:** a)-c) Protruding 3D nanoelectrodes for improved cell-chip coupling via engulfment-like processes and resulting close contact between cell electrode. and d) Nanocavity electrode. Cellular protrusion into the sensor area results in good cellchip coupling. Figure adapted from [21].



3.1 3D Nanoelectrodes

Since their first report in the literature by Spira et al. in 2007 [23], 3D nanoelectrodes have been the focus of extensive research as possible solution for the problem of low sealing resistance. Spira observed the formation of a very tight cell-electrode contact between mushroom-shaped 3D structures and both Aplysia californica neurons and human cardiomyocytes [23]. These, as well as various other cell lines [24], readily engulf the mushroom-shaped gold spines of 850nm stalk width, 1 µm stalk height, 1.8 µm cap width, and 1.6 µm total height. For Aplysia neurons, the cell-electrode cleft was determined to be around $35\pm21\,\text{nm}$ around the spine in contrast to 56 ± 29 nm on flat surfaces [24]. First electrophysiological investigations employing this system were conducted in 2009, with a set of four mushroom-shaped electrodes yielding a 4.5 times increase in signal amplitude as compared to planar MEAs while maintaining the normal electrophysiological behavior of the cells [24]. Since then, a variety of different threedimensional nanoelectrode designs have been reported in the literature. In our group, pillars, mushroom-shaped structures, as well as hollow pillars and hollow mushrooms are prepared via electron-beam (e-beam) lithography in conjunction with electrodeposition, which enables an easy adjustment of the structure design and a bottom-up, parallel, controlled deposition process as depicted in Fig. 7. An e-beam sensitive polymer such as poly(methyl methacrylate) (PMMA) is spincoated onto a MEA (Fig. 7a, exemplary single electrode) and either circle or ring patterns are introduced into the layer via e-beam exposure to produce solid or hollow structures, respectively (Fig. 7 b_1, b_2). Gold is then electrodeposited into the holes, yielding either (hollow) pillars (Fig. 7 c_2 and c_1) or (hollow) mushroom-shaped electrodes (Fig. 7 d_2 and d_1), depending on whether the galvanization is allowed to proceed past the boundaries of the template. Subsequent removal of PMMA yields the free-standing structures (Fig. 7 e_1/e_2).



Fig. 7: Schematic representation of 3D nanoelectrode fabrication. PMMA is spincoated onto a MEA (a, exemplary single electrode) and patterned via e-beam exposure either with circles (b_1) or rings (b_2). Gold is then electrodeposited into the template, yielding either (hollow) pillars (c_2 and c_1) or (hollow) mushroom-shaped electrodes (d_2 and d_1), depending on whether the galvanization is allowed to proceed past the boundaries of the template. PMMA removal yields the free-standing electrodes (e_1/e_2). Figure adapted from [25].

Fig. 8 shows scanning electron microscopy images of the different 3D electrode designs obtained by application of the process in Fig. 7.



Fig. 8: The types of 3D structures fabricated according to the process in Fig. 7: solid pillars (a) and solid mushrooms (b), as well as hollow pillars (c) and hollow mushrooms (d). Images (a,b) were acquired at 65° sample tilt, (c,d) at 55° sample tilt. Scalebar represents 500 nm. Image adapted from [21].

In our group, we have conducted extensive research on the design criteria influencing the cellelectrode distance by means of focused ion-beam (FIB) sectioning and scanning electron microscopy (SEM). We found that the interaction of cells with 3D electrodes is greatly dependent on the geometry of the structure, with shape, size, and aspect ratio highly influencing the behavior of the cell [26]. Using cardiomyocyte-like HL-1 cells as model system, we could show that high aspect ratios between stalk height and width result in an improved cell-electrode contact and that the "cap" on mushroom-shaped structures results in a better engulfment by the cell as compared to plain pillars [26]. Investigations of embryonic rat cortical neurons on hollow pillars showed a good contact between the cell and electrode around the outer perimeter of the structure. However, in contrast to reports in the literature, no cellular protrusion into the hollow interior could be observed (Fig. 9).



Fig. 9: Focused ion-beam cross-section of an embryonic rat cortical neuron interacting with a hollow pillar of 900 nm in outer diameter. The observed tight contact between cell and electrode facilitates good electrical coupling. In contrast to reports in the literature, no cellular protrusion into the hollow opening can be observed. Scalebar in a) represents $2 \mu m$, scalebar in b) is $1 \mu m$. Image adapted from [21].

To date, many different three-dimensional electrode designs as well as associated electrophysiological studies have been reported in the literature. However, the conducted experiments differ in fundamental parameters such as the overall size of the electrode, the type and size of the 3D structure, the electrode material, the number of 3D structures per electrode, or the employed cell type. All of these aspects preclude a direct comparison of the results with respect to the improvement of the recording capabilities of the devices yielded by application of 3D structures. In order to solve this problem, we recently developed a process that enables the fabrication of multiple different structures with respect to both design and size on a single chip and thus parallel testing on the biological system [25]. In the future, we want to employ this system to enable a rapid-prototyping approach to 3D electrodes.

3.2 Nanocavities

Another approach investigated in our group is the usage of nanocavities [27], [28]. They can be considered to be an "inverted" version of 3D electrodes, enabling not an engulfment by the cell but rather a protrusion of the cell into the sensor. These structures are prepared using a stack of either platinum or gold as electrode material, covered by a sacrificial layer of chromium as shown in Fig. 10. Small apertures of up to $24 \,\mu\text{m}$ enable the interaction with the environment. Since chromium can easily be dissolved using commercial etching solutions which leave both platinum and gold layers intact, a selective underetching of the passivation by removal of chromium enables the formation of a nanocavity.



Fig. 10: Left: Schematic representation of the nanocavity concept. The electrode material is covered by a sacrificial layer of chromium. Small apertures of up to $24 \mu m$ in diameter enable the interaction with the environment. Removal of the chromium via an etching step results in the formation of a nanocavity underneath the passivation that facilitates cellular protrusion into the opening and good sealing. Scheme according to [28]. Right: Nanocavity chip after etching. The bright areas surrounding the circular apertures mark the underetched area.

While the vertical dimension of the cavity can be tailored via the thickness of the chromium layer, the horizontal dimension depends on the parameters of the etching process. Due to the increase of the active surface area, the impedance of the device decreases dramatically with increasing etching time and thus larger nanocavity size, as can be seen in Fig. 11.



Fig. 11: Dependence of the electrode impedance on the etching time and thus cavity size. The impedance decreases up to an etching time of approximately 600s, at which the ohmic resistance of the narrow cavity starts to mask the effect of the increased electrodeelectrolyte interface. Image adapted from [29].

A low device impedance is an important aspect for low-noise electrophysiological measurements [29]. In contrast to reducing the impedance via an increase in the overall electrode area, the introduction of nanocavities enables the fabrication of low-impedance devices while maintaining a small geometrical footprint of the interaction site and thus high spatial resolution. Recordings from HL-1 cells using these devices exhibit stable action potentials with an amplitude of around 1 mV and excellent signal-to-noise characteristics [28]. Furthermore, this approach also enables the stimulation of electrical activity [29]. Fig. 12 shows exemplary action potential recordings from embryonic rat cortical neurons employing both planar MEAs as well as nanocavity sensors. When comparing planar and nanocavity MEAs of the same aperture size of $12 \,\mu$ m, the nanocavity MEA yields significantly improved signal

amplitudes with $346 \,\mu V$ peak-to-peak (p2p) as compared to $176 \,\mu V_{p2p}$ on planar MEAs. Due to the larger active electrode area, nanocavity electrodes with an aperture of $3 \,\mu m$ yield a noise level comparable to that of planar $12 \,\mu m$ electrodes, while still enabling signal amplitudes larger than those recorded on the larger, planar electrodes.

Fig. 12: Action potential recordings of embryonic rat cortical neurons on nanocavity electrodes of different aperture size in comparison to a planar electrode. When comparing electrodes of equal aperture size, the nanocavity device yields significantly improved signal amplitudes at lower noise level. Even nanocavity devices with verv small electrode apertures of 3 µm still enable high amplitude signals with a noise level comparable to that of planar electrodes. 12 μm Image courtesy of Johannes Lewen.



In future studies, we aim to employ the nanocavity platform to investigate the development of neuronal networks. Furthermore, we want to combine 3D nanoelectrodes with geometries that enable cellular protrusion into the active area to merge the advantages of engulfment-inducing 3D structures and nanocavities.

4 New materials

4.1 Graphene-based devices

Graphene FETs

Graphene, a single-layer allotrope of carbon consisting of carbon atoms arranged in a hexagonal lattice, was first proposed and implemented as a transistor material by Geim and Novoselov in their Nobel Prize winning work in 2004, showing that graphene is uniquely responsive to an applied external gate potential [30] and can exhibit extremely large mobilities μ of up to $1 \times 10^5 \cdot 1 \times 10^6 \text{ cm}^2 \text{V}^{-1} \text{s}^{-1}$ and consequently high sensitivities. While, these values are obtained for exfoliated graphene, suspended and measured in a close-to-zero temperature [31], [32], the mobility in ambient conditions is usually in the range of 1000-50000 cm² V⁻¹ s⁻¹, depending on the quality of graphene [33]. Despite this being significantly lower than at close-to-zero temperature, the mobility at ambient conditions is still large, especially if compared to that of silicon nanowires, which exhibit mobilities of around $1000 \text{ cm}^2 \text{V}^{-1} \text{s}^{-1}$ [34]. A typical layout and the characteristics of a liquid-gated graphene field effect transistor (GFET) is shown in Fig. 13. The GFET's typical I-V curve shows a cone-like behavior of the current flowing through the graphene channel (left y-axis) upon application of an electrical field to the gate dielectric (x-axis).



Fig. 13: Top: Schematic representation of a liquid-gated GFET. Graphene is employed as channel material and changes in the channel conductance vield information about modulations of the liquid gate. Bottom: Representative transfer curve. The parabolic I-V curve is shown in blue in the region of hole conduction and in red in the region where electron conduction occurs. The first derivative of the I-V curve, representing the transconductance, is shown in green. These characteristics of the device are determined experimentally and then employed to deduce changes occurring at the liquid gate.

Initially, the current is commonly dominated by holes, goes through a minimum at the Dirac point, where the number of charge carriers is minimal, and then increases again, though the conductance is then dominated by electrons. The two key parameters, transconductance (g_m) and mobility (μ), can be computed according to the following equations

$$g_m = \frac{\Delta I_{DS}}{\Delta V_{GS}}$$
 and $\mu = \frac{L}{W} \cdot \frac{g}{c_{ox} \cdot V_{DS}}$

where I_{DS} , V_{DS} , and V_{GS} are the drain-source current, drain-source voltage, and gate-source voltage, respectively, C_{ax} is the gate oxide's capacitance, and W and L are the width and length of the graphene channel, respectively. The transconductance is the actual amplification factor and, consequently, the sensitivity factor and the figure of merit when the GFET is used for sensing applications. The mobility is correlated to the electrical double layer capacitance in the case of liquid gating. Sensing cellular activity occurs utilizing the device characteristics depicted in Fig. 13 to deduce changes occurring at the gate: if capacitive coupling between the cell and transistor is assumed, the changes in the graphene's drain-source current can be easily translated into changes of the extracellular voltage. Based on this principle, different groups have used graphene transistors and their arrays to record cellular action potentials, their propagation and, moreover, perform biochemical studies. In particular, electrical activities of the human kidney cell line HEK-293 and cardiomyocyte-like HL-1 cells have been successfully measured with GFETs [35]-[37]. Fig. 14a shows an optical microscopy image of GFETs as fabricated in our group. When cardiomyocyte-like HL-1 cells are cultured on the chips' surface, they form a continuous and conformally contracting layer. A typical recording of the HL-1 action potentials (AP) from one of the GFETs is shown in Fig. 14b. The cells are contractile, producing repetitive APs that propagate through the cellular layer with rate of 23 ± 3 beats per minute (bpm) and amplitude of 1.2±0.2 mV_{p2p} in the presented case. Typically, the beating rate of an HL-1 cell culture may vary from 15 to 150 bpm and depends on many environmental and chemical conditions. The shape of the APs is related to different coupling mechanisms. An example is shown in Fig. 14b, where more than a hundred consecutive spikes are averaged. The



Fig. 14: (a) Optical microscopy image of a graphene transistor array. (b) A typical time trace of HL-1 activity recorded with a graphene transistor as well as the averaged HL-1 spikes from 115 individual consecutive spikes from the chip, yielding an average amplitude of $1200 \mu V_{p2p}$ (c) Time trace of neuronal recording with the inherent neuronal feature of bursting, when the neurons exhibit alternating periods of high-frequency activity (bursts) and low-frequency intermittent spiking. The averaged AP from 77 individual APs from the neuronal time series yields an amplitude of $630 \mu V_{p2p}$.

shape of the action potential, as in agreement with previous works, represents a very good sealing between the cell and the transistor [38], [39]. Neuronal recordings with GFETs have also been reported to be possible both *in vitro* [37], [40] as well as *in vivo* [41], [42]. When cultured *in vitro*, cortical neurons must be cultured for at least two weeks before they start producing spontaneous and large amplitude action potentials that can propagate through the network. Compared to HL-1 cells, neurons yield smaller extracellular action potentials [6]. One representative neuronal action potential recording is presented in Fig. 14c, with a well-defined bursting pattern. Averaged and reconstructed, the APs exhibit a comparably high amplitude of $600 \,\mu$ V, but an unusual shape and low signal-to-noise ratio (SNR) of the recordings due to large effective gate noise of up to 100-200 μ V. Such a noise level consequently impedes the possibility of recording ultra-low amplitude neuronal signals, which are typically below 500 μ V and would thus barely exceed the noise.

Graphene MEAs

Graphene microelectrode arrays (GMEAs) can be produced with a significantly less complex fabrication process as compared to GFETs, with three to four fabrication steps allowing a wafer-scale process when using specific graphene transfer [43]. Typically, rigid borofloat glass or SiO_2/Si substrates are employed for the device fabrication. Borofloat wafers are used because of their transparency, which aids the long-term monitoring of the cell cultures, particularly in combination with the inherent transparency of graphene. For the fabrication of GMEAs, CVD grown graphene is used due to the need for large crystal sizes. Similarly to GFETs, polyimide or SU-8 photoresists are commonly utilized as passivation polymers. In agreement to previously developed MEAs and in order to utilize the same multichannel measurement schematics, the

GMEA chips are typically $24 \times 24 \text{ mm}^2$ in size and exhibit an array of 64 electrodes. They follow the same layout as depicted in Fig. 5, with the active electrode area made up of graphene rather than a metal layer as shown in Fig. 15.

Topview



Fig. 15: Schematic layout of the fabricated GMEA chips. While they follow the general structure as depicted in Fig. 5, with metal feedlines connecting the electrodes to the bondpads, the active electrode area is prepared entirely from graphene rather than metal. The chips exhibit a size of $24 \text{ mm } x \ 24 \text{ mm}$, with the 64 electrode openings distributed over an area of $1.4 \text{ mm } x \ 1.4 \text{ mm}$ in size.

As with the GFETs, the initial proof-of-concept studies were performed with HL-1 cells. Livedead staining was employed to prove the biocompatibility of the devices. The cells were electrically active and trains of action potentials could be recorded as seen in Fig. 16. There is a visible time-delay between the APs detected at different recording sites, which allows to evaluate the signal propagation within the cell layer. The recorded beating frequency in this case is in the range of 1 ± 0.5 Hz. The recorded action potential amplitudes and their shapes vary from chip to chip (*culture effect*) and from electrode to electrode (*sealing effect*). A large number of HL-1 action potentials from various time series were analyzed in order to compare them with the experimental and simulated data [44]. The resulting pie-chart depicting the different AP shapes observed during this study is shown in Fig. 16.



Fig. 16: Left: Time trace of action potentials recorded from HL-1 cells on six different electrodes on one GMEA chip, showing the repetitiveness of the spikes and clear propagation of the signal. Right: distribution of different signal shapes. The large occurrence of spikes of type A shows that good cell-electrode sealing can be achieved with the developed GMEA devices. Spikes of types B through D result from lower sealing resistances. Figure adapted from [44].

As already discussed, many factors contribute to the shape of the APs recorded with extracellular electrodes. Regardless of the impedance of the electrode itself, there are other physical and physiological parameters that affect the way the signal will be recorded and seen. First of all, the more mature the culture, the larger and more stable the APs [45]. The second important parameter is the *sealing* between the cellular layer and the electrode [39]. For the measurements conducted as part of this study, the action potential shapes A and B were observed with the highest occurrence (see Fig. 16). According to previously simulated data [39], spikes of type A result from a large sealing resistance, large sodium peak, and large amplitude. The difference in pre- and post-spike overshoot, their amplitude, and duration can be modeled by variations in the sealing resistance [22] and current flows of Na⁺, Ca²⁺, and K⁺ ions [9], [38], [39]. Spikes of types B, C, and D are found to be not significantly different from each other, only varying in the absence or presence of the post- and pre-spike overshoot, which can be described by differences in sealing. The last, type E, occurring the rarest, has an extremely slow negative component, which, according to [38], can be dominated by the Ca^{2+} component of the action potential. The remaining spikes, categorized as "other", mostly consist of very unreproducible shapes with double or triple peaks, and usually are the result of a pinhole in the passivation or other defects.

Next, embryonic rat cortical neurons were cultured on the GMEAs until a mature network was established (Fig. 17a). Neuronal cultures typically begin to exhibit spontaneous electrical activity that is propagating through the network after 14 to 21 days *in vitro*. The variety of spontaneous spiking-bursting activities recorded with GMEA chips is very similar to patterns recorded with planar Au-MEAs [7], [46]–[50]. Fig. 17b shows one time trace with distinguishable bursting and non-bursting (random) spikes. Depending on the culture, these bursts occur every 5 to 15 seconds, each burst containing a series of very high amplitude spikes of up to $800 \,\mu V_{p2p}$, followed by a series of evanescent spikes. Between the bursts, nonbursting, random APs with typically lower amplitude (50-150 μV_{p2p}) can be recorded.



Fig. 17: (a) Microscopic image of a neuronal culture on a GMEA chip. (b) Representative time trace showing the characteristic bursting behavior of neuronal action potentials. These bursts occur every 5 to 15s and contain a series of high-amplitude spikes of up to $800 \mu V_{p2p}$. Figure adapted from [44].

In addition to the fabrication on rigid substrates, GMEAs were produced on flexible substrates such as parylene or polyimide. The devices were shown to be both highly flexible as well as stable even after mechanical deformation [51]. High-amplitude action potentials of up to $1.0\pm0.2 \,\mathrm{mV_{p2p}}$ could be recorded from embryonic heart tissue. Recordings from HL-1 cells yielded action potentials of up to 300 μV_{p2p} in amplitude with good signal to noise

characteristics. In the future, these devices could be employed for applications that necessitate conformal tissue contact, as for example needed for *in vivo* application.

4.2 Printed electronics

One major disadvantage of conventional cleanroom technology is the costly and timeconsuming nature of the device fabrication. Furthermore, due to the mask-based production, an adaptive fabrication towards device prototyping is highly expensive, necessitating the production of new mask designs via electron-beam lithography for each device generation. As an alternative, our group investigates inkiet printing for the fabrication of bioelectronic sensors. Inkjet printing offers the advantage of a high-throughput, low cost, and maskless manufacturing process. This approach enables rapid prototyping capabilities merely by changing the digital design file. One major limitation of printed electronics, however, is the lower resolution as compared to cleanroom fabricated devices. During inkiet printing, the lateral resolution is limited by an interplay between the achievable droplet size and its interaction with the surface to around 10 to $20 \,\mu\text{m}$ [52], while photolithography can routinely achieve feature sizes of less than 1 µm. For this reason, different groups have employed a mixed approach with only certain steps of the device fabrication being performed via inkjet printing while classic cleanroom technology is employed for the remaining steps. However, a careful tailoring of the ink characteristics, surface energy of the substrate, and resulting wetting characteristics of the sample surface can enable the fabrication of fully ink-jet printed devices with adequate resolution for bioelectronic applications. In this manner, our group was able to develop an inkjet-printed 64 electrode MEA with electrode openings of 31 µm in diameter and a spacing of 200 µm between the electrodes [53]. Another advantage of inkiet printing is the flexibility not only with respect to design but also substrate material and ink composition. This enables, for example, the fabrication of highly porous electrodes and thus electrodes of high surface area, an aspect that greatly influences the recording capabilities of the devices since an increase in surface area results in a decrease in impedance and thus thermal noise. While increasing the overall electrode size also results in a reduction in impedance, this decrease in impedance occurs at the expense of the spatial resolution. In contrast, rough surfaces enable a significant reduction in impedance while maintaining a high spatial resolution. However, rough or even porous electrodes are not directly accessible via conventional cleanroom fabrication methods. In our group, Schnitker et al. [54] developed a porous carbon-based microelectrode array that exhibits a specific interfacial capacitance of around $880 \,\mu\text{F/cm}^2$, which is approximately a 30-fold increase in capacitance as compared to non-porous carbon surfaces. Fig. 18 depicts the fabrication process for these porous carbon MEAs. In the first step, the outer feedlines are printed onto a polyethylene naphthalate (PEN) substrate with a silver nanoparticle ink. Afterwards, the active electrode area is printed with a carbon nanoparticle ink, thereby preventing any contact between the silver layer and the biologic system. Finally, a polyimide (PI)-based passivation ink is printed to cover all but a small area in the center of the chip. Sintering of the final device results both in the establishment of interconnected feedlines due to a melting of the nanoparticles, as well as removal of solvents from all printed layers.





Fig. 18: Schematic fabrication flow for inkjet-printed MEAs. 1) Printing of the outer feedlines with a silver nanoparticle ink. 2) Printing of the active electrode area with a carbon nanoparticle ink. 3) Passivation of the device with a polyimide (PI) based dielectric layer. Image adapted from [54].

In this manner, electrodes of around $30\,\mu\text{m}$ in diameter and 2 to $3\,\mu\text{m}$ spacing were produced (Fig. 19b). FIB sectioning in conjunction with scanning electron microscopy could prove the highly porous nature of the carbon layer (Fig. 19c).



Fig. 19: a) Flexible, fully inkjet-printed printed carbon MEAs. b) Microscopic image depicting the carbon microelectrodes and PI passivation. c) FIB-SEM cross section of the highly porous carbon electrode. Images courtesy of Jan Schnitker, in parts adapted from [54].

These all-inkjet-printed MEAs could be fabricated within the timeframe of less than one hour at a cost of less than 3 cents per chip, enabling both a rapid prototyping as well as low-cost production. Furthermore, the developed MEAs were shown to be highly biocompatible, yielding a viability of 95% for cardiomyocyte-like HL-1 cells cultured on the devices. Action potential recordings of HL-1 cells yielded signals with amplitudes of up to 0.96 mV_{p2p}, which is comparable to the results obtained with cleanroom fabricated devices.

Another important aspect for the fabrication of MEAs for the application as implants are the mechanical characteristics of the devices. While cleanroom fabricated, silicon technologybased MEAs exhibit a Young's modulus in the GPa range, the mechanical properties of the central nervous system are in the range of 100 Pa to 10 kPa. This mechanical mismatch results in an inflammatory response of the surrounding tissue, glial scar formation, and loss of functional interaction. Therefore, the fabrication of flexible electronics is an important aspect for next-generation neuroprosthetic devices. While cleanroom-fabricated MEAs can be produced on flexible materials, the employed polymers are still fairly rigid. Truly soft polymers such as hydrogels, however, are often not compatible with cleanroom fabrication processes, making the production of highly flexible devices difficult with the available methods. While inkjet printing on highly flexible materials is not simple and necessitates a careful tailoring of the surface characteristics of the substrate as well as the sintering conditions, we could prepare all-inkjet-printed carbon MEAs on substrates such as PDMS, gelatin, agarose, and even commercial gummy bears [55]. We could show functional interaction of the produced devices with HL-1 cells, yielding action potentials of up to $906 \,\mu V_{p2p}$ on PDMS and up to $442 \,\mu V_{p2p}$ on gummy bear MEAs. Fig. 20 shows exemplary HL-1 action potentials recorded with an allinkiet-printed MEA on a gummy bear substrate. In the future, this approach could pave the way for the production of highly conformal, biocompatible, low-cost bioelectronic implants.





5 Towards 3D in vivo applications

In vivo neural probes have been gaining increasing importance for the diagnosis and treatment of neurological diseases due to the possibility of simultaneously targeting different locations within the living tissue. For decades, silicon-based penetrating devices have been primarily used for chronic neuronal recordings. However, these probes suffer from inconsistent long-term performance due to probe corrosion as well as tissue damage and resulting immune responses, which is caused by the mechanical and biological mismatch between the tissue and the probes. Different approaches have been developed in the last years in order to design and fabricate more compliant devices, considering flexible and soft materials such as polymers as substrates for the devices. In the following Section, suitable materials and several design considerations will

be discussed. Here, the main focus lies on penetrating probes, therefore different strategies to insert such flexible systems inside the tissue will be discussed. At the end, two possible applications for flexible and penetrating probes will be explored.

5.1 Flexible and penetrating MEAs

Polymers are the material of choice for ensuring the fabrication of devices which are more compliant and thus conformal with the brain tissue. Several different polymers are used to fabricate flexible and soft probes. Here, we will focus on the following three polymers: polyimide, parylene, and polydimethylsiloxane (PDMS).

Polyimide is widely used as insulation layer due to its high stability towards thermal, chemical, and physical effects. In recent years, it has gained increasing popularity for biological applications. Two major advantages of polyimide are the ease of preparation of layers down to few micrometers in thickness as well as the compatibility with photolithographic processes [56], [57]. Parylene, and in particular parylene-C, is an inert and optically transparent polymer that is commonly used as coating material as it can be deposited easily using low temperature chemical vapour deposition processes. Additionally, parylene is employed for the fabrication of biomedical implants such as stents and pacemakers. The biocompatibility of this polymer has been confirmed by receiving USP Class VI and ISO 10993 compliance [56], [57]. PDMS is a material commonly used for microfluidic applications and non-penetrating electrodes. Advantages of PDMS are its high flexibility, permeability to gases, and low Young's modulus $(\sim 1.8 \text{ MPa})$, which is significantly lower than that of the above-mentioned polymers (parylene ~2.8 GPa, polyimide ~2.5 GPa) [57]. Additionally, the Young's modulus of PDMS can be tuned by varying the pre-polymer and curing agent ratio and optimizing the curing temperature. While not all PDMS formulations meet the requirements of USP Class VI and ISO 10993, special formulations meeting this standard are commercially available. Besides the material selection, another important fact which should be considered during the fabrication process is to adapt the design and the geometry of the probes to the desired application. The general design of flexible, penetrating probes consists of long, thin, sword-like extensions called shanks, with multiple electrodes distributed along the shank and a sharp tip enabling the penetration of the tissue (see Fig. 21).



Fig. 21: General design of a flexible probe consisting of multiple shanks, with several electrodes distributed along the shank, and a sharp tip to facilitate the insertion into the tissue.

The length of the shank should be customized to the desired application. For example, cortical implants with approximately 1 mm shank length enable recording and stimulation of all cellular layers within the neocortex of the mouse. In contrast, a shank length of $200\,\mu\text{m}$ is sufficient for retinal probes, facilitating the interaction with the different layers within the retina, which is

only $400\,\mu$ m thick in humans and around $200\,\mu$ m in rodents. Besides the shank length, the design of the probe has to facilitate the penetration of the tissue. An important parameter influencing this aspect is not only the cross-section of the probe, which influences the force needed to penetrate the tissue, but also the opening angle of the shank tip. It is known that larger opening angles result in larger penetration forces [58]. Additionally, several modifications can be introduced to improve the performance of the *in vivo* probe. For example, anchors have been utilized at the edges of the shanks to minimize micromotion of the probe after implantation [59]. Finally, the probe should also meet the requirements for packaging as the electrical activity of the tissue has to be transmitted to a computer to enable the measurement and analysis of the electrical signals.

5.2 Insertion tools

As the development of compliant neural probes relies on soft, flexible, and penetrating features, a proper strategy to insert such devices into the tissue must be taken into account. In most of the cases, the probes themselves are incapable of penetrating the tissue without any aid, thus requiring a stiff shuttle or a biodegradable coating system that makes the flexible implant temporarily rigid to ensure its penetration into the target tissue [56]. Stiff shuttles usually refer to rigid and structural guides, which could simply be a metal rod/wire or a rigid polymer vehicle bound to the back of the flexible probe (see Fig. 22a). Both, shuttle and probe, can be transitorily fixed together either by surface modifications to enhance physical interactions among the probe and shuttle, such as a self-assembled monolayer [60], or by the use of biodissolvable adhesives, such as polyethylene glycol (PEG) [61], [62]. While it is true that this method ensures enough stiffness for insertion and offers the advantage of keeping the electrode sites intact, the tissue is prone to unnecessary trauma and the implant is susceptible to micromotion and displacement during the stiffener's removal [56].



Fig. 22: Schematic illustration of different insertion tools. a) One approach for a possible insertion system is a stiff shuttle bound to the flexible probes using PEG [61]. Another attempt is the application of biodegradable coatings. Such coatings can be done using b) dipcoating [63], or micro-molding with, for example c) syringe-jetting [64].

In contrast, biodegradable insertion systems embody all approaches that use a biocompatible and degradable polymer coating that stiffens the probe under dry conditions and dissolves when in contact with moisture and aqueous environments. Coating materials such as PEG, silk, polyvinyl alcohol (PVA), and sugars, among others, are currently being used and tuned to achieve specific degradation and physiological absorption rates [56], therefore necessitating a trade-off between polymer molecular weight, polymer-to-water ratio, and thickness. In order to coat the flexible probes, various coating methods such as dip-coating [63], [65] (see Fig. 22b) or micro-molding techniques [64], [66] are being explored in different research groups. When using dip-coating, the active electrode sites are initially covered with the insertion coating, temporarily affecting the recording capabilities of the device until the polymer is completely dissolved. When using micro-molding, further syringe-jetting [64] (see Fig. 22c), spin-coating [66], or even blade-coating techniques are required, thereby allowing the possibility of covering only the back-side of the probe.

The aforementioned coating methods for biodegradable polymers have been tested successfully for *in vivo* applications [56], however, the whole coating process is highly inefficient as each probe is coated manually. Hence, a wafer scale process compatible with microfabrication techniques should be considered. The introduction of a built-in but removable mold capable of casting the biodissolvable polymer shuttle at the back-side of the flexible device on the wafer, ideally as the last step of an up-side down probe fabrication (see Fig. 23), or the use of photopatternable biodegradable polymers [67] are methodologies and technologies that are currently in development.



Fig. 23: Casting a biodissolvable shuttle in a microfabrication process with a built-in and removable mold on the wafer scale. a) top view, b) cross-section, and c) stack of flexible probe with insertion shuttle after release from wafer.

Furthermore, other approaches combine flexible encapsulation layers like parylene with substrate polymers that adapt mechanically to their environment, fabricating in this way implants that are stiff or soft under dry or wet conditions, respectively [68]. In light of the high solubility and degradability featured by the aforementioned bio-polymers, the main challenge for their integration in efficient microfabrication processes remains in the use of dry patterning techniques such as laser micromachining and micro-molding, or dry ways to perform lift-off processes, as well as in the choice of the proper polymers and resists to yield a suitable and cost-effective fabrication process of flexible probes and their corresponding insertion tool.

5.3 Applications

Flexible and penetrating probes can be adapted in their design to match several neuroprosthetic applications. With longer shank lengths, cortical implants can be fabricated and used to record and electrically stimulate the different cellular layers within the neocortex. Having multiple electrode sites at a certain distance and in multiple shanks enables the monitoring of the neuronal interaction in a certain three-dimensional volume within a specific brain area. In the future, such flexible and penetrating probes could enable long-term recordings that can be used to analyze the changes in electrical activity caused by diseases such as epilepsy and Parkinson's disease. Moreover, these flexible devices can be used to improve approaches for retinal implants, aimed at treating people with degenerative retinal diseases. The introduction of

penetrating shanks with multiple electrode sites allows the interaction with the different retinal layers, thereby opening the possibility to record from retinal ganglion cells and electrically stimulate the inner retinal cells simultaneously. Such a system could not only provide information about the success of the electrical stimulation, but also detect abnormal retinal activity due to remodeling processes of the tissue.

References

- [1] C. A. Thomas et al., Exp. Cell Res. 74, 1 (1972)
- [2] D. A. Israel et al., Am. J. Physiol. 247, 4 (1984)
- [3] P. Connolly et al., Biosens. Bioelectron. 5, 3 (1990)
- [4] L. J. Breckenridge et al., J. Neurosci. Res. 42, 2 (1995)
- [5] W. G. Regehr et al., IEEE Trans. Biomed. Eng. 35, 12 (1988)
- [6] J. Pine, J. Neurosci. Methods 2, 1 (1980)
- [7] G. W. Gross et al., J. Neurosci. Methods 5, 1–2 (1982)
- [8] P. Bergveld et al., IEEE Trans. Biomed. Eng. 23, 2 (1976)
- [9] P. Fromherz et al., Science 252, 5010 (1991)
- [10] A. Offenhäusser et al., Biosens. Bioelectron. 12, 8 (1997)
- [11] D. T. Jobling et al., Med Biol Eng Comput 19, 5 (1981)
- [12] A. Offenhäusser et al., J. Vac. Sci. Technol. a 13, 5 (1995)
- [13] A. Cohen et al., Biosens. Bioelectron. 19, 12 (2004)
- [14] S. Meyburg et al., Biosens. Bioelectron. 21, 7 (2006)
- [15] F. Heer et al., IEEE J. Solid-State Circuits 41, 7 (2006)
- [16] K. Imfeld et al., IEEE Trans. Biomed. Eng. 55, 8 (2008)
- [17] W. L. C. Rutten, Annu. Rev. Biomed. Eng. 4 (2002)
- [18] A. Lambacher and P. Fromherz, Appl. Phys. A 63 (1996)
- [19] G. Wrobel et al., J. R. Soc. Interface 5, 19 (2008)
- [20] K. Toma et al., ACS Nano 8, 12 (2014)
- [21] S. D. Weidlich, Schriften des Forschungszentrums Jülich, (2017)
- [22] M. E. Spira and A. Hai, Nat. Nanotechnol. 8, 2 (2013)
- [23] M. E. Spira et al., TRANSDUCERS 2007 2007 Int. Solid-State Sensors, Actuators Microsystems Conf. (2007)
- [24] A. Hai et al., J. R. Soc. 6, 41 (2009)
- [25] S. Weidlich et al., Nanotechnology 28, 9 (2017)
- [26] F. Santoro et al., ACS Nano 8, 7 (2014)
- [27] B. Hofmann et al., Lab Chip 11, 6 (2011)
- [28] A. Czeschik et al., Phys. Status Solidi 211, 6 (2014)
- [29] A. Czeschik et al., Nanoscale 7, 20 (2015)
- [30] K. S. S. Novoselov et al., Science 306, 5696 (2004)
- [31] K. I. Bolotin et al., Solid State Commun. 146, 9–10 (2008)
- [32] S. V. Morozov et al., Phys. Rev. Lett. 100, 1 (2008)
- [33] F. Schwierz, Nat. Nanotechnol. 5, 7 (2010)
- [34] E. B. Ramayya et al., IEEE Trans. Nanotechnol. 6, 1 (2007)
- [35] L. H. Hess et al., Small 11, 14 (2015)
- [36] D. Kireev et al., IEEE Trans. Nanotechnol. 17, 1 (2016)
- [37] D. Kireev et al., Sci. Rep. 7, 1 (2017)
- [38] C. Sprössler et al., Phys. Rev. E 60, 2 Pt B (1999)

- [39] M. Schottdorf et al., Phys. Rev. E 85, 3 (2012)
- [40] F. Veliev et al., Front. Neurosci. 11, August (2017)
- [41] B. M. Blaschke et al., 2D Mater. 4, 2 (2017)
- [42] C. Hébert et al., Adv. Funct. Mater. 1703976 (2017)
- [43] D. Kireev et al., Carbon N. Y. 107 (2016)
- [44] D. Kireev et al., Adv. Healthc. Mater. 6, 12 (2017)
- [45] C. Xie et al., Nat. Nanotechnol. 7, 3 (2012)
- [46] G. W. Gross, IEEE Trans. Biomed. Eng. BME-26, 5 (1979)
- [47] L. Berdondini et al., J. Neurosci. Methods 177, 2 (2009)
- [48] M. Chiappalone et al., Brain Res. 1093, 1 (2006)
- [49] V. Pasquale et al., Neuroscience 153, 4 (2008)
- [50] J. vanPelt et al., IEEE Trans. Biomed. Eng. 51, 11 (2004)
- [51] D. Kireev et al., Biosensors 7, 1 (2016)
- [52] N. Y. Adly et al., RSC Adv. 7, 9 (2017)
- [53] B. Bachmann et al., Flex. Print. Electron. 2, 3 (2017)
- [54] J. Schnitker et al., Adv. Biosyst. Accepted (2018)
- [55] N. Adly et al., Submitted (2018)
- [56] A. Weltman et al., Micromachines. 2016
- [57] J. H. Lee et al., Lab Chip (2016)
- [58] D. J. Edell et al., IEEE Trans. Biomed. Eng. 39, 6 (1992)
- [59] P. Köhler et al., Proc. 31st Annu. Int. Conf. IEEE EMBC 2009 978 (2009)
- [60] T. D. Y. Kozai and D. R. Kipke, J. Neurosci. Methods (2009)
- [61] S. Felix et al., Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc. EMBS (2012)
- [62] B. J. Kim et al., J. Neural Eng. (2013)
- [63] Z. Xiang et al., J. Micromechanics Microengineering (2014)
- [64] A. Lecomte et al., J. Micromechanics Microengineering (2015)
- [65] K. L. Tan et al., IEEE 17th Electron. Packag. Technol. Conf. (2015)
- [66] F. Barz et al., Proc. IEEE Int. Conf. Micro Electro Mech. Syst. February (2015)
- [67] W. Liu et al., Adv. Sci. (2017)
- [68] A. E. Hess et al., J. Micromechanics Microengineering (2011)
E 1 Rheology

M. P. Lettinga Soft Matter Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Introduction	2	
2	Shear flow and extensional flow		
3	Scaling time		
4	Measuring geometries		
5	Examples of linear Rheological characterization		
6	Examples of non-linear Rheological responses6.1Strain hardening and softening6.2Shear thinning and thickening	10 10 12	
7	In situ rheology		
8	Micro-rheology and fluidics	17	
9	Concluding remarks	20	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Everything flows, or $\pi\alpha\nu\tau\alpha$ $\rho\epsilon\iota$ as the ancient greeks used to say (" $\pi\alpha\nu\tau\alpha$ " means "everything" and $\rho\epsilon\iota$ means "flows"). If you think of something that does not flow, then just wait. To put it in the words of the Prophetess Deborah" In the eyes of the lord the mountains flow". To put the words of Deborah into an equation: systems flow when the Deborah number $De = \tau_{system}/\tau_{probe} < 1$, so when the typical time during which the system is probed, $\tau_{probinq}$, is longer than the relaxation time of the system, τ_{system} . However, a force must be exerted on the system in order to induce flow. When the force is removed it might be that the material will return to its original state, in which case it is a solid. It might also be that it continues to flow for some time, in which case it is a fluid. For most materials, something in between will happen. Depending on the exerted force and the particles that constitute the material, a solid can also be turned into a fluid and vise versa. In other words, flow can change the characteristics of the material. Understanding of the flow behavior of the material means understanding its "rheology", the logic of its flow. Thus, rheology covers two main areas. The first area concerns the characterization of the material while the second area concerns the transitions in the material as induced by flow. The physical world consists of materials. Our body is an assembly of complex soft materials, as well as what we pull over our bodies, as what we put into our body, or the materials we use to transport our body. If we want to understand the mechanics of all these processes then we need at least some rheological understanding of these materials.

In this short introduction to rheology I want to familiarize the reader with the basic experiments that rheologists use to characterize and manipulate materials and the jargon to describe the mechanical properties of these materials. Each of the experiments will be elucidated by examples from biology, as this is the theme of this school. One often wonders how much information on a very molecular scale can be obtained, though this might depend on very involved modeling. To show that there is indeed a connection between structure and flow, also a few examples will be given where structural information is obtained *in situ*. Moreover, rheology on the micro-scale will be discussed as this has high biological relevance.

This crash course on rheology is a success when at the end the reader knows what is meant with shear thinning and thickening, with yielding, and with the storage and loss moduli. For those of you who are interested in the a clear and thorough introduction into rheology, I would refer to a few very nice books [1, 2]. This chapter is partly based on another chapter in [3], which focuses on the interplay between phase behavior and shear flow.

2 Shear flow and extensional flow

When testing a material you can can do a few things. You can stretch it and compress it, you can twist it and you can shear it. We focus here on the most common tests namely stretching and shearing, as depicted in Fig. 1.

Stretching and shearing have each there own dimensionless parameter describing the deformation. The deformation when performing an uni-axial elongation or stretch, is given by Chauchy or engineering strain ε which is defined as the increase in length Δl per unit starting length L_0 , see Fig. 1a:

$$\varepsilon = \frac{\Delta l}{L_0}.$$
(1)



Fig. 1: (*a*) Definition of the shear directions in the simplest flat-Couette (or plate-plate) geometry: 1, flow direction; 2, gradient direction; 3, vorticity or neutral direction. (b) Shear flow can be decomposed in compressional-extensional flow and rotational flow.

For the shear strain deformation is given by γ , which is the ratio of the displacement of the wall, Δl , in direction x, relative to an opposing parallel wall at distance h in direction y, see Fig. 1c:

$$\gamma = \frac{\Delta l}{h}.$$
(2)

The force F needed to deform a system is proportional to the area A it is applied to, which depends on the geometries that are used, see below. Therefore one always uses in rheology the stress σ . In the case of uni-axial deformation the stress is simply given by

$$\sigma = \frac{\mathbf{F}}{A}.\tag{3}$$

With this definition we can now quantify the rheological response of a material introducing Young's modulus E as the proportionality between the stress and strain:

$$E = \frac{\sigma}{\varepsilon}.$$
 (4)

This is basically Hook's law for solids where the Young's modulus plays the role of the spring constant.

Things become more complicated when the material start to flow during stretching, which will give a mess. Therefore, if a sample has the tendency to flow, one is better advice to characterize the sample by performing a shear experiment. However, in this case we cannot simply define a stress, but rather we need to consider a stress tensor. To appreciate this complication, consider that the deformation of a small volume element can be written as the gradient in the deformation field u_x , so that $\gamma = \frac{\partial u_x}{\partial y}$. A flow $\mathbf{v} = v_x \hat{x}$ will be induced when the sample is continuously strained. The flow rate depends on the distance from the moving wall y so that a shear rate can be defined as $\dot{\gamma} = \frac{\partial v_x}{\partial y}$. Thus, the y-direction is called the gradient direction. The third direction is unperturbed and is called the neutral direction, or also the vorticity direction since it is set by $\nabla \times \mathbf{v}$. The different directions are schematically shown in Fig. 1a. The resulting deformation rate tensor Γ describes the rate of deformation in three dimensions. For simple shear flow, left hand side of 9, the shear flow field can be decomposed in two components: a

combined compressional and extensional flow field, which is responsible for deformation of structure, and a rotational flow field, which is responsible for reorientation of structure, see the right-hand side of Eq. 9 and Fig. 1b:

$$\mathbf{\Gamma} = \dot{\gamma}\hat{\mathbf{\Gamma}} = \dot{\gamma} \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \frac{\dot{\gamma}}{2} \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \frac{\dot{\gamma}}{2} \begin{pmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} .$$
(5)

When the flow field is described by a tensor, then one should also describe its stress response by a tensor. In rheology one generally measures at least on component of the stress tensor, namely the component that is due to the force that is exerted in the 1 (= x) direction of the moving wall, divided by the area of the wall, and which is transmitted to the bulk of the sample in the 2 (= y) direction, normal to the wall. This stress component is called the shear stress and is denoted by σ_{12} , where the indices refer to the direction of the force (1) and the normal of the plane on which the force is exerted (2). σ_{12} is mostly referred to as the stress σ (or also often τ). Stresses in other directions are also important, especially for complex systems, but are often difficult to measure. We will come back to this later.

When the shear stress is proportional to the strain γ then the sample behaves like a solid and the proportionality constant G is the shear modulus, again similar to Hooke's spring:

$$\sigma_{12} = G\gamma. \tag{6}$$

When the shear stress is proportional to the rate of deformation $\dot{\gamma}$ then the sample behaves like a fluid and the proportionality constant η is the viscosity, which Newton introduced is mimicking 'hooks' between lamellar layers of fluid :

$$\sigma_{12} = \eta \dot{\gamma}. \tag{7}$$

As was suggested in the introduction, the answer to the question if a material is a solid or a liquid depends on the force that is exerted on the material, as well as the time given to respond. Often a material is not a solid and not a liquid, but something in between. This means that G in principle depends on the time the system is given to respond to the exerted field. There are many ways to find this functionality of G(t) and E(t) for these so called visco-elastic materials, but the most straightforward way is by subjecting the material to an oscillatory deformation. In the case of shear flow this is given by a time-dependent strain of

$$\gamma(t) = \gamma_0 \sin(\omega t). \tag{8}$$

This is the so called dynamic test. When the strain amplitude γ_0 is not too high, so the sample does not change its character, then the response can be written as

$$\sigma_{12}(t) = G^* \gamma_0 \sin(\omega t + \delta) = [G' \sin(\omega t) + G'' \cos(\omega t)] \gamma_0.$$
(9)

The sample can thus be characterized by the dynamic modulus G^* and the phase angle δ , or equivalently by the storage modulus G' and loss modulus G''. The system behaves like a solid when $G' \gg G''$ and $\delta = 0^\circ$, and like a fluid when $G' \ll G''$ and $\delta = 90^\circ$. The system is visco-elastic when $G' \approx G''$ or $0 < \delta < 90^\circ$, see Fig. 2. In the same way a complex Youngs modulus can be defined in an oscillatory stretch experiment. Thus, rheology is an often used to probe the state of a material.



Fig. 2: Top panel: Solid response of the stress (red line, like the Hookian knot) which is proportional to an oscillatory strain field (solid blue line). Bottom panel: Fluid response of the stress (red line, like the Newtonian hooks) which is proportional to an oscillatory strain rate field (dashed blue line). Middle Panel: viscoelastic response lies in between both responses.

3 Scaling time

In the introduction the Deborah number De was defined as the ration between the relaxation time of the system and the time during which the system is probed. When performing a dynamic experiment, as described above, the probing time is varied simply by varying the frequency ω in Eq. 9. Therefore, in case of dynamic experiments De takes the form $De = \omega \tau_{system}$. Hence, the system has no time to relax when $1/\omega \ll \tau_{system}$ and the system will behave as a solid, while it behaves as a fluid when $1/\omega \tau_{system}$.

When the system is probed with a constant shear rate, then the probing time is infinite. Therefore another scaled parameter is introduced, relating the applied shear rate with the relaxation time. Such a parameter is needed as the implicit assumption in Newton's law for fluids is that the fluid behaves like a continuum which can be described by a single friction coefficient, namely the viscosity η . The continuum does not change its properties when subjected to shear flow because shear flow does not effect the structure between the particles constituting the fluid. However, to satisfy this condition, the memory of the fluid should be very short compared to the applied deformation rate. Thus, the applied shear rate needs to be compared to the relevant relaxation time of the system τ_{system} in order to estimate the effect of shear flow. This scaled shear rate is called the Péclet number and is defined as $Pe = \dot{\gamma}\tau_{system}$. Changes in the intermolecular structure of the solvent molecules are very quickly lost for simple fluids, due to the diffusion of the molecules, so that the relaxation time is very fast and Pe << 1 for shear rates that can be reasonable easily be applied in the lab. For colloidal spheres, the diffusion rate D is much slower and given by the famous Stokes-Einstein relation (see also Eq. 2 in Chapter A3):

$$D_0 = k_B T / 6\pi \eta R \,, \tag{10}$$

where k_B is Boltzmann's constant, R is the radius of the particle and η is the viscosity. Thus the time for a particle to diffuse its own diameter is $\tau = 1/R^2D = 6\pi\eta R^3/k_BT$. In order for the shear rate to take effect it has to be of the same time scale, so $Pe = \tau_{system}\dot{\gamma} = (6\eta\dot{\gamma})/\langle k_BT/R^3\rangle \approx 1$. For a colloid of $R = 500 \ nm$ in water this means that a shear rate of $0.6 \ s^{-1}$ is sufficient to compete with Brownian motion. When considering colloidal rods, shear flow will compete with the rotational motion of the rods. For slender rods of $1 \ \mu m$ the rotational diffusion in water at room temperature is $D_R = 0.062 \ s^{-1}$. So in order to have an appreciable alignment, a shear rate of $\dot{\gamma} > \mathcal{O}(10 \ s^{-1})$ needs to be applied.

So far we considered particles that can freely move. However, in complex fluids particles will interact with each other. This will affect the dynamics of the particles dramatically. When the applied rate is faster than the new 'effective' relaxation, then the structure of the fluid will be affected and therefore also the rheological response. Thus, the response will become non-linear. Examples of such non-linear behavior will be given in section 6. The shear rate is now scaled by an effective diffusion, so the ratio is called the effective Péclet number.

For polymeric systems the Weissenberg number is generally used, given again by $Wi = \dot{\gamma} \tau_{sustem}$. It relates the effect of flow to the degree of anisotropy or orientation generated by the deformation of the constituent particles, i.e. polymer coils. The Weissenberg number also measures the ratio between elastic stress and viscous stress. Without interactions the relaxation of a polymer is given by the Rouse time λ_{Rouse} , which describes the relaxation time of the polymer conformation. In polymer melts, however, the relaxation is characterized by the time a polymer needs to diffuse over its full length through the 'tube' caused by the presence of other polymers. As this motion is much like the movement of a snake through the grass, it was called 'reptation' by de Gennes, who derived that the reptation time $\lambda_{rep} \propto M$, where M is the molecular weight of the polymer [4]. Hence, when the shear rate is high, the system cannot relax so Wi is high, the entanglements will act like knots until they break, the polymer will act as a Hookian solid with a high anisotropy in the system due to stretching of the polymers. In an oscillatory experiment, when the applied frequency is low and therefore De is high, the knots will relax and act rather as Newtonian hooks, between the sliding layers and the response will be of a fluid. Thus the question is when do Newtonian hooks become Hookian knots? We will come back to this in section 5.

4 Measuring geometries

In choosing a device for performing a rheological experiment one should consider that there are two basic options: either the stress is measured which results from straining system or, vice versa, the strain is measured which results from applying a stress to a system. This is a fundamental difference, as in stress controlled experiments the sample will 'decide' if it will flow, while in strain controlled experiments flow is enforced upon the sample, which might change its structure. One should always keep this difference in mind. Standard geometries that are used to impose shear flow are depicted in Figure 3. The geometries can be categorized in pressure flow, where shear is generated by pushing sample through a channel, and drag flows, where shear is generated between a moving and a fixed wall.



Fig. 3: Shear geometries for drag flow (a,b), pressure driven flows (c) and extensional flow. In (b) the three principle axes are indicate: the flow axis (\mathbf{v}), the gradient axis ($\nabla \mathbf{v}$), and the vorticity axis ($\nabla \times \mathbf{v}$). The same for the indicated cross section in (c).

Pressure-drop flow is fairly easy to set up, as one only needs to push through the fluid and measure the pressure drop over the geometry. Clearly pressure-drop experiments are stress controlled. To obtain the viscosity from a capillary experiments, one only needs to measure the volumetric flow rate and the pressure drop. However, flow is parabolic throughout the cell, so the shear rate is zero in the middle and linear only close to the wall, see Fig. 3c. This means that reliable data can only be obtained when the viscosity does not change with shear rate, though in principle local shear rates can be measured using Particle Imaging Velocimetry.

Drag geometries allow for both stress and strain controlled experiments, depending on the machine, and the shear rate is well defined. The most homogeneous flow is obtained with a flat couette, or plate-plate cell, but it has the disadvantage that under continuous shear the plates move away from each other, reducing the contact surface and exposing the sample to air. Therefore, drag-flow devices are always rotating devices. Both the cone-plate and the concentric cylinders (couette flow) have a fixed shear rate at every position. There are however disadvantages. In the cone-plate geometry sample at the rim can easily become unstable, while in the couette flow there is no continuous stress as the effective surface is decreasing throughout the gap and therefore the stress is decreasing. moving from the inner to the outer cylinder. Couette geometries have more surface and therefore more torque but loading of highly viscous samples can be difficult and a lot of sample is required. Loading the cone-plate is easier and generally less sample is needed. In conclusion, the choice of the geometry depends on many different aspects. For details on the advantages and disadvantages of the different shear geometries from a pure rheological point of few, please refer to Chapter 5 of Ref. [1].



Fig. 4: Dynamic frequency sweeps of the storage modulus and loss modulus of (a) Polystyrene melts for molecular weight varying from 34 kd (PS-1) to 2540 kd, see Ref. [5]; (b) 92.6 kd polystyrene, showing the cross over between fluid to solid like behavior, see Ref. [6].



Fig. 5: Dynamic frequency sweeps of the storage modulus and loss modulus of F-actin dispersions at a concentration of 1.4mg/ml. The storage (triangles) and loss (circles) shear moduli were measured at a strain amplitude of 2% without (a) or with (b) gelsolin at a molar ratio to actin of 1:500, see Ref. [7].

5 Examples of linear Rheological characterization

Having defined why and how rheology can be used, we can now proceed to applications. Here I want to focus on experiments to test the nature of the complex material and, as mentioned in section 3, one prominent example of a complex material materials is polymer melts. In Fig. 4a and b dynamic frequency sweeps are plotted where G' and G'' (see Eq. 9) are measured as a function of frequency for polymers melts of polystyrene with varying molecular weight. Indeed, a plateau in the storage modulus G' forms when increasing the polymer length, which stretches out over many decades for the longest polymer. This is due to the fact that with increasing length, also the length of tube through which the polymer reptates increases, such that the reptation time increases. Note that in order to obtain this plot the so called timetemperature superposition trick has been applied, as no rheometer is able to cover this huge range of frequencies. Assuming that the relaxation time decreases with increasing temperature, either because of an temperature dependent activation energy or viscosity, the frequency sweeps can be shifted along the frequency axis by a multiplication with a shift factor a_T . In Fig. 4b the overlay is plotted of both the G' and G'' curves. This is very instructive as now it is easy to see that there is a well defined frequency where both curves cross, which is exactly at the lowest relaxation time, and hence we can read from this curve the reptation time. For a polymer with a molecular weight of 2540 kd this can be as slow as $4 \cdot 10^{-5}$ Hz at room temperature.

Contrary to polymers, the most important biological building blocks are very stiff. One of the main building blocks is filamentous-actin. In Chapter C.3 the physiological importance and physical characteristics of the system will be explained in detail. For now, we just note that F-actin is a very stiff filament. Dispersions of F-actin in water at a concentration as low as 1 mg/ml, so a volume fraction of around 0.1 % as compared to the 100 % polymer melts, display no relaxation at all within the frequency regime probed, as can be appreciated from Fig. 5a. Of course, for this system time-temperature superposition is hardly possible as F-actin has a very limited temperature range in which is functional. Compared to the polymer melt the modulus is so low that loss modulus is for most frequencies even below the torque limit, below which the rheometer cannot produce a reliable number. The size distribution can be tuned, similar

to the polymeric system, by the addition of Gelsolin, which acts as an initiator for F-actin polymerization. Thus, at the same volume fraction but for much shorter lengths the system can now actually relax as both G' and G'' decay by lowering the probing frequency. However, G' > G'' over the full range, such that still the system could be considered as a solid, even though the modulus is now extremely low.

6 Examples of non-linear Rheological responses

In the previous examples rheology was used to probe dynamics and mechanics of the equilibrium system. Thus, the structure of the fluid is supposed to be undistorted by flow. In case of the F-actin filaments this is obviously an issue as the modulus is low and as an experimentalist one is tempted to apply high amplitudes. However, due to the extremely slow relaxation, flow would immediately introduce alignment resulting in an altered flow behavior. General, when flow affects structure the response to flow will alter and hence the response will be non-linear. We will now introduce strain and strain rate dependent phenomena that are commonly encountered.

6.1 Strain hardening and softening

When a solid material is subjected to a small deformation nothing will change. Classical experiments by Evans and Lebow [8] show, for example, that upon stretching a dry bone, its modulus will be constant as the stress-strain curve is a straight, see Fig. 6a, left, si that Hook's law Eq. 6 is valid. However, when the bone is wet, it will yield when it is stretched more than a fixed yield strain, and the curve will bend down resulting in a softening of the material, see Fig. 6a, right. Generally, a small deformation is a strain γ for which the linear relation in Eq. 6 holds, just as is the case for the bone. If the strain is to large and no time is given to relax then there are basically two options. The first option is that it will oppose the deformation causing a hardening of the material so that the stress-strain curve is higher than the linear curve, see the blue line in Fig. 6b. This is called strain hardening. The second option is that the material will yield causing a softening of the material, so that the stress-strain cure is lower than the linear curve, see the red line in Fig. 6b. This is called strain softening. In principle, yielding can set in immediately, but for many materials there is a finite yield strain γ_{yield} . These are plastic materials, see the green curve in Fig. 6b.

There are many examples of strain hardening and one of the most relevant examples is again due to F-actin solutions. Fig. 6c shows curves of the stress vs an applied strain for different applied shear rates. This hardening depends, however, strongly on the rate with which the deformation is being applied. Clearly there is less hardening when more time is given to relax, so at low applied shear rates. Fig. 9a of Chapter C3 treats the opposite experiment, as here the strain is measured as a function of applied stress for dispersions of F-actin, but also for more stiff filaments such as microtubuli. The dispersion with the very stiff microtubuli is the most ductile, which means that the deformation for the same applied stress is much higher than for F-actin.

Straining of a system that has no time to relax can not continue indefinitely. There are several possible ends to the story. The most dramatic end is that the system fractures, as is the case for the bone where the X in the curves indicate the point of fracture. Fracturing is of course an extremely non-linear event, which requires a whole different approach to understand. When there are relaxation mechanisms, then the system can also undergo yielding and strain softening, as is the case for wet bones, see Fig. 6a, right. F-actin dispersions display dramatic yielding



Fig. 6: Stress-strain curves of (a) dry and a wet bone, see Ref. [8]; (b) a Hookian solid (black curve, see Eq. 6), a strain hardening solid (blue curve) and a strin softening solid (red curve); (c) a F-actin dispersion with an average length of 21 μ m at a concentration of 0.4mg/ml at 21 °C (solid lines) and 25 °C (dashed lines) for $\dot{\gamma} = 0.05, 0.1, 0.2, 0.4 \ s^{-1}$, see Ref. [9]; (d) a F-actin dispersion at a concentration of 0.15mg/ml up to high γ . The shear rate is increased every 5 min, see Ref. [10].

and softening, as can be inferred from Fig. 6d. Here, strain softening is observed between strain $\gamma = 2$ and 8, so beyond the strains used in Fig. 6c, at very low slow shear rates.

A yield strain can also be defined in a dynamic test. As mentioned above, torque is often an issue so that one would like to use a larger strain amplitude to perform the dynamic frequency sweep discussed in Fig. 5. It is good practice to first perform a dynamic strain sweep at a fixed frequency to define a linear regime, as plotted in Fig. 7. Traditionally, the strain amplitude where G' and G'' cross is called the dynamic yield strain γ_{yield} , as at this point the system becomes predominantly fluid. In the case of F-actin this is just below a strain amplitude of $\gamma_0 = 1$, see Fig. 7a, so below the steady test in Fig. 6d, which had a much lower concentration. A strain sweep of a slab of a Bovin brain is plotted in Fig. 7b and c a. In this case there is no cross over, so the brain is predominantly solid although it softens with increasing strain and with decreasing frequency.

The response for $\gamma_0 < \gamma_{yield}$ will already be non-linear, in the sense that it will not be anymore sinusoidal, which raises the question what should be considered as a small strain amplitude. Since about a decade it became fashionable to perform such 'Large Amplitude Oscillatory Shear' (LAOS) experiments, but the interpretation of the response needs to be done with caution [13]. Nice examples of response to LAOS for F-actin dispersions can be found in Ref. [9], where within one oscillation strain hardening and softening can be identified, while over several cycles the system softens. The hardening can actually not be seen in Fig.7, as G' displays a monotonic decay. It is interesting to read in Ref. [11], the possible reasons, which are manifold as the system as well as the experiment are hard do handle.

Where the strain hardening of entangled F-actin can still be described by theory, see Chapter C4, the process of strain softening involves a complete reorganization, as can be inferred from the flat hairpins shape of the filament as we will discuss in section7, for which there is at present no theoretical description. This reorganization involves so called non-affine deformation, which means that deformation also takes place in directions other than that of the applied deformation. This results additional in terms in the stress tensor other than σ_{12} . Forces that are induced in directions normal to the applied deformation, the normal forces, are often essential for the understanding of complex systems, but notoriously difficult to measure. This also couples back to the Weissenberg number which measures the ratio between elastic stress and viscous stress.

6.2 Shear thinning and thickening

Thinking of biological systems, like our body, we could be satisfied with what has been described so far as the word 'body', which suggests a solid. Of course we have our bodily fluids that need to flow, but there is also flow at very small length scales. F-actin filaments, important building blocks of the cyto-skelleton, can form striking non-equilibrium patterns when they to cytoplasmic flows within animal embryos generated by molecular motor activity [14]. Flow can be desirable when material is transported from A to B, but also a problem. A good example is blood which is needed as a transport medium, where highly concentrated dispersions of deformable cells are pushed through narrow channels to supply oxygen to the body [15]. But blood is also needed as a base for the formation of a new solid, in case of an open wound. Hence, the dual character of "non-Newtonian" complex fluids is perfectly exemplified by blood, though the examples in industrial applications are manifold, which is beyond the scope of this chapter.

The reason for complex flow behavior behavior was discussed in section 3, namely that the shear rate can be faster than the relaxation dynamics in the system, so when Pe > 1. Again



Fig. 7: Dynamic strain sweeps for (a) a *F*-actin dispersion at a concentration of 1 mg/ml, showing a cross over of the sotrage G' (•) and loss G'' (•) modulus at a strain amplitude of $\gamma \approx 0.7$, see Ref. [11] (b,c) a slab of brain tissue, see Ref. [12].





Fig. 8: Cartoons of flow curves in the viscosity vs shear rate representation (a) and stress vs shear rate representation (b): Newtonian (black); shear thinning (red); shear thickening (blue). The green curve represents a yield stress fluid, which has infinite viscosity at zero shear rate.



Fig. 9: Flow curves of (a) blood (stress vs shear rate), see Ref. [8]; (b) two F-actin dispersion (viscosity vs shear rate), see Ref. [10].

there are two options. The first option is that it will oppose the deformation causing a hardening of the material so that the viscosity increases with increasing shear rate. This shear thickening behavior is apparent when taking a flow curve where the stress, and hence the viscosity, is measured as a function of shear rate as plotted in Fig. 8. Here the blue line displays the increased viscosity as compared to the constant black Newtonian response, given by Eq. 7. The second option is that the material will yield causing a thinning of the material, so that the viscosity decrease with increasing shear rate, see the red line in Fig. 8. This is called shear thinning. In principle, yielding can set in immediately, but for many materials there is a finite yield stress σ_{yield} . These are, again, plastic materials, see the green curve in Fig. 8b.

There are many more or less phenomenological models that describe the so called constitutive relation between the stress and the shear rate for complex fluids. A very simple example is the Bingham fluid, which describes the flow behavior of a plastic material [16]:

$$\sigma = G\gamma \qquad \text{for} \quad \sigma < \sigma_{yield} \sigma = \eta \dot{\gamma} + \sigma_{yield} \quad \text{for} \quad \sigma \ge \sigma_{yield}$$
(11)

This equation does not describe, however, one of the best studied but still not fully understood system, namely blood. The one thing that blood and tomato ketchup have in common, except the color, is that they are both yielding fluids. Fig. 9a displays a flow curve (in this case stress vs rate) of normal blood in a peculiar way, namely plotting $\sigma^{1/2}$ vs $\dot{\gamma}^{1/2}$. The reason for this way of plotting is that for low shear rates the flow can be described by Casson's two parameters constitutive relation [17], which gives a smooth transition between yielding and Newtonian flow, see Eq. 12:

$$\dot{\gamma} = 0 \qquad \text{for} \quad \sigma < \sigma_{yield}$$

$$\sigma^{1/2} = \sigma_{yield}^{1/2} + (\eta \dot{\gamma})^{1/2} \quad \text{for} \quad \sigma \ge \sigma_{yield}.$$
(12)

Constitutive relations have been developed on phenomomological level up to a molecular level, see also Ref. [2]. Generally, it is a prerequisite to know the interaction between the particles that constitute the complex fluid at the microscopic level in order to develop a fundamental base for phenomenological constitutive relations as given in Eq. 11 and 12. Constitutive relations have been developed for colloidal rods [18, 19], with moderate success, see section 7. For slightly more complex like F-actin the situation is worse. The shear thinning as shown in Fig. 9b for two concentrations can not be theoretically described as it is a highly non-linear response, but we will have some microscopic hints in section 7. For blood the situation is even more complex, given the biological complexity of the interactions. Hence, one needs to resort to microscopic experiments, as described in section 7 or simulations, as described in Chapter E2.

7 In situ rheology

Rheologists are masters in extracting information on the smallest relevant length scale performing experiments on the largest possible length scale, namely the mechanical relations of a bulk sample. Dynamic oscillatory experiments can be complimented by, for example, dynamic light scattering or dielectrics. It is experimentally more challenging to relate the non-linear rheological response to structural changes. In order to make this link, *in situ* experiments need to be



Fig. 10: (a) In situ SANS experiment where the orientational distribution, in this case of rodlike viruses, can be probed in 3D by directing the neutron beam along the gradient and vorticity direction, see Ref. [20]. (b) A counter-rotating cone-plate shear cell in combination with an ultra-fat confocal microscope can be is used to take 3D stacks of images, showing a small fraction of labelled actin filaments embedded in a dark background of unlabelled F-actin (b). Here the green line is an example of a tracked filament with local tangents and binormals indicated by the arrows. Scale bar, 10 μ m, see Ref. [10]. The ellipsoid depicts the orientational distribution that can be obtained by (a) and (b), of which the axes mimic the eigenvalues λ_i . (c) a confocal cross section of labeled red blood cells showing break up and tumbling of rouleaux. Data taken by O. Korculanin at FZJ.



Fig. 11: (a) Head actuator used for stimulating low-frequency shear vibrations in the brain. (b) Illustration of the outcome of a single-slice multifrequency-MRI experiment of a 47-year old male volunteer. Complex modulus images G' and G'' denote real and imaginary part of $G(x, y, \omega)$. Drive frequencies used in experiments are given above the columns, see Ref. [21]

performed. Fig. 10 shows two examples: in situ Small Angle Neutron Scanning (SANS), see also Chapter A4, and confocal microscopy, see also Chapter A1. In choosing the technique one has to consider, as always, the length and time-scale, but also contrast. For example, labeled F-actin filaments in a dispersion of unlabeled F-actin can be imaged by ultra-fast confocal microscopy, taking stacks of images to reconstruct the full 3-d trace, see Ref. [10]. Of course the filament needs to stay in the field of view during the recording and therefore a counter-rotating cone-plate shear cell is designed. The 3-d images already tell the story of strain softening: the filaments, which are entangled in all three direction at zero strain, reform into flat hairpins which orient such that they can slide past each other, without any entanglement in the gradient direction. This explains the strain softening observed in Fig. 6d. Analysis of the trace also gives detailed knowledge on the shear-induced curvature as well as stretching of the filaments. The latter can be converted in an 3d orientational distribution which is biaxial, as the corresponding eigenvalues λ_i are all different. Moreover, the largest eigenvalue increases with increasing shear rate thus explaining the shear thinning shown in Fig. 9b.

Where the level of detail is the advantage of microscopy, the disadvantage is statistics. With SANS we can access the shear-induced orientational distribution of rod-like viruses(length 880 nm, diameter and 6.6 nm) with one 2-d scattering image. This distribution is directly proportional to the intensity profile at a well-chosen scattering angle, as indicated by the dashed lines in Fig. 10a. Again it can be shown that this distribution is biaxial [20], using two scattering geometries, directing the neutron beam through the gradient direction as well as the vorticity direction. The latter is non-trivial and requires tilting of the couette cell, rendering on -line rheology impossible [22]. Of course using SANS all contributions are added up, so that no distinction can be made between curved and stretched segments. However, the afore mentioned viruses are so stiff that this is not an issue. Still, even though this system can be considered as almost ideal rods, there is up to date no theoretical description that full describes its flow behavior [20].

Fig. 10c also shows a confocal cross section of rouleaux. Rouleaux are columns of piled up red blood cells, which cause the yield stress in blood. When sheared, columns will first break apart and then fragment into single cells. To get a complete picture of the yielding process of blood one should consider that blood flow is always a Poiseuille flow, who as a pioneer of blood flow, described the transport through pipes.

Finally, there is a nice resent example of truly *in situ* rheology, where actually the imaging method itself is used to obtain G' and G''. The brain of a 47-year old male volunteer is actuated at different frequencies, while a trigger Magnetic Resonant Imaging machine takes data. As signal is low, data is obtained at different delays, using a trigger from the actuator. Again a Laplace transform is used to obtain $G(x, y, \omega)$ spatial resolved in the frequency domain, see Ref. [21]

8 Micro-rheology and fluidics

There are two issues that make all that we discussed so far a bit problematic if one envisions its use in biology. The first issue is that all the experiments described so far require quite a lot of sample, ranging from at least 200 μl in the smallest cone-plate geometries to about 20 ml in double-wall couette geometries. Many biological samples, however, can only be produced in very small amounts, which makes that rheological experiments become impractical. Maybe even more important is the heterogeneity in many biological samples. Basically one is often



Fig. 12: Dynamic frequency sweeps, as in Fig. 5, using micro-rheological techniques. Laser interferometry on *F*-actin at 1 mg/ml without (squares) and with (triangles) cross-linking: (a) Storage modulus *G*' and (b) loss modulus *G*'', see Ref. [23]. (c) Diffusive wave spectroscopy also at 1 mg/ml *F*-actin, see Ref. [11]. The solid lines give the 3/4 power law. (d) Optical tweezer experiments on single fibroplasts show a cross-over of *G*'' and *G*'''. The insets display fluorescently labeled actin distribution inside the cell, see Ref. [24]



Fig. 13: (a) Images of molecular configurations spaced every 0.13 s at the highest strain rate investigated ($\dot{\epsilon} = 0.86 \ s^{-1}$). The configurations are clasified as (from top to bottom) dumbbell, kinked, halfdumbbell, and folded. The molecular extension of the last image in the first row is 13.9μ m. Sketches of possible molecular configurations are shown on the left. (Inset) A schematic illustration of the flow pattern where the red rectangle indicates the observation region. (b) Individual traces of length versus residence time for dumbbell or folded conformations, see Ref. [25]. (c) Micro-fluidics is used to make an artificial microvascular network, in which the viscosity as a function of red blood cell concentration can be measured for increasing (average) shear rates (d), see Ref. [26]

interested in the local rheology of a cell or somewhere inside a cell and not of the full tissue as described in Fig. 7b. A solution to both issues is the use of tracer particles of which the mean square displacement is registered. The tracers either diffuse because of its Brownian motion (passive) or are driven by some external source (active).

Passive micro-rheology has been theoretically described in e.g. Ref. [27], where it is explained how one can obtain information on G' and G'' as a function of the frequency by a Laplace transformation of the time-dependent mean square displacement. This principle has been applied to F-actin dispersions using e.g. interferometric microscopy [28], Fluorescence Correlation Spectroscopy [29], as treated in Chapter A3, Diffuse wave spectroscopy [11], and laser-interferometric [23]. The latter two experiments demonstrate another huge advantage of this approach, namely that it is possible to access a much broader frequency range than is feasible with classic rheology. This is is important as for biological systems the time-temperature superposition trick does not work. Compared to Fig. 5 both experiments show a second cross over as at high frequencies the shear modulus is entirely controlled by the relaxation of individual polymer chains.

The passive micro-rheological approach does have a few drawbacks. For one, a probe particle is inserted in the medium, creating its own space, which might affect the result. Another disadvantage is that applied forces are purely Brownian and therefore very small. A way around this problem is to use micro-tweezers. For example, it is possible to manipulate a magnetic probe particle with a varying magnetic field, in case of Ref. [30] again in F-actin dispersions. Exploiting differences in refractive index, one can also build optical tweezers [31] that can be used stretch actual cells by tweezering the cell membrane [24]. Fig. 12d demonstrates how in this way a frequency sweep of a fibroplas cell can be obtained, displaying a typical cross over of G' and G'' and therefore typical relaxation time can be obtained. Rheology on the smallest length scale, namely on single filaments, can be performed using the same tools. Magnetic tweezers can be used to perform stretching experiments [32], in principle like in Fig. 3d, and torques experiments to study the mechanics of DNA. Another nice example of single filament experiments is described in Fig. 9b of Chapter C3. Here Atomic force microscopy is used to show that F-actin filemantes can soften when being stretched.

Chain stretching of DNA has also been observed on a single particle level using micro-fluidics. This is a technique where channels structures are designed and imprinted in glass or curable polymer to produce any required flow. It is for example feasible to create a purely extensional flow by crossing two channels, thus basically isolating one of the flow components of shear flow, as explained Fig. 1d. Fig.13a-b shows fluorescently labeled DNA that evolves into different shapes in extensional flow [25]. Microfluidics is also frequently used in the study of blood flow as artificial microvascular networks can be constructed, see for example Fig. 13c. Clearly, the break up of rouleaux is dramatic in such systems, but the fact that the viscosity still depends on the (averaged) shear rate examplifies the importance of the interactions, as well as its strong dependence on the concentration of red blood cells, see Fig. 13d [26].

A review on micro-rheology can be found in Ref. [33], on the use of micro-fluidics to study physics of flow on the nanoliter scale in Ref. [34], and on both in Ref. [35].

9 Concluding remarks

When the reader by now knows what is meant with a loss and storage modulus, with a yield stress, and when the reader also knows how these parameters can be measured even locally in

the system of interest, then the exercise of reading this chapter was a success. What will be clear is that rheological is not only a set of techniques to probe and characterize a system, but also to understand fundamental process in living systems, such as blood flow. When and how things happen all depends on the time is given to the system, the force that is applied and the micro-structure of the system. Clearly the number of experiments that still can be thought of and designed are numerous as well as the number of biological systems that are still waiting to be studied in this way.

References

- [1] C. W. Macosko, *Rheology-Principles, Measurements, and Applications* (Whiley-VCH, 1994).
- [2] R. G. Larson, The Structure and Rheology of Complex Fluids (Topics in Chemical Engineering) (Oxford University Press, USA, 1998).
- [3] M. Lettinga, *Phase behavior of colloidal dispersions in shear flow* (Wiley, Hoboken, New Jersey, 2016), p. 408.
- [4] P. G. Degennes, Molecular Crystals and Liquid Crystals 12(3), 193 (1971).
- [5] A. Schausberger, G. Schindlauer, and H. Janeschitzkriegl, Rheologica Acta 24(3), 220 (1985).
- [6] R. H. Colby, Thesis (1985).
- [7] P. A. Janmey, S. Hvidt, J. Ks, D. Lerche, A. Maggs, E. Sackmann, M. Schliwa, and T. P. Stossel, Journal of Biological Chemistry 269(51), 32503 (1994).
- [8] F. G. Evans and M. Lebow, Journal of Applied Physiology 3(9), 563 (1951).
- [9] C. Semmrich, R. J. Larsen, and A. R. Bausch, Soft Matter 4, 1675 (2008).
- [10] I. Kirchenbuechler, D. Guu, N. A. Kurniawan, G. H. Koenderink, and M. P. Lettinga, Nat Commun 5, 5060 (2014).
- [11] A. Palmer, T. G. Mason, J. Y. Xu, S. C. Kuo, and D. Wirtz, Biophysical Journal 76(2), 1063 (1999).
- [12] L. E. Bilston, Z. Z. Liu, and P. T. Nhan, Biorheology **34**(6), 377 (1997).
- [13] S. A. Rogers, J. Rheol. 56(5), 1129 (2012).
- [14] S. Ganguly, L. S. Williams, I. M. Palacios, and R. E. Goldstein, Proceedings of the National Academy of Sciences of the United States of America 109(38), 15109 (2012).
- [15] E. W. Merrill, Physiological reviews **49**(4), 863 (1969).
- [16] E. Bingham, Bulletin des Bureau of Standards 13, 44 (1916).
- [17] N. Casson, Flow Equation for Pigment Oil Suspensions of the Printing Ink Type. (Pergamon Press, London, U.K., 1959).

- [18] M. Doi and S. F. Edwards, J. Chem. Soc., Faraday Trans. II 74(5), 918 (1978).
- [19] J. K. G. Dhont and W. J. Briels, Colloid Surface A 213(2-3), 131 (2003).
- [20] C. Lang, L. Porcar, J. Kohlbrecher, and M. Lettinga, Polymers 8, 291 (2016).
- [21] I. Sack, B. Beierbach, J. Wuerfel, D. Klatt, U. Hamhaber, S. Papazoglou, P. Martus, and J. Braun, Neuroimage 46(3), 652 (2009).
- [22] M. W. Liberatore, F. Nettesheim, N. J. Wagner, and L. Porcar, Physical Review E 73(2) (2006).
- [23] G. H. Koenderink, M. Atakhorrami, F. C. MacKintosh, and C. F. Schmidt, Phys. Rev. Lett. 96, 138307 (2006).
- [24] F. Wottawah, S. Schinkinger, B. Lincoln, R. Ananthakrishnan, M. Romeyke, J. Guck, and J. Kas, Physical Review Letters 94(9) (2005).
- [25] T. T. Perkins, D. E. Smith, and S. Chu, Science 276(5321), 2016 (1997).
- [26] N. Z. Piety, W. H. Reinhart, J. Stutz, and S. S. Shevkoplyas, Transfusion 57(9), 2257 (2017).
- [27] T. G. Mason and D. A. Weitz, Physical Review Letters 74(7), 1250 (1995).
- [28] F. Gittes, B. Schnurr, P. D. Olmsted, F. C. MacKintosh, and C. F. Schmidt, Physical Review Letters 79(17), 3286 (1997).
- [29] A. Bernheim-Groswasser, R. Shusterman, and O. Krichevsky, Journal of Chemical Physics 125(8) (2006).
- [30] F. Amblard, A. C. Maggs, B. Yurke, A. N. Pargellis, and S. Leibler, Physical Review Letters 77(21), 4470 (1996).
- [31] J. Guck, R. Ananthakrishnan, T. J. Moon, C. C. Cunningham, and J. Kas, Physical Review Letters 84(23), 5451 (2000).
- [32] C. Bustamante, J. F. Marko, E. D. Siggia, and S. Smith, Science 265(5178), 1599 (1994).
- [33] L. G. Wilson and W. C. K. Poon, Physical Chemistry Chemical Physics 13(22), 10617 (2011).
- [34] T. M. Squires and S. R. Quake, Reviews of Modern Physics 77(3), 977 (2005).
- [35] T. A. Waigh, Reports on Progress in Physics 79(7) (2016).

E 2 Modeling blood flow and primary hemostasis in microcirculation

Dmitry A. Fedosov

Theoretical Soft Matter and Biophysics Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Intro	oduction	2		
	1.1	Blood	2		
	1.2	Blood cells	2		
	1.3	Microcirculation	3		
	1.4	Primary hemostasis	4		
2	Met	hods and models	4		
	2.1	Blood cells	4		
	2.2	von Willebrand factor	6		
3	Results				
	3.1	Fahraeus & Fahraeus-Lindqvist effecs	7		
	3.2	Margination	8		
	3.3	VWF adhesion	9		
	3.4	VWF-platelet aggregates	11		
	3.5	Blood flow in microcirculation	11		
4	Conclusions				
A	A Smoothed dissipative particle dynamics				

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

1.1 Blood

Blood is a bodily fluid which performs a number of physiological functions, including the transport of oxygen and nutrients to cells of the body, removal of waste products, and circulation of vital molecules and cells. Furthermore, circulating blood is important for the organism's defense and immune response and plays an essential role in the tissue repair process. Abnormal changes in blood flow are often associated with a broad range of disorders and diseases which include, for instance, hypertension, anemia, atherosclerosis, malaria, and thrombosis. Therefore, understanding of fundamental rheological properties and dynamics of blood cells and blood flow is crucial for many biomedical and bioengineering applications, such as the development of blood substitutes, the design of blood flow assisting devices, and drug delivery. In addition, a detailed understanding of vital blood-related processes in health and disease is likely to result in the development of new effective treatments.

Blood is a suspension of erythrocytes or red blood cells (RBCs), leukocytes or white blood cells (WBCs), thrombocytes or platelets, and various molecules and ions in the blood plasma. RBCs take up about 45% of the total blood volume, WBCs around 0.7%, and the rest corresponds to blood plasma and its substances. One microliter of blood contains about 5×10^6 RBCs, roughly 5000 WBCs, and approximately 2.5×10^5 platelets.

1.2 Blood cells

Figure 1 shows a scanning electron micrograph of blood cells. Human RBCs have a relatively simple structure in comparison to other cells. They have a biconcave shape with an average diameter of approximately $8 \mu m$ and a thickness of about $2 \mu m$ [1]. A RBC membrane consists of a lipid bilayer with an attached cytoskeleton formed by spectrin proteins and linked by short filaments of actin to the bilayer. At the stage of RBC birth, the nucleus and other organelles, which are generally present in eukaryotic cells, are ejected, leaving behind a relatively homogeneous cytosol and no bulk cytoskeleton. RBC cytosol is a hemoglobin-rich solution, which



Fig. 1: A scanning electron micrograph of blood cells. From left to right: a human RBC, activated platelet, and a WBC. Source: The National Cancer Institute at Frederick (NCI-Frederick).



Fig. 2: 3D image of the microvasculature of rat's spinal cord. Reproduced with permission from Ref. [3].

is able to bind oxygen. Therefore, the main function of RBCs is oxygen delivery to the body tissues. RBCs are very deformable and can pass through capillaries with a diameter several times smaller than the RBC diameter.

WBCs are spherical in shape with a diameter ranging between 7 μ m and 20 μ m, and have one or multiple nuclei. Even though WBCs are stiffer than RBCs, they are able to undergo significant deformation when entering the smallest blood capillaries or transmigrating from blood into the surrounding tissue. WBCs are an important part of the body's immune system. They protect the body against invading bacteria, parasites, and viruses by killing these microorganisms through phagocytic ingestion and other antigen-specific cytotoxic mechanisms. There exist different types of WBCs (e.g., neutrophils, eosinophils, basophils, monocytes, and lymphocytes), each of which is designed for a specific task. WBCs may adhere to the vascular endothelium, which is the first step in their transmigration into the surrounding tissue and a part of the immune response.

Not activated platelets possess a disc-like shape with a diameter ranging between $2 \mu m$ and $4 \mu m$, and a thickness of about $0.5 \mu m$ [2]. These thin cells have a bulk cytoskeleton, but no nucleus, and play a pivotal role in hemostasis or the blood-clotting process. Upon activation, platelets change their shape by the formation of multiple tethers, as shown in Fig. 1, and become adhesive. The process of adhesion of numerous activated platelets to a blood-vessel injury site results in the formation of a clot (or a plug), which eventually stops bleeding.

1.3 Microcirculation

The microcirculation (or microvascular network) [4, 5] consists of the smallest vessels (i.e., arterioles, capillaries, venules) with diameters up to about $100 \,\mu$ m. Typically, the microcirculation geometry resembles a tree-like structure (see Fig. 2) formed by microcirculatory vessels, which provide the routes for the circulating blood and enable efficient exchange between blood and the surrounding tissues. Blood flow in microcirculation is extremely complex and diverse with dramatic changes in flow rates and patterns. This complexity is attributed to non-trivial vessel geometries and blood rheological properties [4, 6], which depend on the mechanical properties of blood cells, aggregation interactions, local hematocrit (RBC volume fraction), and flow conditions.

Difficulties with reliable *in vivo* measurements and the complexity of blood flow in microvascular networks in health and disease pose considerable limits on experimental investigations. In contrast, the capabilities of *in silico* modeling to describe the flow behavior of blood in the microcirculation are expanding. Realistic blood flow modeling has a great potential to provide the necessary theoretical tools to better understand blood flow in normal and diseased microcirculation. Some of these approaches will be discussed further.

1.4 Primary hemostasis

Hemostasis is the process of blood clotting, which is required to stop bleeding at an injury site and initiate wound healing. The start of the hemostatic process is usually referred to as primary hemostasis and corresponds to initial platelet-plug formation at an injury site. Low platelet counts may significantly affect the primary hemostasis, resulting in a substantial increase of bleeding times [7]. Thus, timely and frequent enough platelet adhesion is very important in primary hemostasis. Adhesion of activated platelets is mediated by several receptors at the cell's surface, which can interact with subendothelial components, such as collagen exposed at the damaged vessel wall [8]. Additionally, the adhesion can be supported by other proteins, including the von Willebrand factor (VWF) [7]. The role of VWF in mediating platelet adhesion becomes essential at high shear rates, when platelets are not able to bind to the injury site autonomously.

VWF is a large polymer-like protein made of repeating subunits [9]. It is a shear responsive molecule, because in equilibrium or at low shear rates, VWF remains in a globular configuration due to its internal associations [10], while at high enough shear rates, VWF is able to stretch [11]. In a globular configuration, VWF remains non-adhesive such that its adhesive sites are shielded [12], while in a stretched form, the adhesive sites of VWF get exposed and it can adhere to a damaged endothelium and platelets [10]. After adhering to the damaged endothelium at high shear rates, tethered VWF molecules are able to capture flowing platelets. The VWF-platelet interactions significantly slow down platelets, and facilitate their firm adhesion at the damaged vessel wall. Various VWF dysfunctions can lead to uncontrolled bleeding as in the von Willebrand disease or to spontaneous thrombotic events [13].

2 Methods and models

Modeling blood flow requires the representation of several main components, including fluid flow, blood cells, and suspended molecules (e.g., VWF). Fluid-flow modeling is referred here to a proper representation of the motion of blood plasma and cell's cytosol, resulting in hydrodynamic interactions between suspended cells. Both blood plasma and cell cytosol can be considered to be viscous Newtonian fluids. Fluid flow is modeled by the smoothed dissipative particle dynamics (SDPD) method [14, 15], which is a mesoscopic particle-based simulation approach. Detailed description of the SDPD method can be found in Appendix A.

2.1 Blood cells

RBC membrane is modeled by a set of points $\{\mathbf{x}_i\}$, $i \in 1...N_v$ which are the vertices of a two-dimensional triangulated network on the RBC surface [17, 18], as shown in Fig. 3. The vertices are connected by N_s edges which form N_t triangles. The potential energy of the system is defined as

$$U(\{\mathbf{x}_{i}\}) = U_{s} + U_{b} + U_{a} + U_{v}, \tag{1}$$

where U_s is the spring's potential energy to impose membrane shear elasticity, U_b is the bending energy to represent bending rigidity of a membrane, and U_a and U_v stand for the area and



Fig. 3: Mesoscopic representation of a RBC membrane by a triangular network of bonds. Reproduced with permission from Ref. [16].

volume conservation constraints, which mimic area-incompressibility of the lipid bilayer and incompressibility of a cytosol, respectively.

The spring's elastic energy mimics the elasticity of the spectrin network, and is given by

$$U_s = \sum_{j \in 1...N_s} \left[\frac{k_B T l_m (3x_j^2 - 2x_j^3)}{4p(1 - x_j)} + \frac{k_p}{l_j} \right],$$
(2)

where l_j is the length of the spring j, l_m is the maximum spring extension, $x_j = l_j/l_m$, p is the persistence length, k_BT is the energy unit, and k_p is the spring constant. The bending energy is defined as

$$U_b = \sum_{j \in 1...N_s} k_b \left[1 - \cos(\theta_j - \theta_0) \right],\tag{3}$$

where k_b is the bending constant, θ_j is the instantaneous angle between two adjacent triangles having the common edge j, and θ_0 is the spontaneous angle.

The area and volume conservation constraints are expressed as

$$U_{a} = \frac{k_{a}(A - A_{0}^{tot})^{2}}{2A_{0}^{tot}} + \sum_{j \in 1...N_{t}} \frac{k_{d}(A_{j} - A_{0})^{2}}{2A_{0}},$$

$$U_{v} = \frac{k_{v}(V - V_{0}^{tot})^{2}}{2V_{0}^{tot}},$$
(4)

where k_a , k_d and k_v are the global area, local area and volume constraint coefficients, respectively. The terms A and V are the total area and volume of RBC, while A_0^{tot} and V_0^{tot} are the specified total area and volume, respectively. More details on the RBC model can be found in Refs. [18, 19].

Other blood cells, such as WBCs and platelets, can be also represented by the membrane model described above for RBCs [20]. The main differences in comparison to RBCs are the shape and the membrane rigidity, which is normally set to be larger for WBCs and platelets than that for RBCs. More sophisticated models, which explicitly take into account the elasticity of the bulk cytoskeleton, are also available [21].



Fig. 4: Schematic of an attractive polymer with N = 30 beads in shear flow. Different bead colors denote non-activated (blue) or activated (green) polymer beads for adhesion. The colors are assigned using the two criteria from Eqs. (7) and (8), which determine a conformation-dependent bead activation for adhesion. Adapted with permission from Ref. [24].

2.2 von Willebrand factor

VWF is modeled as a shear-activated polymer, which retains its compact shape at low fluid stresses, but is able to stretch at high enough shear rates [11]. The model is based on a bead-spring chain with self-avoiding attractive monomers (see Fig. 4) [22, 23], where beads are connected by springs represented by the finite extensible nonlinear elastic (FENE) potential,

$$U_{\text{FENE}}(r_{ij}) = -\frac{1}{2}k_{\text{s}}l_{\text{max}}^2 \ln\left(1 - \left(\frac{r_{ij}}{l_{\text{max}}}\right)^2\right),\tag{5}$$

where k_s is the spring stiffness, l_{max} is the maximum spring extension, and r_{ij} is the distance between neighboring beads *i* and *j*. Self-avoidance and attractive interactions between beads, which lead to a globular configuration, are modeled by the 12-6 Lennard-Jones (LJ) potential as

$$U_{\rm LJ}(r_{ij}) = 4\epsilon_{\rm LJ} \left[\left(\frac{d_{\rm pol}}{r_{ij}} \right)^{12} - \left(\frac{d_{\rm pol}}{r_{ij}} \right)^6 \right],\tag{6}$$

where ϵ_{LJ} is the depth of the potential well, which controls the attraction strength, and d_{pol} is the bead diameter.

In a globular configuration, the polymer remains non-adhesive, while in a stretched state under flow, adhesive interactions with a surface become possible. This activation mechanism is modeled by tracking the degree of local stretching of the polymer in flow [24], as illustrated in Fig. 4. Two geometrical criteria are considered for the bead activation. The first condition corresponds to the angle between two adjacent bonds linking the bead *i* to its neighboring monomers (i - 1 and i + 1), which is directly related to the degree of local stretching of the polymer. Mathematically, it can be expressed as

$$\theta_{i-1,i,i+1} \ge \theta_{\text{thres}} \quad 2 < i < N-1, \tag{7}$$

where θ_{thres} is a threshold angle. Note that this condition is always assumed to be true for the first and the last bead in the polymer. The second criterion monitors the proximity of non-direct neighboring beads within a polymer and is expressed as

$$r_{ij} \ge R_{\text{thres}} \quad j \ne i, i \pm 1,$$
(8)

where R_{thres} is a threshold radius. This condition prohibits activation of a polymer bead for adhesion within a globule, even if the condition in Eq. (7) is satisfied. Thus, an inactive bead becomes activated when the two conditions above are satisfied. Similarly, an active bead will be deactivated when one or both criteria are not met.

3 Results

3.1 Fahraeus & Fahraeus-Lindqvist effecs

The basis of our current understanding of microvascular blood flow is mainly formed by two observations, the Fahraeus [25] and the Fahraeus-Lindqvist [26] effects, which have first been found in *in vitro* experiments [25, 26, 27, 28] of blood flow in glass tubes. The Fahraeus effect describes a reduced tube hematocrit in comparison with the discharge hematocrit at the tube exit. Here, the discharge hematocrit is defined as a fraction of RBCs exiting a tube per unit time, while the tube hematocrit corresponds to the volume fraction of RBCs in that tube. The Fahraeus effect arises from the property of RBCs to populate mainly the center of a tube in flow, where they move faster on average than the whole cell suspension (or blood) leading to an increased outflow rate of RBCs measured experimentally. The concentration of RBCs in tube center is a consequence of the migration of RBCs away from vessel walls toward the tube center [29, 30]. This migration is governed by hydrodynamic interactions of RBCs with the walls, which is often referred to as *lift force* in the literature [31, 32]. The Fahraeus effect [25] for single vessels directly implies reduced hematocrit values in microcirculation in comparison with a systemic hematocrit, which has been found experimentally [33, 34]. Following the conservation of blood volume within a closed circulatory system, we can conclude that the hematocrit in microvasculature should be lower than that in large vessels, since RBCs in microcirculation flow faster than the average blood flow.

The second well-known property of flowing blood is the Fahraeus-Lindqvist effect [26], which describes a decrease in the apparent blood viscosity with decreasing tube diameter found in experiments of blood flow in glass tubes [27, 28]. The apparent viscosity is calculated as

$$\eta_{app} = \frac{\pi \Delta p D^4}{128QL} = \frac{\Delta p D^2}{32\bar{v}L},\tag{9}$$

where D is the tube diameter, Q is the flow rate, and $\Delta p/L$ is the pressure drop in a tube of length L. For higher hematocrits, the apparent viscosity increases, since higher cell crowding leads to a larger flow resistance. For convenience, we define the relative apparent viscosity as

$$\eta_{rel} = \frac{\eta_{app}}{\eta_o},\tag{10}$$

where η_o is the plasma viscosity. Figure 5 presents simulation results [35] in comparison with the empirical fit to experiments [28] for tube diameters from 10 µm to 40 µm and hematocrit values in the range 0.15 to 0.45. The empirical fit [28] shows that the effective blood viscosity has its minimum in tubes with diameters of about 7 – 8 µm, while for smaller and larger tube diameters the apparent viscosity is markedly higher. This effect is also directly associated with the property of RBCs to migrate toward the tube center leading to the formation of two phases [29, 30]: a RBC-rich core and a RBC-free layer (RBC-FL) next to the tube walls. The RBC-FL leads to a decrease in the apparent viscosity when its thickness is significant in comparison with the tube diameter, i.e. for not too wide tubes.



Fig. 5: Simulated relative apparent viscosity of blood [35] in comparison with experimental data [28] for different hematocrit values and tube diameters. Reproduced with permission from Ref. [16].

Another effect which may introduce considerable variations of the hematocrit in microcirculation is related to the splitting of RBC volume fraction at bifurcations and branching points [36, 37]. For instance, many side branching vessels may receive only a small fraction of RBCs due to the plasma skimming effect, i.e. such vessels would be fed primarily by blood plasma, since RBCs are located in the center of the feeding vessel. Also, heterogeneous flow rates within the microvasculature result in an elevated hematocrit at daughter vessels with higher flow rates [36, 37]. Thus, the effective discharge hematocrit for a network defined as the flow-weighted mean of the discharge hematocrits in the individual vessels is higher than the simple mean of the individual discharge hematocrits for this vessel network. This property of microcirculatory flow leads to a further reduction of the mean microvascular tube hematocrit and has been referred to as the "network" Fahraeus effect in the literature [36, 37].

3.2 Margination

Another interesting phenomenon, which takes place in blood flow in the microcirculation is margination of WBCs, platelets, and drug carriers, as illustrated in Fig. 6. Margination is referred to the process of migration of these cells or suspended particles toward vessel walls. In case of WBCs and platelets, the margination process appears to be important for their function, since these cells may need to adhere to vessel walls. Current understanding of the margination process rests on the two mechanisms: (i) particle/cell migration due to hydrodynamic interactions (i.e., lift force) with the walls [31, 32], and (ii) shear-induced diffusion due to collisions/interactions between cells in flow [38, 39, 40]. The former mechanism may lead to the competition between the lift forces on different cells, generally resulting in the effect that more deformable cells (e.g., RBCs) migrate faster away from the wall in comparison to WBCs, for



Fig. 6: A simulation snapshot of RBCs and a marginated WBC at a hematocrit of 0.3. The vessel diameter is 20 µm and the flow is from the left to the right. Reproduced with permissions from Ref. [20].

example. The latter mechanism may lead to gradients in the apparent particle/cell diffusivities, which can be interpreted as effective driving forces considering the drift-diffusion formalism [38, 39].

Simulations of WBC margination in blood flow [41, 42, 20] have shown that the most efficient margination of WBCs is achieved at an intermediate range of hematocrits ($H_t = 0.2 - 0.4$) consistent with microcirculatory values [4, 6] and low enough flow rates characteristic for venular blood flow. Margination of platelets has been investigated in a number of numerical studies [43, 38] showing that a certain hematocrit value is required for platelet margination. The margination dynamics of platelets is found to be diffusional [38] in large enough vessels and might be slow in comparison to small vessels such as capillaries. Margination of micro- and nano-particles in blood flow has been also investigated in simulations [44, 45], indicating that microparticles possess much better margination properties in comparison to their nano-scale counterparts.

3.3 VWF adhesion

VWF molecules become adhesive upon stretching [10]. This process is represented by the shear-responsive polymer model described in section 2.2, and activated adhesive beads of the polymer can form bonds with ligands distributed on a wall or platelet surface. The adhesive interactions between two available sites can be formulated as 'chemical' reactions between the unbound and bound states with the rates k_{on} and k_{off} called association and dissociation rates, respectively. The reaction rates determine the frequency of state change, which is modeled by the transition probabilities P_{on} and P_{off} [46] as

$$\frac{\mathrm{d}P_{\mathrm{on}}}{\mathrm{d}t} = -k_{\mathrm{on}}P_{\mathrm{on}}, \text{ if } r \le r_{\mathrm{cut}}^{\mathrm{on}}, \tag{11}$$

while $P_{\text{on}} = 0$, if $r > r_{\text{cut}}^{\text{on}}$, and

$$\frac{\mathrm{d}P_{\mathrm{off}}}{\mathrm{d}t} = -k_{\mathrm{off}}P_{\mathrm{off}},\tag{12}$$

where $r_{\text{cut}}^{\text{on}}$ is the cutoff range for bond association. The association rate k_{on} is often assumed to be constant, while k_{off} might be dependent on the force applied to a bond.

It is intuitive that bonds should rupture under strong enough forces and therefore, the slip-bond model assumes that the detachment probability of a bond increases as the applied force on this





Fig. 7: Schematic illustration of VWF adhesion to a substrate: (a) formation of initial adhesion to a surface by one of its ends, (b) stretching of an initially tethered VWF, (c) partially adhered chain after stretching, and (d) fully adhered VWF.

bond is elevated. An alternative characteristic to the detachment probability is the lifetime of a bond, which for a slip bond decreases with increasing applied force. However, the lifetime of certain ligand-receptor interactions may increase as the force is elevated [47, 48] and such bonds are referred to as catch bonds [47]. Clearly, any physical bond will eventually rupture when the applied force becomes large enough. Therefore, the lifetime of an initially catch bond has to start decreasing when a certain value of the applied force is exceeded. As a result, the catch bond behavior should be rather considered as a dual catch-slip behavior. For example, adhesion of VWF to platelet GPIb α receptor exhibits the catch-slip behavior [49, 50].

The stretching of a self-attractive polymer (e.g., VWF) in shear flow starts by pulling a tether from the globule or equivalently one of the polymer ends [23]. Thus, the initial adhesion of VWF is very likely to proceed by binding one of its ends first. Figure 7 illustrates a typical adhesion process of a VWF. At first, one end of a partially stretched polymer adheres to the wall (Fig. 7(a)) and the polymer becomes tethered. Next, the tethered polymer is subjected to significant shear forces exerted by the fluid flow, leading to a further unfolding of the globule (Fig. 7(b)). Note that stretching of an immobilized polymer occurs at much lower shear rates than those required for the stretching of a freely-suspended globule in shear flow [51]. Thus, after initial tethering, unfolding of a globule proceeds rather rapidly (Fig. 7(c)). Following a rapid unfolding of the tethered globule, the VWF exposes new adhesive sites, which can bind to a surface as the shear flow pushes it closer to the wall. Finally, the VWF chain fully adheres (Fig. 7(d)).

It is important to emphasize that the formation of an initial adhesion (Fig. 7(a)) plays a deciding role for the overall adhesion to a surface. Clearly, this initial interaction has to be strong enough, and the formed bonds need to possess a long enough lifetime in order to sustain large forces exerted by the fluid flow and to allow enough time for polymer stretching and further adhesion to occur. Therefore, this process has to be facilitated by a relatively fast association rate and need to involve preferentially a catch-like bond behavior. In fact, recent simulations [52] suggest



Fig. 8: An aggregate formed by spherical particles (mimicking platelets) and VWF chains. Such aggregates form at high enough shear rates and are reversible when the flow is stopped.

that the adhesion of VWF to a surface at high shear rates is controlled by long-lived (catchlike) bonds. Similarly, adsorption of homopolymeric globules to a surface has been found to be enhanced at high shear rates by catch-bond interactions [53].

A very interesting behavior of VWF is the dissociation of VWF-platelet aggregates, which form at high shear rates, but disaggregate when the flow is stopped [49, 54]. Thus, adhesion of VWF is reversible under flow cessation. An adhered polymer with catch-slip surface interactions remains stable in shear flow unless the shear stresses become very large or equivalently when formed bonds stretch significantly such that they reach their slip part, where the lifetime strongly drops. At an intermediate shear stress, most of the bonds actually have very long lifetime due to the catch part of bond interactions. When the shear stress is removed by stopping the flow, the adhesive bonds return to their nearly unstressed state, and their lifetime significantly drops because of their catch-bond characteristic. Hence, the polymer starts losing its bonds to the substrate and forms a globule, mainly due to its internal attractive interactions. Eventually, the globule completely dissociates from the surface, since all polymer beads become inactive.

3.4 VWF-platelet aggregates

Similarly to the VWF adhesion at a wall, VWF chains bind to flowing platelets, resulting in the formation of VWF-platelet aggregates, as illustrated in Fig. 8. Aggregate formation is supported by the rotational motion in shear flow, which facilitates the wrapping of stretched and adhesive VWFs around platelets. Here, also high enough stresses are required for the activation of VWF. Furthermore, these aggregates are reversible and they dissociate when the flow strength is reduced or the flow is stopped. The reversibility of aggregates is directly associated with the reversible adhesion of VWF discussed in section 3.3.

3.5 Blood flow in microcirculation

Blood-flow modeling on the level of single cells in large microvascular networks still remains very limited. In order to overcome these problems, current models of blood flow in microcirculation [55, 56] implement a network approach of a coupled system of Poiseuille-flow equations in individual segments of the microvasculature. Thus, following the Poiseuille law the conduc-

tance G of each segment can be described similar to Eq. (9) as

$$G = \frac{Q}{\Delta p} = \frac{\pi D^4}{128\eta L},\tag{13}$$

where η is the fluid's dynamic viscosity. In a network, mass conservation (the sum of inflows must be equal to the sum of outflows) is employed to obtain the equations

$$\sum_{i,j=1}^{N} Q_{ij} = \sum_{i,j=1}^{N} G_{ij}(p_i - p_j) = 0,$$
(14)

where p_i is the pressure at different segment nodes (e.g., bifurcations, vessel junctions), Q_{ij} is the flow rate at a segment (i, j), and G_{ij} is its hydraulic conductance. Such a system of equations also requires consistent boundary conditions at all inflows and outflows of the selected network.

The hydraulic conductance G_{ij} is defined through the parameters D_{ij} , L_{ij} , and η_{ij} , where the diameters and lengths of segments simply come from a geometrical structure of the network of interest. As a first approximation, a constant viscosity in all segments could be assumed; however, in practice, this assumption is far from realistic due to significant hematocrit heterogeneities within a vessel network. Therefore, a better assumption is to make η_{ij} hematocrit dependent with the values following the empirical fits to available experimental data [28]. The appearance of hematocrit as a new variable in this system of equations implies that we have to determine physical separation of blood cells at branching points or equivalently to define bifurcation laws. Currently, the state of the art for the bifurcation laws is based on a phenomenological model, which is an average fit to available experimental data [6].

4 Conclusions

Theoretical models of blood flow in microvasculature are very powerful tools, because they are able to represent blood flow in a large network of vessels and capture various microcirculatory processes. However, these models depend entirely on the input assumptions which include blood-viscosity dependence, bifurcation laws, margination effects, 'non-hydrodynamic' processes (e.g., the adhesion of WBCs, platelets, and VWF), etc. This input should ideally come from experimental measurements to make the blood-flow simulations realistic. Even though some experimental data about the mean blood-flow characteristics are readily available, a systematic and detailed description of blood flow properties in microvasculature requires new experimental measurements in real and/or artificial microvascular networks, which can provide the basis for model validation. Furthermore, the combination of a detailed model on the level of single cells and a simplified model capable of capturing blood flow in a large microvascular network appears to be very promising. The integration of these different approaches may make it possible to build advanced blood-flow models capable of describing realistically the flow in microcirculation.
Appendices

A Smoothed dissipative particle dynamics

The original smoothed dissipative particle dynamics (SDPD) method [14] is derived through a Lagrangian discretization of the Navier-Stokes (NS) equations similar to the smoothed particle hydrodynamics method [57], while thermal fluctuations in SDPD are introduced similarly to the dissipative particle dynamics method [58, 59]. Each SDPD particle represents a small volume of fluid, instead of individual atoms or molecules, and corresponds to a thermodynamically consistent and well-defined physical volume [60]. The original SDPD method [14] does not conserve angular momentum, which might be crucial for some problems. A more general version of the SDPD method with angular momentum conservation [15] has been proposed recently, where every particle possesses a spin variable ω and a moment of inertia *I* in addition to its mass *m*, coordinate **r**, and linear velocity **v**. The new SDPD method [15] is obtained by a Lagrangian discretization of the continuity equation and the NS equation with spin [61].

Particle interactions in the SDPD method with angular momentum conservation [15] are represented by three deterministic forces (conservative (C), translational dissipative (D), and rotational dissipative (R)) and a random force (\sim) given by

$$\mathbf{F}_{ij}^{C} = \left(\frac{p_{i}}{\rho_{i}^{2}} + \frac{p_{j}}{\rho_{j}^{2}}\right) F_{ij}\mathbf{r}_{ij}, \\
\mathbf{F}_{ij}^{D} = -\gamma_{ij} \left(\mathbf{v}_{ij} + \mathbf{n}_{ij}(\mathbf{n}_{ij} \cdot \mathbf{v}_{ij})\right), \\
\mathbf{F}_{ij}^{R} = -\gamma_{ij}\frac{\mathbf{r}_{ij}}{2} \times (\boldsymbol{\omega}_{i} + \boldsymbol{\omega}_{j}), \\
\tilde{\mathbf{F}}_{ij} = \sigma_{ij} \left(d\overline{\boldsymbol{\mathcal{W}}}_{ij}^{S} + \frac{1}{3}tr[d\boldsymbol{\mathcal{W}}_{ij}]\mathbb{1}\right) \cdot \frac{\mathbf{n}_{ij}}{\delta t},$$
(15)

where p is the particle pressure, which follows a selected equation of state, ρ is the particle density obtained as $\rho_i = \sum_{j=1}^{N} m_j W_{ij}$ for particle i with $W_{ij} = W(\mathbf{r}_{ij}) = W(\mathbf{r}_i - \mathbf{r}_j, h)$ being an interpolation kernel, h the smoothing radius, and $\nabla_i W_{ij} = -\mathbf{r}_{ij} F_{ij}$. Furthermore, $\mathbf{n}_{ij} = \mathbf{r}_{ij}/|\mathbf{r}_{ij}|$, $\gamma_{ij} = 20\eta F_{ij}/(7\rho_i\rho_j)$ is the dissipative coefficient, $\sigma_{ij} = 2\sqrt{\gamma_{ij}k_BT}$ is the random coefficient with T being the temperature and k_B the Boltzmann constant. \mathcal{W}_{ij} is a matrix of independent Wiener increments, $tr[d\mathcal{W}_{ij}]$ is its trace, $d\overline{\mathcal{W}}_{ij}^S$ is the traceless symmetric part, 1 is the unity matrix, and δt is the integration time step.

The forces in Eq. (15) lead to torques \mathbf{L}_{ij} acting on a particle *i* from surrounding particles *j* given by $\mathbf{L}_{ij} = \frac{1}{2}\mathbf{r}_{ij} \times \mathbf{F}_{ij}$. The equation of motion for a particle *i* is described by the Newton's second law as

$$\dot{\mathbf{r}}_i = \mathbf{v}_i, \qquad \dot{\mathbf{v}}_i = \sum_j \frac{1}{m_j} \mathbf{F}_{ij}, \qquad \dot{\boldsymbol{\omega}}_i = \sum_j \frac{1}{I_j} \mathbf{L}_{ij}.$$
 (16)

Time evolution of particle positions, translational and angular velocities is integrated using the velocity-Verlet algorithm [62].

References

 Y. C. Fung, *Biomechanics: Mechanical properties of living tissues* (Springer-Verlag, New York, 1993), 2nd ed.

- [2] G. Pocock, C. D. Richards, and D. A. Richards, *Human physiology* (Oxford University Press, 2006).
- [3] Y. Cao, T. Wu, Z. Yuan, D. Li, S. Ni, J. Hu, and H. Lu, Sci. Rep. 5, 12643 (2015).
- [4] A. S. Popel and P. C. Johnson, Annu. Rev. Fluid Mech. 37, 43 (2005).
- [5] T. W. Secomb and A. R. Pries, J. Physiol. 589, 1047 (2011).
- [6] A. R. Pries, T. W. Secomb, and P. Gaehtgens, Cardiovasc. Res. 32, 654 (1996).
- [7] A. J. Reininger, Haemophilia 14, 11 (2008).
- [8] B. P. Nuyttens, T. Thijs, H. Deckmyn, and K. Broos, Thromb. Res. 127, S26 (2011).
- [9] W. E. Fowler, L. J. Fretto, K. K. Hamilton, H. P. Erickson, and P. A. McKee, J. Clin. Invest. 76, 1491 (1985).
- [10] T. A. Springer, Blood 124, 1412 (2014).
- [11] S. W. Schneider, S. Nuschele, A. Wixforth, C. Gorzelanny, A. Alexander-Katz, R. R. Netz, and M. F. Schneider, Proc. Natl. Acad. Sci. USA 104, 7899 (2007).
- [12] H. Ulrichts, M. Udvardy, P. J. Lenting, I. Pareyn, N. Vandeputte, K. Vanhoorelbeke, and H. Deckmyn, J. Biol. Chem. 281, 4699 (2006).
- [13] R. Schneppenheim, U. Budde, F. Oyen, D. Angerhaus, V. Aumann, E. Drewke, W. Hassenpflug, J. Häberle, K. Kentouche, E. Kohne, K. Kurnik, D. Mueller-Wiefel, *et al.*, Blood **101**, 1845 (2003).
- [14] P. Español and M. Revenga, Phys. Rev. E 67, 026705 (2003).
- [15] K. Müller, D. A. Fedosov, and G. Gompper, J. Comp. Phys. 281, 301 (2015).
- [16] D. A. Fedosov, H. Noguchi, and G. Gompper, Biomech. Model. Mechanobiol. 13, 239 (2014).
- [17] H. Noguchi and G. Gompper, Proc. Natl. Acad. Sci. USA 102, 14159 (2005).
- [18] D. A. Fedosov, B. Caswell, and G. E. Karniadakis, Biophys. J. 98, 2215 (2010).
- [19] D. A. Fedosov, B. Caswell, and G. E. Karniadakis, Comput. Meth. Appl. Mech. Eng. 199, 1937 (2010).
- [20] D. A. Fedosov and G. Gompper, Soft Matter 10, 2961 (2014).
- [21] M. L. Rodriguez, P. J. McGarry, and N. J. Sniadecki, Appl. Mech. Rev. 65, 060801 (2013).
- [22] A. Alexander-Katz, M. F. Schneider, S. W. Schneider, A. Wixforth, and R. R. Netz, Phys. Rev. Lett. 97, 138101 (2006).
- [23] A. Alexander-Katz and R. R. Netz, Macromolecules 41, 3363 (2008).
- [24] B. Huisman, M. Hoore, G. Gompper, and D. A. Fedosov, Med. Eng. Phys. 48, 14 (2017).
- [25] R. Fåhraeus, Physiol. Rev. 9, 241 (1929).
- [26] R. Fåhraeus and T. Lindqvist, Am. J. Phys. 96, 562 (1931).
- [27] W. Reinke, P. Gaehtgens, and P. C. Johnson, Am. J. Physiol. 253, H540 (1987).
- [28] A. R. Pries, D. Neuhaus, and P. Gaehtgens, Am. J. Physiol. 263, H1770 (1992).
- [29] H. L. Goldsmith, G. R. Cokelet, and P. Gaehtgens, Am. J. Physiol. 257, H1005 (1989).
- [30] G. R. Cokelet and H. L. Goldsmith, Circ. Res. 68, 1 (1991).
- [31] I. Cantat and C. Misbah, Phys. Rev. Lett. 83, 880 (1999).
- [32] M. Abkarian, C. Lartigue, and A. Viallat, Phys. Rev. Lett. 88, 068103 (2002).
- [33] B. Klitzman and B. R. Duling, Am. J. Physiol. 237, H481 (1979).
- [34] G. W. Schmid-Schönbein, R. Skalak, S. Usami, and S. Chien, Microvasc. Res. 19, 18 (1980).
- [35] D. A. Fedosov, B. Caswell, A. S. Popel, and G. E. Karniadakis, Microcirculation 17, 615 (2010).
- [36] A. R. Pries, K. Ley, and P. Gaehtgens, Am. J. Physiol. 251, H1324 (1986).
- [37] A. R. Pries, K. Ley, M. Claassen, and P. Gaehtgens, Microvasc. Res. 38, 81 (1989).

- [38] H. Zhao, E. S. G. Shaqfeh, and V. Narsimhan, Phys. Fluids 24, 011902 (2012).
- [39] A. Kumar, R. G. Henriquez-Rivera, and M. D. Graham, J. Fluid Mech. 738, 423 (2014).
- [40] K. Vahidkhah, S. L. Diamond, and P. Bagchi, Biophys. J. 106, 2529 (2014).
- [41] J. B. Freund, Phys. Fluids 19, 023301 (2007).
- [42] D. A. Fedosov, J. Fornleitner, and G. Gompper, Phys. Rev. Lett. 108, 028104 (2012).
- [43] L. Crowl and A. L. Fogelson, J. Fluid Mech. 676, 348 (2011).
- [44] K. Müller, D. A. Fedosov, and G. Gompper, Sci. Rep. 4, 4871 (2014).
- [45] S. Fitzgibbon, A. P. Spann, Q. M. Qi, and E. S. G. Shaqfeh, Biophys. J. 108, 2601 (2015).
- [46] G. I. Bell, Science **200**, 618 (1978).
- [47] M. Dembo, D. C. Torney, K. Saxman, and D. Hammer, Proc. R. Soc. Lond. B 234, 55 (1988).
- [48] F. Kong, A. J. García, A. P. Mould, M. J. Humphries, and C. Zhu, J. Cell. Biol. 185, 1275 (2009).
- [49] T. A. Doggett, G. Girdhar, A. Lawshé, D. W. Schmidtke, I. J. Laurenzi, S. L. Diamond, and T. G. Diacovo, Biophys. J. 83, 194 (2002).
- [50] J. Kim, C.-Z. Zhang, X. Zhang, and T. A. Springer, Nature 466, 992 (2010).
- [51] C. E. Sing and A. Alexander-Katz, Macromolecules 44, 9020 (2011).
- [52] C. E. Sing, J. G. Selvidge, and A. Alexander-Katz, Biophys. J. 105, 1475 (2013).
- [53] M. Radtke and R. R. Netz, Eur. Phys. J. E 38, 69 (2015).
- [54] H. Chen, M. A. Fallah, V. Huck, J. I. Angerer, A. J. Reininger, S. W. Schneider, M. F. Schneider, and A. Alexander-Katz, Nat. Commun. 4, 1333 (2013).
- [55] A. R. Pries, T. W. Secomb, T. Gessner, M. B. Sperandio, J. F. Gross, and P. Gaehtgens, Circ. Res. 75, 904 (1994).
- [56] A. R. Pries, T. W. Secomb, P. Gaehtgens, and J. F. Gross, Circ. Res. 67, 826 (1990).
- [57] J. J. Monaghan, Rep. Prog. Phys. 68, 1703 (2005).
- [58] P. J. Hoogerbrugge and J. M. V. A. Koelman, Europhys. Lett. 19, 155 (1992).
- [59] P. Español and P. Warren, Europhys. Lett. 30, 191 (1995).
- [60] A. Vázquez-Quesada, M. Ellero, and P. Español, J. Chem. Phys. 130, 034901 (2009).
- [61] D. W. Condiff and J. S. Dahler, Phys. Fluids 7, 842 (1964).
- [62] M. P. Allen and D. J. Tildesley, *Computer simulation of liquids* (Clarendon Press, New York, 1991).

E 3 Principles of Active Matter

Roland G. Winkler Theoretical Soft Matter and Biophysics Institute for Advanced Simulation Forschungszentrum Jülich GmbH

Contents

1	Intr	oduction	2	
2	Active Brownian particle (ABP)			
	2.1	Equations of motion	4	
	2.2	Active Ornstein-Uhlenbeck particle (AOUP)	5	
	2.3	Mean square displacement	5	
	2.4	Surface accumulation	6	
	2.5	Active pressure	7	
	2.6	Motility-induced phase separation (MIPS)	8	
	2.7	Collective motion—the Vicsek model	9	
3	Life at low Reynolds numbers			
	3.1	Hydrodynamics	10	
	3.2	Solution of Stokes equation	10	
	3.3	Force dipole	11	
	3.4	Squirmer—a model hydrodynamic microswimmer	12	
	3.5	Squirmer cooperative locomotion	13	
	3.6	Squirmer cluster formation	14	
4	Con	clusions	15	
A	Fokker-Planck equation of ABP			
B	Fok	ker-Planck equation of AOUP	17	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Active matter, whose constituents (agents) consume internal energy or extract energy from the environment and are thus far from thermal equilibrium, comprises a vast range of systems ranging from the nano- and microscale up to macroscopic length scales (cf. Fig. 1) [1–5]. On the nanoscale, proteins and other macromolecules in the interior of a cell undergo cyclic conformational changes and stir the surrounding fluid [6,7]. Since the appearing fluctuating flows are non-thermal, work can be performed by extraction of energy leading to, e.g., an enhanced diffusive motion [6,8] and chemotaxis [9]. In motility assays, biological semiflexible polar filaments, such as actin and microtubules, are propelled on carpets of motor proteins anchored on a substrate, which results in a directed motion [10-17]. Propulsion of such biological filaments in the cell cytoskeleton due to tread-milling and dimeric or tetrameric motor proteins is ubiquitous. Mixtures of active and passive components are a characteristics of eukaryotic cells with the active cytoskeleton on the one hand and an embedded large variety of passive colloidal and polymeric objects on the other hand. Here, an enhanced random motion of tracer particles has been observed [18]. Moreover, an influence of the active microtubule [19] or actin-filament [20] dynamics on the motion of chromosomal loci [21,22] or that of chromatin has been found [23]. On larger scales, there is plethora of biological microswimmers such as spermatozoa, bacteria, protozoa, and algae [2]. They use flagella—whip-like structures protruding from their bodies for their propulsion (cf. Fig. 1). Swimming of uni- and multicellular organisms is essential for their search for food (chemotaxis), the reaction to light (phototaxis), the orientation in the gravitation field (gravitaxis), or for reproduction [2,24]. A paradigmatic example is *Escherichia* coli (E. coli), which propels itself by helical filaments (cf. Fig. 1). A flagellum is rotated by a motor complex consisting of several proteins, and is anchored in the bacterial cell wall [25,26]. The (counterclock-wise) rotating flagella self-organize in helical bundle(s), which push the cell forward [24, 27, 28]. E. coli and other bacteria swim in a "run-and- tumble" motion, where they change the rotation direction of some or all flagella, which results in a deterioration of the bundle and a reorientation (tumble) [25, 29, 30]. By returning to the original rotation direction, the bundle is reestablished and the bacterium swims again (run). This leads to a diffusive motion determined by the activity of the cell on long time scales.

Flagellated microorganisms not only swim as individuals, but exhibit collective behavior at a surface or in a thin liquid film in form of swarming [2]. Here, bacteria cooperativity reaches a new level as they exhibit highly organized movements with remarkable large-scale patterns such as networks, complex vortices, swarms, or turbulence [31, 32]. Remarkably, similar collective phenomena are observed by active systems on the macroscale such as schools of fish, flocks of birds (cf. Fig. 1), mammalian herds, or crowds of humans [33].

Aside from biological active matter, the design of artificial nano- and microswimmers is highly desirable to perform a multitude of tasks in technical and medical applications. Consequently, various design strategies are explored and different propulsion mechanisms have been proposed [2, 3, 34]. In particular, experimental studies on self-propelled Janus particles and computer simulations reveal motility induced cluster formation and phase separation (MIPS) [2, 3, 35–42] (cf. Fig. 1).

The physics ruling activity and swimming on the micrometer scale is very different from that applying to the macroworld. Swimming at the micrometer scale is swimming at low Reynolds numbers, where viscous damping by far dominates over inertia [43]. Hence, swimming concepts of the high-Reynolds number macroworld are ineffective on small scales. In the evolutionary process, microorganisms acquired propulsion strategies, which successfully overcome



Fig. 1: (Top left) An active nematic liquid of filamentous microtubules driven by molecular motors exhibits collectively drive mesoscale turbulent-like dynamics [44, 45]. (Top right) Swarming E. coli bacteria [46]. The inset shows a Salmonella enterica with the cylindrical cell body and several flagella [47]. (Bottom left) Living crystals phoretically assembled from a homogeneous distribution of bimaterial colloids—a TPM polymer colloidal sphere (3-methacryloxypropyl trimethoxysilane) with protruding hematite cube (inset)—under illumination by blue light [38]. (Bottom right) Swarm of starlings [48].

and even exploit viscous drag. Hence, the rules governing swarming on the macroscale of flocks of birds or schools of fish are very different from those applying to the microscale of bacteria. Nevertheless, similar features appear on very distinct length scales.

A variety of models for active matter or agents has been developed to describe the aspect of interesting in sufficient detail. Examples are models for sperm [49] or *E. coli* bacteria [28] (see also references therein). From the theoretical side, rather generic models have been proposed, which lack details of real agents, but allow for a systematic study of active matter over a wide range of parameters. The prototype of a generic model is the active Brownian particle (ABP), a hard-sphere colloid, which actively moves in a prescribed direction, where the latter changes in a diffusive manner [2, 3]. This model has developed into the standard model to study the nonequilibrium statistical properties of active matter. A system of ABPs is a so called dry active system, because the ABPs are not embedded in a fluid, which is typically the case for biological and synthetic active matter. Fluid-mediated interactions (hydrodynamics) are captured by the so-called squirmer model for an active particle. Thereby, the squirmer is propelled by a prescribed flow velocity (slip velocity) on the colloid surface [50–54]. In the

hydrodynamic far field, the squirmer flow field is determined by the weakest decaying point multipoles, which are the source dipole and the force dipole flow field. Specifically, by the latter, the difference, e.g., between bacteria (pusher) and algae (puller) can be captured by a squirmer representation [52–54].

In this contribution, various specific aspects of active matter, compared to passive systems, will be discussed. First of all, active Brownian particles will be introduced and their properties be studied. Secondly, the low-Reynolds number hydrodynamics of microswimmers will be addressed. Specifically, the squirmer model is presented, and the difference of its collective behavior compared to ABPs is discussed.

2 Active Brownian particle (ABP)

The active Brownian particle captures essential aspects of a self-propelled object [2, 3, 35, 37, 41, 55–59]. It is typically represented as a repulsive spherical colloid (rigid body) propelled by a constant (external) force in the direction of its instantaneous orientation, which is changing in a diffusive manner. However, no hydrodynamic interactions are taken into account, an aspect to be kept in mind.

2.1 Equations of motion

The Langevin equations for the center-of-mass position r and the orientation e of an ABP are given by [10]

$$\dot{\boldsymbol{r}}(t) = \boldsymbol{v}(t) + \frac{1}{\gamma} \left(\boldsymbol{F}(t) + \boldsymbol{\Gamma}(t) \right), \tag{1}$$

$$\dot{\boldsymbol{e}}(t) = \boldsymbol{\xi}(t) \times \boldsymbol{e}(t), \tag{2}$$

where $v = v_0 e$, with e a unit vector, is the propulsion velocity, F a force exerted on the particle, Γ and ξ are Gaussian and Markovian processes (white noise) with zero odd moments and the second moments

$$\langle \Gamma_{\alpha}(t)\Gamma_{\alpha'}(t')\rangle = 2\gamma k_B T \delta_{\alpha\alpha'} \delta(t-t'), \tag{3}$$

$$\langle \xi_{\alpha}(t)\xi_{\alpha'}(t')\rangle = (d-1)D_R\delta_{\alpha\alpha'}\delta(t-t'). \tag{4}$$

Here, k_B is the Boltzmann constant, T the temperature, γ the translational friction coefficient, which is related to the translational diffusion coefficient D_T via $D_T = k_B T/\gamma$, D_R the rotational diffusion coefficient, d the spatial dimension, and $\alpha, \alpha' \in \{x, y, z\}$. For a particle in a viscous fluid in three dimensions (3D), $\gamma = 6\pi\eta R$, with η the viscosity and R the particle radius, hence, D_R and D_T are related according to $D_T/D_R R^2 = 4/3$. However, in general, D_R can be independent of D_T and be of nonthermal origin, e.g., tumbling of bacteria. Equations (1) and (2) describe the solid-body translation and rotation, respectively. Thereby, we neglect the inertia terms, consistent with the fact that propulsion and motility on the nano- and microscale is typically governed by low-Reynolds number hydrodynamics (cf. Sec. 3) [2, 24]. The rotational motion (Eq. (2)) is independent of the colloid translation. As a particular result, the correlation function

$$\langle \boldsymbol{v}(t) \cdot \boldsymbol{v}(t') \rangle = v_0^2 e^{-(d-1)D_R |t-t'|}$$
(5)

is obtained (d > 1) [60–62]. In the above description, we consider the velocity v as an intrinsic property of the ABP. Alternatively, we can consider Eq. (1) only, with v as an external stochastic process with the exponential correlation (colored noise) of Eq. (5) [2,41].

A general solution of Eqs. (1) and (2) is difficult to determine, even a stationary state solution, because of the violation of detailed balance [14] (cf. App. A). However, various aspects can be calculated directly, specifically for a zero force F = 0.

2.2 Active Ornstein-Uhlenbeck particle (AOUP)

The propulsion direction e, with Eq. (4), obeys the strict condition |e(t)| = 1. By lifting this condition and considering the equation of motion for the propulsion velocity

$$\dot{\boldsymbol{v}}(t) = -\gamma_R \boldsymbol{v}(t) + \boldsymbol{\eta}(t), \tag{6}$$

a model is obtained, which is analytically easier tractable [2, 57, 63–67]. Here, the Cartesian velocity components are independent. The damping factor γ_R is related to the rotational diffusion coefficient according to $\gamma_R = (d-1)D_R$, and η is a Gaussian and Markovian stochastic process with zero mean and the second moments

$$\langle \eta_{\alpha}(t)\eta_{\alpha'}(t')\rangle = \frac{2(d-1)}{d}v_0^2 D_R \delta_{\alpha\alpha'}\delta(t-t').$$
(7)

The Langevin equations (6) with noise amplitudes $(v_0^2 D_R)$, which are independent of the stochastic variables (v), describe a process denoted as Ornstein-Uhlenbeck process [68]. Hence, an active particle obeying Eqs. (3) and (6) is denoted as active-Ornstein-Uhlenbeck particle (AOUP) [66]. Most importantly, an analytical solution and a stationary-state distribution function is obtained for a linear force F (harmonic potential) (cf. App. B) [67].

2.3 Mean square displacement

Integration of Eq. (1) in the force-free case (F = 0) yields

$$\boldsymbol{r}(t) - \boldsymbol{r}(0) = \int_0^t \left(\boldsymbol{v}(t') + \frac{1}{\gamma} \boldsymbol{\Gamma}(t') \right) dt', \tag{8}$$

from which the mean square displacement $\Delta r(t)^2 = \langle (r(t) - r(0))^2 \rangle$, with

$$\Delta \boldsymbol{r}^2 = \frac{6k_BT}{\gamma} t + \int_0^t \int_0^t \left\langle \boldsymbol{v}(t') \cdot \boldsymbol{v}(t'') \right\rangle dt' dt'' = 2dD_T t + \frac{2v_0^2}{\gamma_R^2} \left(\gamma_R t + e^{-\gamma_R t} - 1 \right), \quad (9)$$

is obtained by insertion of the correlation function (5). Equation (9) turns, for small and large times, respectively, into

$$\Delta \boldsymbol{r}^2 = 2dD_T t + v_0^2 t^2, \quad \gamma_R t \ll 1; \qquad \Delta \boldsymbol{r}^2 = \left(2dD_T + \frac{2v_0^2}{\gamma_R}\right) t, \quad \gamma_R t \gg 1.$$
(10)

Hence, activity leads to a ballistic motion at short times, and a diffusive motion for $\gamma_R t \gg 1$ due to rotational diffusion of the propulsion direction, with the effective translational diffusion coefficient $D_{T,\text{eff}} = D_T + v_0^2/d(d-1)D_R$. The latter is much larger than D_T for $v_0 \gg \sqrt{D_T D_R}$.



Fig. 2: Surface accumulation of ABPs. (Left) Separation of the active velocity $\mathbf{v} = v_0 \mathbf{e}$ parallel (\mathbf{v}_{\parallel}) and perpendicular (\mathbf{v}_{\perp}) to the surface. (Right) Probability distribution of an ABP in 3D confined between two parallel flat walls. The walls are located at z = 0 and z = 200R. The Péclet numbers are Pe = 0, 5, 10, 20, and 40 (bottom to top).

Thus, activity leads to a drastically increased diffusive motion. The theoretical prediction of Eq. (9) has been confirmed experimentally, e.g., for phoretic Janus particles [55]. The activity of an ABP is usually characterized by the dimension less Péclet number

$$Pe = \frac{v_0}{2RD_R}.$$
(11)

Therefore, $D_{T,\text{eff}}/D_T$ is much larger than unity for $Pe \gg 1$. This corresponds to the regime, where activity dominates over thermal fluctuations.

2.4 Surface accumulation

Activity gives raise to various unusual and *a priori* unexpected effects. An example is the accumulation of ABPs even at purely repulsive walls [71]. The mechanism is illustrated in Fig. 2. The propulsion velocity (v) of an ABP located at the wall can be separated in a component parallel (v_{\parallel}) and normal (v_{\perp}) to the wall. The parallel component leads to a translational motion of the center-of-mass parallel to the wall, since there is no friction between wall and ABP. The component normal to the wall pushes the swimmer toward the wall as long as $v_{\perp} \cdot e_z < 0$ (e_z is the unit vector along the *z*-axis), and the APB stays at the wall. Only after the change of the propulsion direction and for $v_{\perp} \cdot e_z > 0$, the particles leaves the wall again. This is different from thermal motion, where the velocity v_{\perp} is inverted instantaneously upon a collision with a (hard) wall. The appearing force onto the surface is calculated in Sec. 2.5. Figure 2 shows density profiles from simulations at various Péclet numbers for ABPs confined between to parallel walls in 3D, where the ABP-wall interaction is described by the Lennard-Jones potential (wall at z = 0)

$$U_{LJ}^{w} = \begin{cases} \epsilon \left[\left(\frac{R}{z}\right)^{48} - \left(\frac{R}{z}\right)^{24} \right] &, r < r_{c} \\ 0 &, r > r_{c} \end{cases}$$
(12)

with $r_c = R\sqrt[6]{2}$. Hence, we use a soft, but steep potential. The different behavior compared to a thermal system is evident, as well as the density increase with increasing *Pe*. As a consequence, the distribution of active agents, e.g., in a convex container is dramatically affected



Fig. 3: Rectification of bacteria motion. (Top left) Bacterial driven micromotor. A nanofabricated asymmetric gear immersed in an active bath of motile E. coli cells rotates clockwise visualized by the yellow circle [69]. (Right) Steady-state distribution of colloidal particles in a bacterial bath exposed to an asymmetric square saw-tooth structure (scale bar 20 μ m). There is a preferred transport in the direction of f_0 (bottom left), which leads to an accumulation of colloids inside (blue) or outside (red) of the squares [70].

by the boundary shape in the limit, in which the container size is small compared to the active persistence length, $l_p = RPe$, the distance a particle travels before its orientation decorrelates. In particular, the particles are confined at the boundary and their steady-state distribution is proportional to the local curvature [72]. This effect can be exploited for rectification or trapping of active agents [73, 74]. Examples are provided in Fig. 3, where bacteria produce a spontaneous and unidirectional rotation of a nano-fabricated gear [69], or passive colloids are spatially organized by a suspension of swimming bacteria [70].

2.5 Active pressure

During the encounter with a wall (cf. Fig. 2), an ABP exerts a force on the wall due to propulsion [62]. The force can be estimated from Eq. (1). For simplicity, we consider a cubic volume, within ABPs are confined. The velocity of an ABP at a wall and in the direction of the wall normal \boldsymbol{n} (\boldsymbol{n} points outward of the volume) is zero, i.e., $\dot{\boldsymbol{r}} \cdot \boldsymbol{n} = 0$, which implies that $v_0 \boldsymbol{e} \cdot \boldsymbol{n} + \boldsymbol{F}^s \cdot \boldsymbol{n}/\gamma = 0$. Here, \boldsymbol{F}^s is the force between the ABP and the wall, which we assume to be short ranged. Since $F_n = -\boldsymbol{F}^s \cdot \boldsymbol{n}$, we find

$$F_n(z_{max}) = \gamma v_0 \boldsymbol{e} \cdot \boldsymbol{n} \tag{13}$$

for the particle of Fig. 2 (left), where z_{max} is approximately equal to the z-position of the maximum in Fig. 2 (right). Defining surface stress (pressure) as force per area, A, Eq. (13) yields a the mechanical stress $\gamma v_0 e_z/A$ per particle. Extending the concept to a three-dimensional volume and taking into account N ABPs in the volume V, the viral approach yields the more general expression for the pressure $3pV = -\sum_{i=1}^{N} \langle F_i^s \cdot r_i \rangle$ [62]. For a short-range surface interaction, the ABP *i* is at a wall and r_i can be taken out of the sum, which leads to above stress (pressure) in the special case of single wall rather than six.

In systems with periodic boundary conditions, there are no confining walls. Here, as for any confined system, pressure can be calculated from the virial formulation exploiting the actual



Fig. 4: Motility-induced phase separation of ABPs. (Left) Two-dimensional well-ordered cluster in contact with a gas of ABPs. The inset illustrates the blockage of three particles. (Right) Highly dynamic cluster in three dimensions in contact with a gas of ABPs. In inset shows the three-dimensional structure.

positions and velocities of the ABPs. This approach leads to the pressure [62]

$$3pV = 3Nk_BT + 3Nk_BT \frac{v_0^2}{6D_TD_R} + \sum_{i=1}^N \langle \boldsymbol{F}_i^s \cdot \boldsymbol{e}_i \rangle + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^{N'} \langle \boldsymbol{F}_{ij} \cdot (\boldsymbol{r}_i - \boldsymbol{r}_j) \rangle.$$
(14)

The first term on the right hand side is the ideal gas contribution due to the Brownian motion of the ABP center of mass. The second term is the swim pressure, which is proportional to the square of the propulsion velocity [75]. The third terms captures corrections due to interactions with the walls; the term vanishes for infinite system sizes. Finally, the last term accounts for interparticle interactions. Compared to passive systems, the second and third term appear additionally.

Active systems are out of equilibrium and equilibrium statistical mechanics and thermodynamics cannot strictly be applied in their description. This is, e.g., reflected by the fact that temperature is ill defined in active systems [67], hence, there is no equation of state, in general, relating properties such as density and temperature with other thermodynamic bulk properties [76]. As shown above, there is an equation of state for the pressure of spherical ABPs confined in an orthorhombic volume. However, as discussed in Ref. [76], no such equation exists for an active fluid in general, since small orientation-dependent interactions (whether wall-particle or particle-particle) immediately destroy the equation of state. Such interactions are typically present in every experimental system. Yet, the mechanical pressure is well defined and can be calculated.

2.6 Motility-induced phase separation (MIPS)

The mechanism responsible for accumulation of ABPs at walls (Sec. 2.4) leads also to novel bulk phenomena. In dilute systems, clusters emerge and a motility-induced phase separation [3], as illustrated in Fig. 1. When two or three particles collide, they block each other due to their persistent motion as sketched in Fig. 4 (left). The particular cluster would resolves after a time $t \sim 1/\gamma_R$ due to the rotational diffusion of the orientation *e*. Interactions and collisions with other particles are controlled by the density and velocity v_0 . Hence, if other particles collide before the original cluster resolved, the cluster grows. At higher concentrations, a phase



Fig. 5: (*Left*) Illustration of the Vicsek dynamics. The red central particle with original orientation (red) aligns with neighbours (green) inside the circle of radius R_0 (red) (top), and then moves along its heading direction (bottom). (Middle) Magnitude S = |S| of the order parameter as function of noise strength ξ_0 for various system sizes L. (Right) Snapshot in the ordered phase. Points represent the individual particles and the red arrow indicates the global direction of motion. Periodic boundary conditions are applied for the 2D system. The simulation parameters can be found in Ref. [77].

transition in a high-density and a low density phase can appear, as depicted in Fig. 4. Thereby, the structure of the high density phase depends on the spatial dimension. In 2D, this phase is typically well ordered in a hexagonal fashion [3, 35], whereas in 3D, the high-density phase is rather mobile [39].

2.7 Collective motion—the Vicsek model

The three-dimensional APB system of Sec. 2.6 (Fig. 4) exhibits collective motion without any alignment rule between neighbors. Collective motion is often studied by implementing an alignment rule for spherical particles, or emerges naturally due to the anisotropy of the active particles, e.g., rods [78]. A prototypical example for a system with an explicit alignment rule is the Vicsek model [14, 79, 80].

The translational motion of a force free (F = 0) particle within the Vicsek model is given by Eq. (1), however the thermal noise is neglected, i.e, $\Gamma = 0$. Particles are assumed to align their direction of motion with their local neighbors. Hence, the orientation e_i of particle *i* depends on the average direction of all particles (including particle *i*) in a spherical neighborhood of radius R_0 centered at r_i (cf. Fig. 5). In fact, alignment within a neighborhood is almost perfect, its only perturbed by thermal-type white noise. For simplicity, we focus on 2D systems, where the direction of motion is determined by the angle θ_i , with $e_i = (\cos \theta_i, \sin \theta_i)^T$. The angle itself changes in time according to

$$\theta_i = \operatorname{Arg}\left(\sum_j \Theta(R_0 - |\boldsymbol{r}_i(t) - \boldsymbol{r}_j(t)|)\boldsymbol{e}_j(t)\right) + \xi_0 \Delta \xi_i(t).$$
(15)

The argument Arg yields the angle between the effective vector $\sum_{j} \Theta(R_0 - |\mathbf{r}_i(t) - \mathbf{r}_j(t)|) \mathbf{e}_j(t)$ and an arbitrary direction. Here, $\Theta(x)$ is the Heaviside step function, $\Delta \xi_i$ white noise with the correlation function $\langle \Delta \xi_i(t) \Delta \xi_j(t) \rangle = 2\delta_{ij}$, and ξ_0 is the noise amplitude. Simulations show that the Vicsek model exhibits a transition from disordered to ordered collective motion. The degree of ordering depends on the parameter ξ_0 . A measure for the degree of alignment is the order parameter $S(t) = \sum_{i=1}^{N} e_i(t)/N$, where N is the number of active particles. Figure 5 (middle) displays the dependence of the order parameter on the noise strength for various system sizes. For large systems and below a certain ξ_0 , a first order phase transition appears from a disordered to an ordered collective dynamic state [77]. The emerging sharp bands are illustrated in Fig. 5 (right).

3 Life at low Reynolds numbers

3.1 Hydrodynamics

Typically, the dynamics of the incompressible (isothermal) fluid flow field surrounding a microswimmer is described by the Navier-Stokes equations

$$\rho\left(\frac{\partial}{\partial t}\boldsymbol{v} + (\boldsymbol{v}\cdot\nabla)\,\boldsymbol{v}\right) = -\nabla p + \eta\nabla^2\boldsymbol{v} + \boldsymbol{f}\,, \qquad \nabla\cdot\boldsymbol{v} = 0, \tag{16}$$

where v(r, t), p(r, t), and f(r, t) are the velocity, pressure, and volume-force-density fields, respectively. At small Reynolds numbers $Re = \rho u L/\eta \ll 1$, where ρ is the fluid mass density, u the characteristic velocity, L the size of the micoswimmer, and η the fluid viscosity, the inertia terms on the left-hand side of Eq. (16) can be neglected, and the equations reduce to the Stokes or creeping flow equations

$$\nabla p(\mathbf{r}) - \eta \nabla^2 \mathbf{v}(\mathbf{r}) = \mathbf{f}(\mathbf{r}), \qquad \nabla \cdot \mathbf{v} = 0.$$
(17)

For illustration, the Reynolds number in water of a swimmer of length $L = 10 \ \mu$ m, a velocity of $u = 50 \mu$ m/s, and the kinematic viscosity $\nu = \eta/\rho = 10^{-6}$ m²/s is $Re \approx 10^{-3}$. The Stokes equation (17) is linear and time independent. The consequences of this intrinsic symmetry under time reversal for microswimmers undergoing periodic shape changes was first expressed in Ref. [43] by Purcell, and is now know as "scallop theorem", which can be stated as: if the shape changes displayed by a swimmer are identical when viewed in reverse order, it will generate an oscillatory, but no directed motion [2, 24, 43, 81]. Thus, just by opening and closing its two shells, a mussel (scallop) cannot move forward at $Re \ll 1$. Microswimmers developed various strategies to beat the scallop theorem. Aside from many (elastic) degrees of freedom, they use specific propulsion mechanisms—bacteria such es *E. coli* are propelled by rotating helical flagella bundles, sperm use sinusoidal bending waves propagating from head to tail, and algae, e.g., *Chlamydomonas* uses an non-reciprocal stroke pattern.

3.2 Solution of Stokes equation

The linear Stokes equations (17) are easily solved analytically for an unbounded fluid. The respective velocity field is

$$\boldsymbol{v}(\boldsymbol{r}) = \int \mathbf{Q}(\boldsymbol{r} - \boldsymbol{r}') \boldsymbol{f}(\boldsymbol{r}') \, d^3 r' \,, \quad Q_{\alpha\alpha'}(\boldsymbol{r}) = \frac{1}{8\pi\eta r} \left[\delta_{\alpha\alpha'} + \frac{r_{\alpha}r_{\alpha'}}{r^2} \right], \tag{18}$$

where $\mathbf{Q}(\mathbf{r})$ is the well-know Oseen tensor, with the Cartesian components $Q_{\alpha\alpha'}$ ($\alpha, \alpha' \in \{x, y, z\}$) and $r = |\mathbf{r}|$ [82, 83]. The Oseen tensor, also denoted as Stokeslet, shows that hydrodynamic interactions are long ranged, with a 1/r decay like the Coulomb potential, and it is



Fig. 6: (*Left*) Flow lines of a hydrodynamic dipole oriented horizontally, Eq. (20) [2]. The separatrices between the inflow and outflow regions are shown by thick red lines. Flow field of a E. coli bacterium from (Middle) experiment [84] and (Right) simulations [28]. In simulations, a system with periodic boundary conditions is considered, which yields closed flow lines in contrast to the flow lines of the experimental bulk system. The logarithmic color scheme (right) indicates the magnitude of the flow speed scaled by the bacterial swimming velocity.

anisotropic due to the incompressibility of the fluid. The Oseen tensor is the Green's function of the Stokes equation (17), which is evident, when the point force $f(r) = f_0 \delta(r) e$ in the direction e(|e| = 1) is inserted. Then, Eq. (18) yields

$$\boldsymbol{v}(\boldsymbol{r}) = \frac{f_0}{8\pi\eta r} \left[\boldsymbol{e} + \frac{(\boldsymbol{r} \cdot \boldsymbol{e})\boldsymbol{r}}{r^2} \right].$$
(19)

The magnitude of the flow field is evidently twice larger in the force direction than perpendicular to it.

3.3 Force dipole

Most swimmers move autonomously, with no external force or torque applied, and hence the total interaction force/torque of the swimmer on the fluid, and *vice versa*, vanishes. In the simplest case, which actually applies to many microswimmers like bacteria, spermatozoa, or algae, the far-field hydrodynamics (at distances from the swimmer much larger than its size) can well be described by a force dipole [24,85]. This has been confirmed experimentally for *E. coli* [84,86] and in simulations [28]. The flow field of *Chlamydomonas* is well reproduced by three Stokeslets [86].

Mathematically, the flow field $v_d(r - r_0)$ of a hydrodynamic force dipole located at r_0 follows by a superposition of two Stokeslets (18) with opposite forces $f_0 = \pm f_0 e$ of equal magnitude at $r_0 \pm l/2$, where l = le. Taylor expansion to leading order in $|l|/|r - r_0|$ yields

$$\boldsymbol{v}_d(\boldsymbol{r}) = \frac{P}{8\pi\eta r^3} \left[-1 + 3\frac{(\boldsymbol{r}\cdot\boldsymbol{e})^2}{r^2} \right] \boldsymbol{r},\tag{20}$$

where $P = \pm f_0 l$ is the dipole strength. Note that the flow field of a force dipole decays as $1/r^2$ from the center of the dipole, faster than the force monopole (Stokeslet) (18). The flow lines of a hydrodynamic dipole oriented vertically (x-direction) are displayed in Fig. 6. There are two inflow and two outflow regions in the xy-projection, which are separated by the separatrices $y = \pm \sqrt{2}x$. In three dimensions, the outflow region is a cone.



Fig. 7: Flow streamlines of isolated squirmers in the laboratory reference frame for (Left) a pusher ($\beta < 0$), (Middle) a neutral squirmer ($\beta = 0$), and (Right) a puller ($\beta > 0$). The inset indicates the definition of the angle ϑ and the tangential vector \mathbf{e}_{ϑ} in the squirmer-fixed reference frame.

Two classes of dipole swimmers can be distinguished, A swimmer with its "motor" in the back, and a passive body dragging along the surrounding fluid in front, creates a "pusher" flow field (cf. Fig. 6 (left)). Similarly, a swimmer with its "motor" in front, and the passive body dragging along the fluid behind, develops a "puller" flow field. The flow fields of pushers and pullers look similarly, but with opposite flow directions. This has important consequences for the interactions between swimmers and of swimmers with walls.

The flow field of an *E. coli* bacterium obtained from experiment [84] and simulations [28] is presented in Fig. 6. In both cases, the far field is well described the force dipole field of Eq. (20) [28]. However, there is also a distinct near field determined by the shape of the bacterium.

3.4 Squirmer—a model hydrodynamic microswimmer

A prototype of a swimmer capturing hydrodynamics is the squirmer, which was introduced by Lighthill [50] and revised by Blake [51]. Originally, it was intended as a model for ciliated microswimmers such as Paramecia. Nowadays, it is considered as a generic model for a broad class of microswimmers, ranging from diffusiophoretic particles [3, 55, 87, 88] to biological cells, and has been applied to study collective effects in bulk [52, 53, 89–93], at surfaces [53, 94, 95], and in narrow slits [96, 97]. In its simplest form, a squirmer is represented as a spherical rigid colloid with a prescribed surface velocity [50, 51, 90, 97]. Restricting the surface (slip) velocity to be tangential, the spherical squirmer is typically characterized by two modes accounting for its swimming velocity (B_1) and its force dipole (B_2). Explicitly, the slip velocity is then given by (cf. Fig. 7) [2, 53, 90, 97]

$$\boldsymbol{v}_{sq} = (B_1 \sin \vartheta + B_2 \sin \vartheta \cos \vartheta) \boldsymbol{e}_{\vartheta} = B_1 (\sin \vartheta + \beta \sin \vartheta \cos \vartheta) \boldsymbol{e}_{\vartheta}. \tag{21}$$

The parameter $B_1 = 2v_0/3$ is related to the swimming velocity, v_0 , and $\beta = B_2/B_1$ accounts of the force dipole. Higher order term can easily be taken into account [2,53]. By the term with B_2 (or β), we can distinguish between pushers ($\beta < 0$), pullers ($\beta > 0$), and neutral squirmers ($\beta = 0$), corresponding, e.g., to *E. coli*, *Chlamydomonas*, or *Volvox*, respectively.

The far field of a squirmer is well described by the flow fields of a force dipole (FD), a source dipole (SD), and a source quadrupole (SQ)

$$\boldsymbol{v}(\boldsymbol{r}) = \kappa^{FD} \boldsymbol{v}^{FD}(\boldsymbol{r}) + \kappa^{SD} \boldsymbol{v}^{SD}(\boldsymbol{r}) + \kappa^{SQ} \boldsymbol{v}^{SQ}(\boldsymbol{r}) + \mathcal{O}(r^{-5}), \qquad (22)$$



Fig. 8: Flow field of a spheroidal squirmer with the aspect ration two. For each of the three pairs, the left image is the flow field in the laboratory reference frame and the right image in the body-fixed frame. (Left) Pusher with $\beta = -3$, (Middle) neutral squirmer ($\beta = 0$), and (Right) puller with $\beta = 3$. The magnitude of the velocity field is color coded logarithmically.

where

$$\boldsymbol{v}^{FD}(\boldsymbol{r}) = \frac{\boldsymbol{r}}{r^3} \left(\frac{3z^2}{r^2} - 1\right),\tag{23}$$

$$\boldsymbol{v}^{SD}(\boldsymbol{r}) = \frac{1}{r^3} \left(-\boldsymbol{e}_z + \frac{3z\boldsymbol{r}}{r^2} \right), \tag{24}$$

$$\boldsymbol{v}^{SQ}(\boldsymbol{r}) = \frac{3}{r^4} \left(\frac{5z^2 \boldsymbol{r}}{r^3} - \frac{2z \boldsymbol{e}_z + \boldsymbol{r}}{r} \right), \tag{25}$$

which decay like r^{-2} , r^{-3} , and r^{-4} [98], respectively, and $\kappa^{FD} = P/8\pi\eta$, $\kappa^{SD} = -v_0R^3/2$, and $\kappa^{SQ} = PR^3/8\pi\eta$. Note that in Eqs. (23)-(25) the swimming direction *e* points along the positive *z* axis, and *e_z* is the unit vector along that axis. Figure 7 depicts flow fields of the various kinds of microswimmers. In the far field, the flow fields of pushers and pullers are given by the force-diploe field (23) (or Eq.(20)).

The assumption of a spherical shape is adequate for swimmers like, e.g., *Volvox*, however, the shapes of other microswimmers (*E. coli, Chlamydomonas, Paramecium*) are nonspherical. Here, an extension of the squirmer concept to spheroidal objects has been proposed [97, 99]. Figure 8 depicts flow fields of a spheroidal squirmer with the aspect ratio of two for the various kinds of dipolar terms (Eqs. (23), (24)) in the laboratory and body-fixed reference frame. The near-field modifications by the finite-size swimmer is clearly visible in comparison with Fig. 6 (left). Moreover, pusher and puller exhibit a stagnation point in front or back, respectively, in the body-fixed reference frame.

3.5 Squirmer cooperative locomotion

Hydrodynamic interactions substantially affect the properties of active particles. An example is the cooperative motion of two spheroidal squirmers confined in a thin slit by two no-slip walls (cf. Fig. 9). Initially, the squirmer's surface-to-surface distance is $d_s = 3.5a$ and the angle between their swim directions $\theta_0 = 3\pi/8$. Due to the setup, the squirmers initially approach each other and collide at $tv_0/\sigma \approx 0.5$ (cf. Fig 9 (right)). The (persistence) Péclet number $Pe = v_0/(2D_R^{\perp}\sigma) \approx 60$ is sufficiently large, such that the squirmer orientation has hardly changed before collision. When the neutral squirmers collide, they initially align parallel ($\cos \theta \approx 1$ at $tv_0/\sigma \approx 1$ in Fig. 9), but their trajectories start to diverge immediately thereafter. Pushers remain parallel for an extended time window, which is expected as pushers





Fig. 9: (*Left*) Definition of the geometry of the confined squirmers, their orientation, and distance. (Middle) Flow streamlines of two pullers swimming cooperatively (laboratory reference frame). The magnitude of the velocity field is color coded logarithmically. (Right) Average surface-to-surface distance d_s and orientation of squirmers, where $\cos(\theta) = e_1 \cdot e_2$, as function of time. The solid blue, dashed black, and dotted red lines correspond to pullers ($\beta = 4$), neutral squirmers, and pushers ($\beta = 4$). The standard deviation of the blue line is indicated by the cyan shaded region [54].

are known to attract each other [89], but at $tv_0/\sigma \approx 3$ (cf. Fig. 9) their trajectories diverge as well. This is probably due to noise, since we observe several realizations where pushers remain parallel for an extended time. Interestingly, pullers, which are known to repel each other when swimming in parallel [89] (cf. Fig. 7 for the flow field), swim cooperatively and reach a stable orientation with $\langle \cos(\theta) \rangle \approx 0.77$ shortly after they collided $(tv_0/\sigma \approx 1)$. Thereby, their cooperative swimming velocity is about $0.8v_0$. The flow field of this stable state, determined by MPC simulations, is shown in Fig. 9 (middle). Note that the velocity field in the swimming plane is left-right symmetric, and that there is a stagnation point in the center behind the swimmers. This point actually corresponds to a line normal to the walls [97]. In contrast, in Ref. [90] a cooperative swimming mode for spherical squirmers has been observed (see Fig. 22(c) of Ref. [90]). However, this cooperative swimming—termed pair-swimming by the authors—is unstable to perturbations that displace one swimmer out of the swimming plane [90]. Since our simulations and those of Ref. [89] include thermal fluctuations, we consequently do not observe the cooperative swimming mode of Ref. [90].

Hence, the stable close-by cooperative swimming of pullers is governed by the squirmer anisotropy, by the hydrodynamic interactions between them and, importantly, between pullers and confining surfaces.

3.6 Squirmer cluster formation

As discussed in Sec. 2.6, activity leads to MIPS in ABPs systems. Detailed studies reveal substantial changes in the phase behavior of active particles in the presence of hydrodynamic interactions [96,100–103]. Hydrodynamic simulation studies of disks (2D system) [100], where MIPS is most pronounced for APBs, show no evidence for a bulk phase separation. The qualitative different behavior is attributed to an emergent faster decorrelation of the squirmers swim-



Fig. 10: (Left) Time dependence of the orientation correlation function $C_e(t)$ for a system of spherical squirmers and active Brownian particles. The dashed line indicates the exponential decay with the rotational diffusion coefficient D_R . (Right) Snapshot of a configuration of squirmers. The two-dimensional packing fraction is $\phi^{2D} = 0.6$ and the Péclet number $P_e = 115$.

ming direction due to HIs compared to the Brownian rotational motion of ABPs. Similar results have been found for spherical squirmers confined in a narrow slit [104]. Figure 10 displays a snapshot of the structure of a system of neutral squirmers. Neither large clusters nor a pronounced order is present for neutral squirmers (compare with Fig. 4). Similarly, no MIPS is observed for pusher and pullers. In addition, Fig. 10 shows the time dependence of the orientation correlation function of the propulsion direction $C_e(t) = \sum_{i=1}^{N_s} \langle e_i(t) \cdot e_i(0) \rangle / N_s$, where N_s is the number of squirmers in the system. The significant faster decay of the correlation function function function of an individual squirmer decays similarly to the correlation function of an APB. Hence, the enhanced decay is a consequence of inter-squirmer interactions.

Structure formation is decisively affected by the shape of a microswimmer (squirmer). On the one hand, it is expected that an elongated shape enhances parallel alignment due to steric interactions, an effect already present for elongated ABPs [78, 105]. On the other hand, hydrodynamic interactions prevent stable aligned states, at least for pusher and neutral squirmers. The question is how hydrodynamics ultimately affects MIPS of elongated squirmers keeping in mind that hydrodynamics suppresses MIPS for spherical squirmers. Our computer simulations show that hydrodynamics enhances cluster formation for spheroidal squirmers. This is illustrated in Fig. 11. At the same Péclet number, Pe = 12, spheroids with the aspect ratio three exhibit evidently a highly compact and large-scale aggregate in contrast to the gas-like spherical system. Interestingly, hydrodynamic interactions enhance cluster formation, as shown in Fig. 11 (right). ABPs show the least tendency to form clusters for Pe = 12, followed by pushers and neutral squirmers; pullers show the highest tendency for cluster formation, already at rather low densities.

4 Conclusions

Active matter exhibits a spectrum of unusual and fascinating phenomena, and offers many promising avenues for creating novel materials with tunable properties. Most remarkably is the intriguing collective behavior with emerging large-scale turbulence-like flow, or motility-





Fig. 11: (*Left*) Snapshots of spherical and spheroidal (aspect ratio three) squirmer configurations. (*Right*) Density-aspect ratio state diagram for ABPs, pullers ($\beta = 1$), neutral squirmers ($\beta = 0$), and pushers ($\beta = 1$). Here, b_z denotes the long and b_x the short axis of the spheroid. The density is $\phi^{2D} = 0.6$ and $P_e = 12$.

induced phase separation. Here, active Brownian particles are an excellent model system to unravel the underlying generic principles of out-of-equilibrium systems. Remarkably, hydrodynamic interactions are essential for active matter, specifically biological microswimmers, and are able to completely alter the steady-state behavior of interacting motile particles. Hydrodynamic interactions are not only fundamental for the propulsion of microswimmers (see contribution E4), but also determine their behavior next to surfaces as well as the emergent collective dynamics and structures. Hydrodynamic interactions imply a very rich dynamics, which depends on the detailed swimming mechanism. We are only at the beginning of our strive to elucidate the properties of active matter. Specifically, the design of novel synthetic active agents and the control of existing (mostly biological) agents requires intensive studies in the future, both from the theoretical and the experimental side.

Appendices

A Fokker-Planck equation of ABP

Equivalently to the Langevin equations (1) and (2), the dynamics of an ABP is described by the Fokker-Planck equation for the distribution function $\psi(\mathbf{r}, \mathbf{e}, t)$,

$$\frac{\partial \psi}{\partial t} = -\frac{\partial}{\partial \boldsymbol{r}} \left(\left[v_0 \boldsymbol{e} + \frac{1}{\gamma} \boldsymbol{F} \right] \psi \right) + D_T \frac{\partial^2}{\partial \boldsymbol{r}^2} \psi + D_R \frac{\partial^2}{\partial \boldsymbol{e}^2} \psi.$$
(26)

Here, $\partial^2/\partial e^2$ is the Laplace operator in polar (2D) or spherical (3D) coordinates. There is typically no simple way to find a stationary-state $(\partial \psi/\partial t = 0)$ solution of this equation, because detailed balance is violated [14]. According to Refs. [68, 106], the stationary state distribution function, or the generalized potential [68], is required to examine detailed balance. A necessary and sufficient condition for the existence of a potential for Eq. (26) is $v_0 \partial e_\alpha/\partial e_{\alpha'} = 0$, $\forall \alpha, \alpha'$, (natural boundary conditions are assumed), which is only satisfied for $v_0 = 0$. Note that there is no drift term related to the orientation *e*. Hence, there is in general no potential and a stationary-state solution is difficult to obtain [68].

B Fokker-Planck equation of AOUP

The Fokker-Planck equation for the distribution function $\psi(\mathbf{r}, \mathbf{v}, t)$ of the AOUP dynamics (6) reads [67]

$$\frac{\partial}{\partial t}\psi = -\frac{\partial}{\partial \boldsymbol{r}}\left(\left[\boldsymbol{v} + \frac{1}{\gamma}\boldsymbol{F}\right]\psi\right) + 2D_R\frac{\partial}{\partial \boldsymbol{v}}\left(\boldsymbol{v}\psi\right) + D_T\frac{\partial^2}{\partial \boldsymbol{r}^2}\psi + \frac{2D_Rv_0^2}{3}\frac{\partial^2}{\partial \boldsymbol{v}^2}\psi.$$
 (27)

Due to the Gaussian nature of the stochastic processes, and the fact that a Gaussian is determined by its first and second moment, the full time-dependent solution of Eq. (27) can be obtained for the special case of a harmonic potential, i.e., a linear force [67, 68].

References

- [1] M. E. Cates and F. C. MacKintosh, Soft Matter 7, 3050 (2011).
- [2] J. Elgeti, R. G. Winkler, and G. Gompper, Rep. Prog. Phys. 78, 056601 (2015).
- [3] C. Bechinger, R. Di Leonardo, H. Löwen, C. Reichhardt, G. Volpe, and G. Volpe, Rev. Mod. Phys. 88, 045006 (2016).
- [4] D. Needleman and Z. Dogic, Nat. Rev. Mater. 2, 201748 (2017).
- [5] S. Ramaswamy, J. Stat. Mech. Theor. Exp. 2017, 054002 (2017).
- [6] A. S. Mikhailov and R. Kapral, Proc. Natl. Acad. Sci. USA 112, E3639 (2015).
- [7] R. Kapral and A. S. Mikhailov, Physica D 318-319, 100 (2016).
- [8] H. S. Muddana, S. Sengupta, T. E. Mallouk, A. Sen, and P. J. Butler, J. Am. Chem. Soc. 132, 2110 (2010).
- [9] K. K. Dey, S. Das, M. F. Poyton, S. Sengupta, P. J. Butler, P. S. Cremer, and A. Sen, ACS Nano 8, 11941 (2014).
- [10] R. G. Winkler, J. Elgeti, and G. Gompper, J. Phys. Soc. Jpn. 86, 101014 (2017).
- [11] Y. Harada, A. Noguchi, A. Kishino, and T. Yanagida, Nature 326, 805 (1987).
- [12] V. Schaller, C. Weber, C. Semmrich, E. Frey, and A. R. Bausch, Nature 467, 73 (2010).
- [13] F. Jülicher, K. Kruse, J. Prost, and J.-F. Joanny, Phys. Rep. 449, 3 (2007).
- [14] M. C. Marchetti, J. F. Joanny, S. Ramaswamy, T. B. Liverpool, J. Prost, M. Rao, and R. A. Simha, Rev. Mod. Phys. 85, 1143 (2013).
- [15] J. Prost, F. Jülicher, and J.-F. Joanny, Nat. Phys. **11**, 111 (2015).
- [16] A. Cordoba, J. D. Schieber, and T. Indei, RSC Adv. 4, 17935 (2014).
- [17] Y. Sumino, K. H. Nagai, Y. Shitaka, D. Tanaka, K. Yoshikawa, H. Chate, and K. Oiwa, Nature 483, 448 (2012).
- [18] C. P. Brangwynne, G. H. Koenderink, F. C. MacKintosh, and D. A. Weitz, J. Cell. Biol. 183, 583 (2008).
- [19] C. P. Brangwynne, G. H. Koenderink, F. C. MacKintosh, and D. A. Weitz, Phys. Rev. Lett. 100, 118104 (2008).
- [20] C. A. Weber, R. Suzuki, V. Schaller, I. S. Aranson, A. R. Bausch, and E. Frey, Proc. Natl. Acad. Sci. USA 112, 10703 (2015).

- [21] S. C. Weber, A. J. Spakowitz, and J. A. Theriot, Proc. Natl. Acad. Sci. USA 109, 7338 (2012).
- [22] A. Javer, Z. Long, E. Nugent, M. Grisi, K. Siriwatwetchakul, K. D. Dorfman, P. Cicuta, and M. Cosentino Lagomarsino, Nat. Commun. 4, 3003 (2013).
- [23] A. Zidovska, D. A. Weitz, and T. J. Mitchison, Proc. Natl. Acad. Sci. USA 110, 15555 (2013).
- [24] E. Lauga and T. R. Powers, Rep. Prog. Phys. 72, 096601 (2009).
- [25] H. C. Berg, E. Coli in Motion, Biological and Medical Physics Series (Springer, New York, 2004).
- [26] C. Brennen and H. Winet, Ann. Rev. Fluid Mech. 9, 339 (1977).
- [27] S. Y. Reigh, R. G. Winkler, and G. Gompper, Soft Matter 8, 4363 (2012).
- [28] J. Hu, M. Yang, G. Gompper, and R. G. Winkler, Soft Matter 11, 7843 (2015).
- [29] R. M. Macnab, Proc. Natl. Acad. Sci. USA 74, 221 (1977).
- [30] L. Turner, W. S. Ryu, and H. C. Berg, J. Bacteriol. 182(10), 2793 (2000).
- [31] D. B. Kearns, Nat. Rev. Microbiol. 8, 634 (2010).
- [32] H. H. Wensink, J. Dunkel, S. Heidenreich, K. Drescher, R. E. Goldstein, H. Löwen, and J. M. Yeomans, Proc. Natl. Acad. Sci. USA 109, 14308 (2012).
- [33] G. Popkin, Nature **529**, 16 (2016).
- [34] F. Wong, K. K. Dey, and A. Sen, Annu. Rev. Mater. Res. 46, 407 (2016).
- [35] J. Bialké, T. Speck, and H. Löwen, Phys. Rev. Lett. 108, 168301 (2012).
- [36] I. Buttinoni, J. Bialké, F. Kümmel, H. Löwen, C. Bechinger, and T. Speck, Phys. Rev. Lett. 110, 238301 (2013).
- [37] G. S. Redner, M. F. Hagan, and A. Baskaran, Phys. Rev. Lett. 110, 055701 (2013).
- [38] J. Palacci, S. Sacanna, A. P. Steinberg, D. J. Pine, and P. M. Chaikin, Science 339, 936 (2013).
- [39] A. Wysocki, R. G. Winkler, and G. Gompper, EPL 105, 48004 (2014).
- [40] M. E. Cates and J. Tailleur, Annu. Rev. Condens. Matter Phys. 6, 219 (2015).
- [41] M. C. Marchetti, Y. Fily, S. Henkes, A. Patch, and D. Yllanes, Curr. Opin. Colloid Interface Sci. 21, 34 (2016).
- [42] A. Zöttl and H. Stark, J. Phys.: Condens. Matter 28, 253001 (2016).
- [43] E. M. Purcell, Am. J. Phys. 45(1), 3 (1977).
- [44] K.-T. Wu, J. B. Hishamunda, D. T. N. Chen, S. J. DeCamp, Y.-W. Chang, A. Fernández-Nieves, S. Fraden, and Z. Dogic, Science 355, eaal1979 (2017).
- [45] URL http://blogs.brandeis.edu/science/files/2016/01/flat.jpg.
- [46] URL http://blogs.brandeis.edu/science/files/2016/01/flat.jpg.
- [47] URL https://microbewiki.kenyon.edu/index.php/ Salmonella_enterica_NEU2011.
- [48] URL https://justinsomnia.org/2006/06/ photo-of-the-day-black-sun-in-denmarkthose-are/.
- [49] J. Elgeti, U. B. Kaupp, and G. Gompper, Biophys. J. 99, 1018 (2010).

- [50] M. J. Lighthill, Comm. Pure Appl. Math. 5, 109 (1952).
- [51] J. R. Blake, J. Fluid Mech. 46, 199 (1971).
- [52] T. Ishikawa and T. J. Pedely, J. Fluid Mech. 588, 437 (2007).
- [53] I. Llopis and I. Pagonabarraga, J. Non-Newtonian Fluid Mech. 165, 946 (2010).
- [54] M. Theers, E. Westphal, G. Gompper, and R. G. Winkler, Phys. Rev. E 93, 032604 (2016).
- [55] J. R. Howse, R. A. L. Jones, A. J. Ryan, T. Gough, R. Vafabakhsh, and R. Golestanian, Phys. Rev. Lett. 99, 048102 (2007).
- [56] F. Peruani, L. Schimansky-Geier, and M. Bär, Eur. Phys. J. Spec. Top. 191(1), 173 (2010).
- [57] P. Romanczuk, M. Bär, W. Ebeling, B. Lindner, and L. Schimansky-Geier, Eur. Phys. J. Spec. Top. 202, 1 (2012).
- [58] Y. Fily and M. C. Marchetti, Phys. Rev. Lett. 108, 235702 (2012).
- [59] B. ten Hagen, R. Wittkowski, D. Takagi, F. Kümmel, C. Bechinger, and H. Löwen, J. Phys.: Condens. Matter 27, 194110 (2015).
- [60] M. Doi and S. F. Edwards, *The Theory of Polymer Dynamics* (Clarendon Press, Oxford, 1986).
- [61] M. Raible and A. Engel, Appl. Organometal. Chem. 18, 536 (2004).
- [62] R. G. Winkler, A. Wysocki, and G. Gompper, Soft Matter 11, 6680 (2015).
- [63] R. G. Winkler, Soft Matter 12, 3737 (2016).
- [64] T. Eisenstecken, G. Gompper, and R. G. Winkler, Polymers 8, 304 (2016).
- [65] A. P. Solon, M. E. Cates, and J. Tailleur, Eur. Phys. J. Spec. Top. 224, 1231 (2015).
- [66] É. Fodor, C. Nardini, M. E. Cates, J. Tailleur, P. Visco, and F. van Wijland, Phys. Rev. Lett. 117, 038103 (2016).
- [67] S. Das, G. Gompper, and R. G. Winkler, New J. Phys. in press (2017).
- [68] H. Risken, The Fokker-Planck Equation (Springer, Berlin, 1989).
- [69] R. Di Leonardo, L. Angelani, D. Dell'Arciprete, G. Ruocco, V. Iebba, S. Schippa, M. P. Conte, F. Mecarini, F. De Angelis, and E. Di Fabrizio, Proc. Natl. Acad. Sci. USA 107, 9541 (2010).
- [70] N. Koumakis, A. Lepore, C. Maggi, and R. Di Leonardo, Nat. Commun. 4, 2588 (2013).
- [71] J. Elgeti and G. Gompper, EPL 101, 48003 (2013).
- [72] Y. Fily, A. Baskaran, and M. F. Hagan, Soft Matter 10, 5609 (2014).
- [73] N. Koumakis, C. Maggi, and R. Di Leonardo, Soft Matter 10, 5695 (2014).
- [74] O. Sipos, K. Nagy, R. Di Leonardo, and P. Galajda, Phys. Rev. Lett. 114, 258104 (2015).
- [75] S. C. Takatori, W. Yan, and J. F. Brady, Phys. Rev. Lett. 113, 028103 (2014).
- [76] A. P. Solon, Y. Fily, A. Baskaran, M. E. Cates, Y. Kafri, M. Kardar, and J. Tailleur, Nat. Phys. 11, 673 (2015).
- [77] H. Chaté, F. Ginelli, G. Grégoire, and F. Raynaud, Phys. Rev. E 77, 046113 (2008).
- [78] M. Abkenar, K. Marx, T. Auth, and G. Gompper, Phys. Rev. E 88, 062314 (2013).
- [79] T. Vicsek, A. Czirók, E. Ben-Jacob, I. Cohen, and O. Shochet, Phys. Rev. Lett. 75, 1226 (1995).

- [80] F. Ginelli, Eur. Phys. J. Spec. Top. 225(11), 2099 (2016).
- [81] M. Theers and R. G. Winkler, Soft Matter 10, 5894 (2014).
- [82] S. Kim and S. J. Karrila, *Microhydrodynamics: principles and selected applications* (Butterworth-Heinemann, Boston, 1991), ISBN 0486317676.
- [83] J. K. G. Dhont, An Introduction to Dynamics of Colloids (Elsevier, Amsterdam, 1996).
- [84] K. Drescher, J. Dunkel, L. H. Cisneros, S. Ganguly, and R. E. Goldstein, Proc. Natl. Acad. Sci. USA 10940, 108 (2011).
- [85] T. Ishikawa, J. R. Soc. Interface 6, 815 (2009).
- [86] K. Drescher, R. E. Goldstein, and I. Tuval, Proc. Natl. Acad. Sci. USA 107, 11171 (2010).
- [87] A. Erbe, M. Zientara, L. Baraban, C. Kreidler, and P. Leiderer, J. Phys.: Condens. Matter 20, 404215 (2008).
- [88] G. Volpe, I. Buttinoni, D. Vogt, H. J. Kümmerer, and C. Bechinger, Soft Matter 7, 8810 (2011).
- [89] I. O. Götze and G. Gompper, Phys. Rev. E 82, 041921 (2010).
- [90] T. Ishikawa, M. P. Simmonds, and T. J. Pedley, Journal of Fluid Mechanics 568, 119 (2006).
- [91] A. A. Evans, T. Ishikawa, T. Yamaguchi, and E. Lauga, Phys. Fluids 23, 111702 (2011).
- [92] F. Alarcón and I. Pagonabarraga, J. Mol. Liq. 185, 56 (2013).
- [93] J. J. Molina, Y. Nakayama, and R. Yamamoto, Soft Matter 9, 4923 (2013).
- [94] T. Ishikawa and T. J. Pedley, Phys. Rev. Lett. 100, 088103 (2008).
- [95] K. Ishimoto and E. A. Gaffney, Phys. Rev. E 88, 062702 (2013).
- [96] A. Zöttl and H. Stark, Phys. Rev. Lett. 112, 118101 (2014).
- [97] M. Theers, E. Westphal, G. Gompper, and R. G. Winkler, Soft Matter 12, 7372 (2016).
- [98] S. E. Spagnolie and E. Lauga, J. Fluid Mech. 700, 105 (2012).
- [99] S. R. Keller and T. Y. Wu, J. Fluid Mech. 80, 259 (1977).
- [100] R. Matas-Navarro, R. Golestanian, T. B. Liverpool, and S. M. Fielding, Phys. Rev. E 90, 032304 (2014).
- [101] J. Blaschke, M. Maurer, K. Menon, A. Zöttl, and H. Stark, Soft Matter 12, 9821 (2016).
- [102] F. Alarcón, C. Valeriani, and I. Pagonabarraga, Soft Matter 13, 814 (2017).
- [103] N. Yoshinaga and T. B. Liverpool, Phys. Rev. E 96, 020603 (2017).
- [104] M. Theers, K. Qi, R. G. Winkler, and G. Gompper, *Hydrodynamically and sterically induced clusters of shperiodal squirmers* (2018), in preparation.
- [105] F. Ginelli, F. Peruani, M. Bär, and H. Chaté, Phys. Rev. Lett 104, 184502 (2010).
- [106] C. W. Gardener, Handbook of Stochastic Methods (Springer, Berlin, 1983).

E 4 Motility of Cells and Microorganisms

Thorsten Auth, Jens Elgeti, Roland G. Winkler, Gerhard Gompper Theoretical Soft Matter and Biophysics Institute of Complex Systems and Institute for Advanced Simulation Forschungszentrum Jülich GmbH

Contents

1	Introduction 1.1 Swimming Cells and Microorganisms 1.2 Biomimetic Microswimmers 1.3 Cell Migration	2 4 7 8
2	Cell Crawling on Surfaces 2.1 Actin Polymerization and Thermal Ratchets 2.2 A Minimal Model of Cell Motility	9 9 10
3	Life at Low Reynolds Numbers 3.1 Swimming at Low Reynolds Numbers 3.2 Microswimmer Flow Fields and Hydrodynamic Interactions	13 13 14
4	Swimming due to Flagellar Motion4.1Anisotropic Hydrodynamic Friction of Slender Bodies4.2Swimming Velocity of Beating Flagella and Sperm4.3Sperm Steering and Navigation4.4Bacteria Swimming near Surfaces	16 16 17 18 20
5	Hydrodynamic Synchronization: Metachronal Waves in Cilia Carpets	22
6	Summary and Conclusions	24
A	Oseen Tensor of Low-Reynolds-Number Hydrodynamics	25
B	Hydrodynamic Force Dipoles	27
С	Hydrodynamics near Planar Surfaces	27

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Motility and locomotion is an essential part of life. The search for food, for favorable environmental conditions (like temperature or acidity), and for partners for mating and reproduction, and the bisexual reproduction process itself all require locomotion. This is the case in the world of our daily experience, but also in the world of cells and micro-organisms, such as sperm, bacteria, algae, and parasites. Since these micro-organisms often live in an aqueous environment, locomotion usually means swimming or sliding.

Cell migration is also a central process in the development and maintenance of multicellular organisms. Tissue formation during embryonic development, wound healing and immune responses all require the orchestrated movement of cells in particular directions to specific locations. Cells often migrate on surfaces or in three-dimensional meshes or tissues, where migration usually means crawling. Also, cells migrate in response to specific external chemical signals and mechanical signals. Thus, an understanding of the mechanisms of cell migration holds promise for the development of novel therapeutic strategies for controlling invasive tumor cells.

An important common feature of all these systems is that they are persistently out of equilibrium, a state which is maintained by continuous energy consumption. This leads to many novel phenomena, which are not known in equilibrium systems. From the physical side, the term "active matter" has been coined for this new state of matter. Of course, non-equilibrium is the essence of all living matter.

Recently, there have been several attempts to utilize the insight into the motility of biological microswimmers to construct artificial microswimmers and self-propelled particles. These microswimmers sometimes mimic biological micro-organisms in their swimming mechanisms, sometimes employ entirely new types of propulsion. Synthetic microswimmers are interesting for applications as future microrobots to perform tasks on the micrometer scale – such as directed motion in the human body for medical applications. They are also interesting model systems to study the behavior of microswimmers without biological complications (such as an active response to external stimuli).

From the physics point of view, the field of microswimmers has been pioneered by Purcell [1] in 1977, who showed that swimming in the microworld is quite different from our macroscopic experience, because at the length scale of nano- and micrometers, swimming occurs effectively at very high viscosity – which would feel for us like swimming in honey. During the 40 years, which have passed since then, the understanding of swimming microorganisms has increased extensively. Recent review articles provide a good overview of the present state of research and the current research directions. Generic aspects of the emergent large-scale behavior of self-propelled particles and active matter have been discussed in Refs. [2–7]. The hydrodynamics of swimming has been reviewed in Refs. [7–11]. The behavior of microswimmers and self-propelled particles near surfaces has been summarized in Refs. [7, 12]. Aspects of bacterial motility have been addressed in Refs. [7, 13–15]. An overview of sperm motility and chemotaxis has been given in Refs. [16, 17]. The behavior of active colloids has been reviewed in Refs. [18, 19]. The dynamical properties of active Brownian particles have been discussed in Refs. [7, 20, 21]. Finally, the propulsion of synthetic swimmer on the nanoscale and the development of nanomachines have been addressed in Refs. [22, 23].



Fig. 1: (top) Peritrichious bacteria, like Salmonella, swim by a bundle of rotating helical flagella. From Ref. [27]. (bottom) Each bacterial flagellum is driven by a rotary motor embedded in the bacterial cell wall. The motor has a series of rings, each about 20 nanometers in diameter, with a rod inside. Attached to the rod is a curved "hook" protein linked to the flagellum. The flagellum is 5 to 15 micrometers long and made of thousands of repeating units of the protein flagellin. The motor is powered by the flow of sodium or hydrogen ions across the cell wall. This generates rotation frequencies of up to 15 Hz. From Ref. [28].



Fig. 2: Bacteria like E. coli and salmonella move by a "run-and-tumble" motion. During the "run" phase, the flagella form a bundle, and the bacterium moves forward in one direction. In the "tumble" phase, one or a few flagella reverse their rotational direction and leave the bundle. This induces a tumbling motion which changes the orientation of the bacterium randomly. CW denotes clockwise, CCW counter-clockwise rotation. From Ref. [31].

1.1 Swimming Cells and Microorganisms

Bacteria. Some bacteria, like Escherichia coli (E. coli) [29] and Salmonella, swim by rotating helical hairlike filaments called flagella. Several of such flagella are attached to the bacterial body, as shown in Fig. 1(top). The flagella are rotated by a motor complex, which consists of several proteins, and which is anchored in the bacterial cell wall [27, 29, 30]; the motor is connected to the flagellum by a flexible hook, see Fig. 1(bottom).

Bacteria like E. coli and salmonella swim in a "run-and-tumble" motion illustrated in Fig. 2. In the "run" phase (stage 1 in Fig. 2), the helical winding of all flagella is left-handed, and they are rotating counter-clockwise. The flagella form a bundle (see also Fig. 1(top)), and the bacterium moves forward in a direction determined by its long axis. At the beginning of the "tumble" phase, one flagellum reverses its rotational direction to clockwise (stage 2 in Fig. 2). This flagellum leaves the bundle, which implies a random reorientation of the bacterial orientation (stage 3 and 4). The reversal of the rotational direction is accompanied by a change of the helical handedness from left-handed to right-handed. At the end of the "tumbling" phase, all flagella start to rotate again in the same counter-clockwise direction (stage 5), the bundle reforms (stage 6), and the bacterium returns to a directional motion (stage 7 and 8).

The purpose of the "run-and-tumble" motion is to detect gradients in the concentration of chemicals (e.g. food) or temperature (to avoid regions of too high or too low temperature). This is achieved by extending the "run" phase in case of improving environmental conditions, and by shortening it in case of worsening conditions.

Sperm. Sperm cells consist of a (roughly spherical) head, which contains the genetic material, a short midpiece, which contains many mitochondria for energy production, and the eukaryotic flagellum, see Fig. 3. The structure of the flagellum is shown in Fig. 4. It consists of two central microtubules, which are surrounded by 9 double microtubules (microtubules are rather stiff filaments with a persistence length of about 1 mm). The microtubules are connected by many proteins (nexin links, central spokes, ...), which stabilize the structure. Motor proteins (dynein) connecting neighboring double microtubules cause an active bending of the structure



Fig. 3: Sperm swim by a snake-like motion of their flagellum. The time sequence (from left to right) of snapshots of the beat of sea-urchin sperm shows a sinusoidal traveling wave on the flagellum. Arrow indicate the wave propagation from the head to the tip. From Ref. [7]

by sliding the microtubules relative to each other.

A sperm cell is propelled though a fluid by a snake-like wiggling of its tail, as shown in Fig. 3. The flagellar beat is a propagating bending wave, with wavelengths smaller than the flagellar length. The beat can either be planar, as in Fig. 3, in particular for sperm swimming near surfaces, or be three-dimensional with a conical envelope [17, 33].

Paramecium and Opalina. *Paramecia* is a group of unicellular ciliated protozoa, which range from about 50 μ m to 350 μ m in length. Many hair-like extrusions — called cilia — cover the body, see Fig. 5(left), which allow the cell to move with a synchronous motion (like a caterpillar) at speeds of approximately 2,700 μ m/sec (12 body lengths per second). They generally feed on bacteria and other small cells. *Opalina* is a genus of protozoa found in the intestines of frogs and toads. It is without a mouth or contractile vacuole. The surface of Opalina is covered uniformly with cilia.

The structure of cilia and flagella of eukaryotic cells is very similar, see Fig. 4. The main structural difference between cilia and flagella is their length. A typical cilium is $10\mu m$ long, while sperm flagella are about $50\mu m$ long. The second large difference concerns their beat patterns. The ciliar beat has two distinct phases. During the *power stroke*, the cilium is stretched out straight and moves rather fast in one direction, while it bends, twists a little sideways and slowly retracts in the *recovery stroke*, as shown in Fig. 5(right). A particularly interesting feature of the beat pattern of cilia arrays on the surface of Opalina and Paramecium is the formation of "metachronal waves", as can be see in Fig. 5(left); the cilia do not beat synchronously, but in a pattern resembling a wheat field in the wind.

Different physical mechanisms for the beat of cilia and flagella have been suggested. Either the beat arises from a coupling of the activity of motor proteins with the curvature of the flagellum, where large curvature implies detachment of motors from the neighboring filament [36, 37], or it arises from the cooperative behavior of several motors pulling in opposite directions, with the "winners" pulling along the "losers" for a while, as in a tug-of-war [38, 39].

Chlamydomonas reinhardtii. Chlamydomonas reinhardtii is a single-celled green algae, about 10 μ m in diameter, that swims with two flagella, see Fig. 6(left). They have an "eye-spot" that



Fig. 4: The axoneme consists of 9 double microtubules, arranged around the perimeter, and two central microtubules. The microtubules are connected by many linker proteins. Motor proteins connecting neighboring double microtubules lead to active bending. From Ref. [32].



Fig. 5: (*Left*) The surface of Opalina is covered by many hairlike filaments called cilia. Opalina swims forward by a stroke-like wiggling of the cilia. The formation of metachronal waves is clearly visible. From Ref. [34]. (Right) The beat of cilia has two distinct phases, the power stroke and the recovery stroke. Snapshots are shown from Volvox somatic cells, imaged at 1000 frames per second. During the power stroke, the cilium is stretched out straight and moves rather fast in one direction (frames 1 - 11), while during the recovery stroke, it bends and slowly retracts (frames 13-27). Adapted from Ref. [35].



Fig. 6: (*left*) Chlamydomonas *swims by a breast-stroke-like motion of its two flagella. Steering in reaction to a light stimulus is facilitated by additional beats of one flagellum. Adapted from Ref. [43]. (right) Volvox is a colony of several thousand of cells, which form a hollow sphere. Daughter colonies inside are visible. Each cell has two flagella, which beat in a coordinated fashion to rotate or move the colony. From Ref. [44].*

senses light. When illuminated, C. reinhardtii can grow in a medium lacking carbon and energy sources. Widely distributed worldwide in soil and fresh water, C. reinhardtii is primarily used as a model organism in biology in a wide range of subfields [40]. The swimming motion of C. reinhardtii resembles the human breast-stroke: the flagella are pulled back in a nearly straight shape, and are then bend over and pushed forward again. The oscillatory velocity field induced by swimming C. reinhardtii in stabilized thin liquid films has been observed directly in time-resolved measurements recently [41,42].

Volvox. Volvox is another green algae. It forms spherical colonies of up to 50,000 cells [44]. Each mature Volvox colony is composed of numerous flagellate cells similar to Chlamy-domonas, embedded in the surface of a hollow sphere, see Fig. 6(right). The cells swim in a coordinated fashion, with distinct anterior and posterior poles. The cells have eye-spots, more developed near the anterior, which enable the colony to swim towards light. Recently, the flow field around freely swimming Volvox has been measured directly [41,45].

1.2 Biomimetic Microswimmers

Locomotion on the nanoscale through a fluid environment is one of the grand challenges confronting nanoscience today [22]. The vision is to synthesize, probe, understand, and utilize a new class of motors made from nanoscale building blocks that derive on-board or off-board power from in-situ chemical reactions. The generated mechanical work allows these motors to move through a fluid phase while simultaneously or sequentially performing a series of tasks. A large variety of such swimmers have been constructed recently, from bimetallic nanorods [22] and Janus colloids [46] to artificial sperm [47,48].

Artificial microswimmers can be constructed by using similar design principles as those found in biological systems. A by now classical example is a swimmer which mimics the propulsion mechanism of a sperm cell [47]. The flagellum is constructed from a chain of magnetic colloid particles and is attached to a red blood cell, which mimics the sperm head. This artificial swimmer is set into motion by an alternating magnetic field, which generates a sidewise oscillatory deflection of the flagellum. However, the swimming motion is not the same as for a real sperm cell; the wiggling motion is more a wagging than a traveling sine wave, and generates a swimming motion toward the tail end, opposite to the swimming direction of sperm.

A more recent example of a biohybrid swimmer, which mimics the motion of sperm consists of a polydimethylsiloxane filament with a short, rigid head and a long, slender tail on which cardiomyocytes (heart-muscle cells) are selectively cultured [48]. The cardiomyocytes contract periodically and deform the filament to propel the swimmer. This is a true microswimmer because it requires no external force fields.

Also, biomimetic systems have been designed to use artificial cilia for transport. For example, externally-actuated semi-flexible strings, like chains of magnetic beads [47, 49, 50], self-assembled microtubule bundles [51], ionic-polymer metal composite actuators [52], and even semi-flexible micro-pillars mounted an a basis of piston-like actuators [53] have been proposed to employ the cilia propulsion mechanism in artificial nano-machines and microfluidic devices [54, 55].

1.3 Cell Migration

Both during development and in the developed organism, there are specialized cell which have a high disposition for migration, closely connected to their function within the organism [25, 26]. A good example is the repair of open wounds in the skin. The bleeding is first stopped passively by von Willebrand factor proteins and platelets, which arrive via the blood flow (see Chapter E2 by D. Fedosov), followed by the coagulation cascade. Then, the immune systems responds and white blood cells remove damaged tissue, and fight possible intruders like bacteria. Finally, fibroblast cells crawl toward the wound, excrete collagen, and thereby regenerate the extracellular matrix. At the same time, the epithelial cells around move collectively to close the wound. Other important examples are the wiring of the brain and the neuronal network (see Chapter E7 by M. Diesmann), and cancer growth in its late stage, as metastasis takes place and malignant cells invade the body.

The ability of cells to move has fascinated biologists for a long time, and more recently has captured the attention of physicists [26]. Many aspects of cell motility are not yet well understood, like (i) how a cell polarizes to head in a certain direction, (ii) how it maintains its integrity and shape during its motion, (iii) how it transfers the force to the substrate, which may have different properties and various parts of the organism, and (iv) how it is able to move at all.

It is now well accepted that cell motility is driven by the activity of the cytoskeleton, the scaffold of actin filaments and microtubules (see Chapters C3 by S. Köster and C6 by K. Kruse), which is kept in a persistent out-of-equilibrium state due to continuous energy consumption (with ATP as the main energy source), to generate the propulsion forces [56]. Of course, the detailed understanding the motility machinery of a single cell is only the first step to understand their individual motion patterns (persistence, crawling speed, reaction to environmental stimuli), as well as their collective motion.



Fig. 7: Cells crawl over a surface in four subsequent steps: (i) protrusion of the leading edge, (ii) adhesion of the leading edge to the substrate, (iii) de-adhesion of the trailing edge, and (iv) movement of the cell body. From Ref. [57].

2 Cell Crawling on Surfaces

2.1 Actin Polymerization and Thermal Ratchets

Filamentous actin (F-actin) forms from its monomeric form, globular actin (G-actin), by polymerization. It is a unusual polymer, because it is a polar filament. Under normal physiological conditions, it has different polymerization and depolymerization rates at its two ends. Therefore, the usual steady state is not passive, but rather a balance between polymerization at one end, and depolymerization at the other end, which is called "treadmilling". Now, for the cytoskeleton near the leading edge of the cell, the free ends of the actin filaments are polymerizing and thereby pushing against the plasma membrane. Since at the same time the cytoskeleton is anchored to the substrate at various adhesion sites, this implies a force which propels the whole cell forward (see Fig. 7).

Let us take a closer at this propulsion mechanism [58]. As the filaments polymerize, their Brownian motion impinges on the load (the plasma membrane) exerting a pressure. However, to add a monomer to the free filament end, a thermal fluctuation must create a big enough gap to permit intercalation of an additional monomer. For a filament approaching the load under an angle θ , the required fluctuation amplitude is $\Delta = \delta \cos(\theta)$, where $\delta \simeq 2.7$ nm is half the size of an actin monomer, see Fig. 8. The frequency with which gaps appear, together with the local monomer concentration, M, determines whether, and how fast, the cytoskeleton edge can advance. A freely polymerizing tip would grow at velocity

$$v_p = \Delta(k_{\rm on}M - k_{\rm off}),\tag{1}$$

with polymerization on- and off-rates k_{on} and k_{off} , respectively. However, because of the load, the actual propagation velocity of the front is less than v_p .

The dominant mode of thermal motion of a single filament is a bending undulation. For a filament of length ℓ and persistence length ξ_p , these bending undulations can be shown to be equivalent to an effective one-dimensional spring, see Fig. 8, with an elastic spring constant [58]

$$\kappa(\theta) = 4\xi_p k_B T / (\ell^3 \sin^2(\theta)), \tag{2}$$

i.e., the effective spring is very stiff for perpendicular, and very soft for parallel orientation of the filament to the membrane. The statistical motion of the filament tip is subject to the



Fig. 8: (*left*) Schematic of a free actin filament tip of length ℓ impinging on a load at an angle θ . A filament tip can add a monomer only by a bending fluctuation of amplitude A. The polymerization rate is $k_{on}M - k_{off}$, where M is the monomer concentration. The actin network behind the last cross-link is regarded as a rigid support. (right) The simplified mechanical equivalent of (a). The bending elasticity is mimicked a spring constant, κ , given by Eq. (2). y is the equilibrium distance of the tip from the load, and x is the deviation of the tip from its equilibrium position. From Ref. [58].

harmonic restoring force of the effective spring as well as the load force, which can be described by a Fokker-Planck equation for the probability distribution of the tip position at time t. The analysis of this equation, in the limit that the thermal fluctuations of the tip are much faster than the polymerization rate, yields the load-velocity relation [58]

$$v_p = \Delta(k_{on} M p(\kappa(\theta), f) - k_{off}), \qquad (3)$$

where $p(\kappa, f)$ is the steady-state probability of a gap of width Δ to open up between the tip and the load with force f. A crucial feature is that the growth velocity is a non-monotonic function of the angle θ , which passes through a maximum at a critical angle θ_c . The reason is that thermal fluctuations may not be able to bend a stiff filament directed perpendicular to the membrane sufficiently to permit intercalation of G-actin, whereas filaments growing parallel to the membrane cannot exert any thrust.

Actin forms a branching network near the leading edge, with a branching angle of about 70°; this results in actin-membrane angles $\theta = 35^{\circ}$ in the symmetric case and $\theta = 0$ and $\theta = 70^{\circ}$ in the very asymmetric case [56, 59]. This is consistent with the result (3) of the ratchet model above that finite angles are best suited for high pushing efficiency.

2.2 A Minimal Model of Cell Motility

The shapes and motion of motile cells are determined by the dynamics of the cytoskeletal filaments and the cell membrane, cell-substrate adhesion [60, 61], and cell-cell [62] or cell-boundary interactions [63]. In order to understand these intricate interplay of the different contributions, it is necessary to construct models of various complexity and sophistication [64–66]. A simple two-dimensional model system [66] to theoretically and numerically study cell motility driven by the actin treadmilling, is a composite active particle that consists of self-propelled


Fig. 9: Run-and-circle motion of a mobile ring (radius R_m) with inner structure of selfpropelled rods. System with $N_r = 40$ attached-pushing and attached-pulling rods, Péclét number Pe = 25, repulsion strength $E_r/k_BT = 10$ and penetrability coefficient Q = 2.5. (left) The red curve represents the trajectory of the center of the ring, the black dot represents the position of the center of the ring at the last simulation frame. (right) The rod orientations and positions inside the ring for the last simulation frame. Arrows indicate pulling and pushing rods. The rods are colored according to their angle with respect to the radial direction. From Ref. [66].

rods, which are confined inside a rigid ring, see Fig. 9. Here, the basic idea is to mimic the dynamics and forces of cytoskeletal filaments by rod-like self-propelled particles: actin and microtubule treadmilling and polymerization forces are modeled by the self-propulsion of rigid rods with constant length. The ring represents the plasma membrane. Rods and membrane interact with a repulsive interaction potential that confines the rods inside the ring and also allows the control of the sliding friction between rods and membrane. Such complex self-propelled particles display several different types of motility, such as persistent motion, random motion, circling, and run-and-circle motion, depending on the rod-alignment interaction and the number of rods – as will be explained in more detail below.

Both the self-propelled rods and the ring-like membrane are described by chains of beads (diameter σ). Rod-beads interact via a repulsive, separation-shifted Lennard-Jones potential

$$\phi(r) = 4\epsilon \left[\left(\frac{\sigma^2}{\alpha^2 + r^2} \right)^6 - \left(\frac{\sigma^2}{\alpha^2 + r^2} \right)^3 \right] + \phi_0 \tag{4}$$

for $r \leq r_{cut}$, and zero otherwise. This potential generates a *finite* energy barrier E_r for two beads to completely overlap. Rods are propelled by a force F_p on each bead along the instantaneous rod (polar) orientation. The behavior of the rods is thus characterized by the dimensionless quantities Péclét number Pe and penetrability coefficient Q,

$$Pe = L_r F_p / (k_B T), \quad Q = L_r F_p / E_r, \tag{5}$$

where L_r is the number of beads of a rod. Alternatively, Pe and E_r/K_BT can be employed. In addition, the ring feels a viscous friction with the embedding medium characterized by the viscosity η .

It is important to emphasize that the rods are attached with their front or tail bead to the membrane, because the purpose of the model is to study the interaction of the cytoskeleton with the



Fig. 10: Complex self-propelled rings with $N_r = 16$ and $\eta = 0.005k_BT\tau_0/L_r$. The first row shows snapshots of the inner structure of the cell at the end of the simulation runs, with each rod colored based on its orientation with respect to the radial direction. The middle row shows kymographs of the rods dynamics, the radial axis represents the time axis and the tangential axis represents the angular axis; the color indicates the rod density along the ring. The third row shows trajectories of the ring centers. The final position is represented by a black dot. From Ref. [66].

membrane. This implies that both pushing and pulling rods are considered. Here, pushing rods correspond to the polymerizing actin filaments at the leading edge of a cell, compare Sec. 2.1. However, the cytoskeleton can also pull on the membrane via anchoring proteins [67].

Results of Brownian Dynamics simulations of this model are shown in Figs. 10 and 11. Depending on the propulsion strength compared to thermal motion, characterized by the Péclét number Pe, and the propulsion strength compared to the repulsive interaction between rods, characterized by the penetrability coefficient Q, rods either cross each other easily (low E_r/k_BT , high Pe) or align and form clusters (high E_r/k_BT , low Pe). This leads to random motion in the former, to persistent directed motion in the latter case, see Figs. 10. However, other more complex types of motion are also possible, such as circling and run-and-circle motion. This behavior is seen, for example, when clusters of pulling and pushing rods compete with each other.

The interesting aspect of this model is that the internal degrees of freedom can react to the environmental conditions. Thus, the model can easily be extended to include shape deformations, membrane proteins which interact with the cytoskeleton, and the collective behavior of many cells.



Fig. 11: Phase diagrams for mobile rings with $N_r = 40$ rods and $\eta = 0.005k_BT\tau_0/L_r$, for different Peclet numbers Pe and energy barriers E_r . In mixed systems, half of the rods are pushing rods, the other half pulling rods. From Ref. [66].

3 Life at Low Reynolds Numbers

3.1 Swimming at Low Reynolds Numbers

The dynamics of fluids — hydrodynamics — is described by the (incompressible) Navier-Stokes equation,

$$\rho\left(\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla)\mathbf{v}\right) = \eta \nabla^2 \mathbf{v} - \nabla p + \mathbf{f}_{ext}$$
(6)

where ρ is the fluid density, η the fluid viscosity, $\mathbf{v}(\mathbf{r}, t)$ the position- and time-dependent fluid velocity field, $p(\mathbf{r}, t)$ the pressure field, and $\mathbf{f}_{ext}(\mathbf{r}, t)$ an external body force. The Navier Stokes equation can be written in dimensionless form by scaling all lengths with a characteristic length L, velocities by a characteristic velocity v_0 , and time by L/v_0 . This implies

$$\mathfrak{Re}\left(\frac{\partial \mathbf{v}'}{\partial t'} + (\mathbf{v}' \cdot \nabla)\mathbf{v}'\right) = \nabla^2 \mathbf{v}' - \nabla p' + \mathbf{f}'_{ext}$$
(7)

where the prime denotes dimensionless quantities. Here, $\Re \mathfrak{e}$ is the Reynolds number, which is found to be

$$\mathfrak{Re} = \rho v_0 L / \eta . \tag{8}$$

For microswimmers, the Reynolds number is typically very small, because the characteristic lengths and velocities are small. For a swimmer of length $L = 50\mu$ m and a velocity of 10 body lengths per second (which is already a large velocity), the Reynolds number in water (with a kinematic viscosity $\eta/\rho = 10^{-6}m^2/s$) is $\Re e = 0.025$. In this case, the nonlinear contributions on the left-hand-side of Eq. (7) can be neglected, leading to the simpler Stokes equation,

$$\nabla p - \eta \nabla^2 \mathbf{v} = \mathbf{f}_{ext} \tag{9}$$

Although this equation is much simpler than the Navier-Stokes equation, it is not possible to solve analytically for microswimmer motion – even in simple geometries.

At surfaces, the fluid velocity is typically very small, because the collisions of fluid molecules with the surface imply that the molecules are scattered backwards and thereby transfer momentum parallel to the wall. Thus, no-slip boundary conditions, with $\mathbf{v}(\mathbf{r}) = 0$ at the surface, are usually employed.



Fig. 12: Schematics of the flow field of dipole swimmers, (a) pusher and (b) puller. From *Ref.* [7].

It has been recognized by Purcell [1] that life at low Reynolds numbers is quite different from the life of our everyday experience. Since low Reynolds numbers can either be achieved by small length scales (and small velocities) or by very high viscosities, a swimming microorganism experiences similar dynamics as a human swimming in honey. Since inertia effects are negligible, forward motion stops immediately when the "motor" stops working. Another important consequence of low Reynolds number is that the Stokes equation is time-reversible. Therefore, a time-reversible motion of a micro-organism will generate an oscillatory, but no directed motion. This has been called the "scallop theorem" by Purcell [1]. It means that just by opening and closing its two shells, a mussel cannot move forward. As explained above for the Purcell-swimmer, a second degree of freedom is required to generate a sequence of moves which is not time reversible.

3.2 Microswimmer Flow Fields and Hydrodynamic Interactions

Most microswimmers move autonomously, with no external force applied, and hence the total interaction force of the swimmer on the fluid, and vice versa, vanishes. In the simplest case, which actually applies to many microswimmers like bacteria, spermatozoa, or algae, the farfield hydrodynamics (at distances from the swimmer much larger than its size) can well be described by a force dipole [8,9]. This has been confirmed experimentally for *E. coli* [68]. Two classes of such dipole swimmers can be distinguished, as shown schematically in Fig. 12. If the swimmer has its motor in the back, and the passive body drags along the surrounding fluid in front, the characteristic flow field of a "pusher" emerges, see Fig. 12(a). Similarly, if the swimmer has its motor in the front, and the passive body drags along the surrounding fluid behind, the characteristic flow field of a "puller" develops, see Fig. 12(b). Thus, sperm, E. coli and salmonella are pushers, because they have active propelling flagella in the rear, and a passive head in the front. In contrast, C. reinhardtii is a puller, because the flagella in the front push the fluid backwards, while the body in the rear remains passive. It is important to notice that the flow fields of pushers and pullers look similar, but with opposite flow directions. This has important consequences for the interactions between swimmers, and of swimmers with walls, as will be explained below.

A first idea about the effect of hydrodynamic interactions on the dynamics of microswimmers



Fig. 13: Velocity fields for fixed parallel pairs of squirmers, for (left) pusher (with $B_2/B_1 = -3$) and (right) puller (with $B_2/B_1 = +3$), with Péclet number Pe = 1155. Swimmers move to the right. Streamlines serve as a guide to the eye. Only a fraction of the simulation box is shown. (Due to the finite resolution of the measured velocity field, some streamlines end on the squirmers' surfaces.) From Ref. [70].

near surfaces can be obtained again from a far-field approximation, which applies when the size of the swimmer is much smaller than its distance from the surface, so that it can be approximated by a force dipole, compare Appendix C. In this limit, analytical expressions for the torque on a swimmer and its drag velocity towards the surface can be obtained [69]. For a swimmer with an orientation angle θ of the dipole with respect to the vertical (z) direction (i.e. $\theta = 90^{\circ}$ when the swimmer moves parallel to the wall), the induced velocity at a distance z away from the no-slip wall is given by [69]

$$u_z(\theta, z) = -\frac{3P}{64\pi\eta z^2} \left(1 - 3\cos^2\theta\right).$$
(10)

As shown in Appendix B, the dipole strength P is given by the product of dipole length and force. The dipole force of sperm can be estimated from the friction force of the head, which — for a sufficiently large head — balances the pushing force of the tail, and is thus proportional to ηvL (compare Eqs. (15) and (18)), where v is the swimming velocity and L the sperm length. The dipole strength therefore is

$$P \sim \eta v L^2 \,. \tag{11}$$

Equation (10) allows several interesting predictions. First, the hydrodynamic interaction is attractive for swimmers oriented nearly parallel to the wall (with θ near 90°), but becomes repulsive for swimmers oriented nearly perpendicular to the wall (with θ near 0°). Second, the hydrodynamic interaction is long-ranged, with a $1/z^2$ decay with increasing distance from the surface. Finally, the hydrodynamic interaction has been shown not only to generate a force on the swimmer, but also a torque [69]. It turns out that this torque always acts to align the swimming cell parallel to the surface – the orientation in which the hydrodynamic force is attractive.

Similarly, the interaction of two microswimmers at long distances is determined by their dipole

flow fields. The dipole approximation predicts that the interactions of pushers and pullers have opposite sign, because their dipole strengths P have opposite signs (compare Eq. (56) of Appendix C). This behavior can be seen explicitly by considering the flow fields of two parallel-swimming pushers or pullers [70]. The results of mesoscale simulations of two squirmers shown in Fig. 13 demonstrate that for *pushers*, the fast backward flow velocity in the rear part extracts fluid from the gap between the swimmers, and thereby induces attraction (Fig. 13a); in contrast, for *pullers*, the fast backward flow velocity in the front part injects fluid into the gap, and thereby induces repulsion (Fig. 13b). At the small squirmer separation show in Fig. 13, the interaction is dominated by the hydrodynamic near-field, and the far-field approximation does not allow quantitative predictions anymore.

4 Swimming due to Flagellar Motion

4.1 Anisotropic Hydrodynamic Friction of Slender Bodies

A microorganism is able to swim forward in a fluid by wiggling or rotating a flagellum, because the hydrodynamic friction of a long, slender body in a viscous environment is *anisotropic*. This can be demonstrated easily for a long and thin rod of diameter d and length L: it experiences less friction when pulled along its axis than perpendicular to it.

Let us consider the motion of the rod, which may be oriented in the x-direction, with an applied force in direction \hat{e} . We approximate the rod as a sequence of touching spheres of diameter d. The friction coefficient ζ of a sphere (with no-slip boundary conditions) moving in a viscous fluid under a force F is given by

$$\mathbf{F} = \zeta_s \mathbf{v}$$
, where $\zeta_s = 3\pi \eta d$ (12)

is the Stokes friction coefficient. The dynamics of bead n in a chain with other beads is determined in general by the Langevin equation [71]

$$\mathbf{v}_{n} = \frac{\partial}{\partial t} \mathbf{R}_{n} = \sum_{m} \mathbf{H}_{nm} \cdot \left[-\frac{\partial U}{\partial \mathbf{R}_{m}} + f_{m} \right]$$
(13)

where $U(\mathbf{R}_1, \mathbf{R}_2, ...)$ is the interaction potential among the beads (including, in particular, the bonds to nearest neighbors), \mathbf{H}_{nm} for $n \neq m$ is the Oseen tensor (see Eqs. (44) and (50) derived in Appendix A), $\mathbf{H}_{nn} = \mathbf{I}/\zeta_s$, and f_m is an external force on bead m. Lets us now assume that the force **f** on all beads is the same, and that the interaction between beads can be neglected (this is only valid for short times, but should not affect the calculation of the friction coefficient) [72]. Then,

$$\mathbf{v}_{rod} \equiv \mathbf{v}_n = \left[\mathbf{I} / \zeta_s + \sum_{m \neq n} \mathbf{H}_{nm} \right] \cdot \mathbf{f} \ . \tag{14}$$

The anisotropic friction coefficients of a rod are defined by

$$\mathbf{F} = \zeta_{\parallel} \mathbf{v}_{\parallel} + \zeta_{\perp} \mathbf{v}_{\perp} \ . \tag{15}$$

Calculations based in Eq. (15), with constant friction coefficients ζ_{\parallel} and ζ_{\perp} , are denoted "resistive-force theory". To calculate ζ_{\parallel} and ζ_{\perp} for a rod, it is easiest to consider the special cases of a rod

pulled parallel and perpendicular to its long axis. Then, with the sum replaced by an integral, for a rod oriented in the x-direction, Eq. (51) gives for the central bead

$$\mathbf{v}_{rod} \equiv \mathbf{v}_n = f \left[\frac{\hat{\mathbf{e}}}{\zeta_s} + \frac{2}{8\pi\eta d} \int_d^{L/2} \mathrm{dx} \, \frac{1}{x} \left(\hat{\mathbf{e}} + (\hat{\mathbf{e}}_x \cdot \hat{\mathbf{e}}) \hat{\mathbf{e}}_x \right) \right] \,. \tag{16}$$

with $\mathbf{f} = f\hat{\mathbf{e}}$. Because $(\hat{\mathbf{e}}_x \cdot \hat{\mathbf{e}})\hat{\mathbf{e}}_x = 1$ and 0 for parallel and perpendicular orientation, respectively, we find immediately that

$$\zeta_{\perp} = 2\zeta_{\parallel} \tag{17}$$

in the limit of very long rods. It is therefore a factor 2 easier to pull a long rod along its axis than perpendicular to it. Finally, the integration in Eq. (16) together with F = (L/d)f gives

$$\zeta_{\perp} = \frac{4\pi\eta L}{\ln(L/d)} \tag{18}$$

to leading order in L/d. The logarithmic divergence is a result of the long-ranged hydrodynamic interaction of different parts of the rod, which *reduce* the friction coefficient compared to that of a chain of non-interacting beads.

4.2 Swimming Velocity of Beating Flagella and Sperm

The result (17) together with Eq. (15) can now be used to calculate the swimming velocity of an sinusoidally beating flagellum. In this case, the shape as a function of time t is described by

$$y(x,t) = A\sin(kx - \omega t) \tag{19}$$

where A is the beat amplitude, ω the beat frequency, and $k = 2\pi/\lambda$ the wave number. The velocity of a segment of the flagellum at x is then

$$v_y(x,t) = \frac{\partial y}{\partial t} = -A\omega\cos(kx - \omega t) \tag{20}$$

With the local tangent and normal vectors

$$\mathbf{t}(x,t) = N^{-1}(1,\partial y/\partial x) = N^{-1}(1,Ak\cos(kx-\omega t))$$
(21)

$$\mathbf{n}(x,t) = N^{-1}(-\partial y/\partial x, 1) = N^{-1}(-Ak\cos(kx - \omega t), 1)$$
(22)

respectively, with normalization $N^2 = 1 + (\partial_x y)^2$, the velocity $\mathbf{v}(x,t) = (0, v_y(x,t))$ can be decomposed into $\mathbf{v}_{\parallel} = (\mathbf{v} \cdot \mathbf{t})\mathbf{t}$ and $\mathbf{v}_{\perp} = (\mathbf{v} \cdot \mathbf{n})\mathbf{n}$ with

$$\mathbf{v}_{\parallel} = -\frac{A^2 \omega k \cos(kx - \omega t)^2}{1 + A^2 k^2 \cos(kx - \omega t)^2} \mathbf{t}$$
(23)

According to Eq. (15), this generates a force F_x in the swimming direction with

$$F_x = (\zeta_{\parallel} - \zeta_{\perp}) \int dx \frac{A^2 \omega k \cos(kx - \omega t)^2}{1 + A^2 k^2 \cos(kx - \omega t)^2}$$
(24)

while the force in the perpendicular direction vanishes when averaged over the whole flagellum. For small beat amplitudes, Eq. (24) can easily be integrated to give

$$F_x = \frac{1}{2} (\zeta_{\parallel} - \zeta_{\perp}) A^2 \omega k \tag{25}$$

The swimming velocity then follows from $v_x \simeq F_x/\zeta_{\parallel}$ to be

$$v_{flag} = -\frac{1}{2} \left(\frac{\zeta_{\perp}}{\zeta_{\parallel}} - 1 \right) A^2 \omega k .$$
⁽²⁶⁾

This result shows several interesting features. First, it shows that swimming is only possible due to the friction anisotropy, since the velocity is proportional to $(\zeta_{\parallel} - \zeta_{\perp})$. Second, for a traveling wave in the positive *x*-direction, the flagellum moves in the negative *x*-direction, i.e. movement is opposite to the direction of the traveling wave. Third, the swimming velocity increases linearly with the beat frequency ω and the wave vector *k*, but quadratically with the beat amplitude *A*. And finally, the swimming velocity is independent of the fluid viscosity.

A more refined calculation has been performed in Ref. [73], also employing resistive force theory, to determine the swimming velocity of sperm. For the same sinusoidal beat pattern as in Eq. (19) and $\zeta_{\perp}/\zeta_{\parallel} = 2$, this leads to [73]

$$v_{sperm} = \frac{1}{2}A^2\omega k \left[1 + A^2k^2 + \sqrt{1 + \frac{1}{2}A^2k^2} \frac{3r_h}{L_{flag}} (\ln(kd/4\pi) - 1/2) \right]^{-1}$$
(27)

Here L_{flag} is the length of the flagellum, and r_h is the radius of the head. The general conclusions of Eq. (26) remain valid, but additional effects appear. The second term in the brackets of Eq. (27) — whose origin is already recognizable in Eq. (23) — arises from the finite beat amplitude, and implies that the velocity *saturates* for large beat amplitudes A. The last term in the brackets describes the reduction of velocity due the drag of the passive head. The swimming of sperm has also been analyzed by slender-body theory (taking into account the hydrodynamic interactions of different parts of the deformed flagellum, as in Sec.4.1 for slender rods) [74]. Results agree with the resistive-force results of Ref. [73] within about 10% deviation.

An exact solution, taking into account the full hydrodynamics, is possible for an infinitely long flagellum in *two* spatial dimensions (where hydrodynamics is more long-ranged than in three dimensions) — corresponding to an infinite sheet with a propagating lateral wave with transverse oscillations in three dimensions. Here, the swimming velocity is obtained in pioneering work by G.I. Taylor to be [75]

$$v = \frac{1}{2}A^2\omega k \left(1 - \frac{19}{16}A^2k^2\right) \,. \tag{28}$$

This result confirms all qualitative features discussed above, but shows somewhat different numerical coefficients (which is in part due to the different dimensionality).

The sperm structure or beat pattern is typically not completely symmetric, but has some chirality. In this case, sperm swim on helical trajectories [76, 77]. For example, the helicity of the swimming trajectories is very pronounced for sea urchin sperm [78, 79].

The same methods can be employed to calculate the swimming velocity of bacteria with rotating helical flagella [80].

4.3 Sperm Steering and Navigation

For sperm to achieve its main objective, to reach the egg and to deliver its DNA cargo, it has to sense the location of the egg ad redirect its swimming path accordingly. It has been proposed that the flagellum also serves as antenna that registers sensory cues as diverse as chemoattractant



Fig. 14: (top,left) Experimental snapshots of a tethered human sperm that rotates clockwise around the tethering point with rotation velocity $\Omega(t)$. Scale bar represents $5\mu m$ (top,right) Experimental power spectrum of the curvature at positions $15\mu m$ and $25\mu m$ along the flagellum from the tethering point. The fundamental frequency is $\omega_0 = 20Hz$. (bottom,left) Comparison (stroboscopic view) of experimental (red) and simulated (blue) beat pattern. Time interval between snapshots (fading lines) is $\Delta t = 2ms$. (bottom,right) Simulated sperm trajectory resulting from a slowly changing phase shift ψ over time (phase indicated by the color of the trajectory). By modulating the phase, sperm swim on curvilinear paths. From Ref. [84].

molecules, fluid flow, or temperature. The sensory cues modify the flagellar beat pattern and, thereby, guide sperm to the egg. Whereas chemotaxis, the directed movement in a chemical gradient, has been firmly established in sperm from marine invertebrates and plants, it is debated which sensory cues guide mammalian, in particular human sperm to the egg [17]. Two different mechanisms have been proposed that could generate a flagellar asymmetry: a dynamic buckling instability resulting from flagellar compression by internal forces [81,82] or an average intrinsic curvature [33, 77, 83]. A new mechanism has been suggested recently, which relies on the presence of a second harmonic frequency in the beat pattern of sperm [84].

High-speed, high-precision video microscopy of flagellar beat pattern of tethered human sperm (with its head attached to the surface of a planar substrate) indeed reveals that the beat pattern is characterized by a superposition of two bending waves – with a fundamental frequency ω and its second harmonic – traveling down the flagellum, see Fig. 14(top). This leads to a rotation of the sperm around the anchoring point with rotation velocity $\Omega(t)$.

The calculations presented in Sec. 4.2 to determine the swim speed of sperm can be extended to include a second-harmonic contribution. Instead of Eq. (19), we now have [84]

$$y(x,t) = A_1 \sin(kx - \omega t) + A_2 \sin(kx - 2\omega t + \phi)$$
 (29)

with a phase shift ϕ . The calculation of swim forces – based on resistive force theory – then proceeds along the same lines as in Sec. 4.2, which results in the force densities (to leading

order)

$$f_x(x,t) = (\xi_{\perp} - \xi_{\parallel})\partial_t y \partial_x y \tag{30}$$

$$f_y(x,t) = -\xi_\perp \partial_t y + (\xi_\perp - \xi_\parallel) \partial_t y (\partial_x y)^2$$
(31)

The propulsion force F_x remains essentially unchanged due to the (small) 2nd harmonic component. The first term $\xi_{\perp} \partial_t y$ in f_y averages out to zero over one beat period. Thus, the force T_a , which is responsible for a rotation around the tethering point, becomes (in the case that there is one wavelength on the flagellum, i.e. $k = 2\pi/L$) [84]

$$T_a = \frac{\omega}{2\pi} \int_0^{2\pi/\omega} dt \int_0^L dx \ x f_y(x,t)$$
(32)

$$= 2\pi\omega(\xi_{\perp} - \xi_{\parallel})A_{1}^{2}A_{2}\sin(\phi)$$
 (33)

The theoretical prediction of rotation velocity Ω to be proportional to the 2nd-harmonic amplitude $A_2 \sin(\phi)$ indeed agrees very well with the experimental observations [84].

Two other points are worth mentioning. First, the beat shape in Fig. 14(top,right) doesn't look very much like a simple sine wave. However, the experimental beat shape can be reproduced very well, if not the shape is prescribed, as in Eq. (29), but rather active bending torques

$$T(s,t) = T_1 \sin(ks - \omega_0 t) + T_2 \sin(ks - 2\omega_0 t + \psi)$$
(34)

on an elastic, semi-flexible filament. The result beat shapes of simulation and experiment are then hardly distinguishable, see Fig. 14(bottom,left). Second, the variation of the phase shift ψ between the two torque modes in Eq. (34) provides an obvious possibility for steering, as demonstrated by simulations in Fig. 14(bottom,right).

4.4 Bacteria Swimming near Surfaces

In a bulk fluid, the bacteria move in a straight manner (run), with all flagella forming a bundle, interrupted by abrupt changes of the swimming direction (tumble) induced by disintegration of the bundle [29], as discussed in Sec. 1.1. The presence of a surface drastically alters the swimming behaviour. For instance, the non-tumbling mutant of *E. coli* swims in a clockwise (CW) circular trajectory close to a solid boundary [85, 86] and a counterclockwise (CCW) trajectory close to a liquid-air interface [87, 88]. Hence, bacteria are able to "sense" the properties of a nearby surface, an aspect of great importance for surface selection and attachment in the early stages of biofilm formation or infection [85, 89]. Also, a theoretical understanding of hydrodynamic interactions between swimming bacteria and surfaces not only sheds light on selective surface attachment, but opens an avenue for the design of microfluidic devices to control and guide bacterial motion [90] for separation, trapping, stirring, etc. [91].

The swimming behaviour of bacteria near surfaces is governed by hydrodynamic forces [7, 8] and, hence, the CW and CCW circular trajectories of *E. coli* have to be explained in terms of hydrodynamic interactions. The basic physical mechanism is as follows. The rotary motors generate rotations of the flagellar bundle and of the cell body in opposite directions, such that the whole microswimmer experience no net (external) torque. Now, a rotating body near a wall, with the axis of rotation parallel to the wall, experiences two forces induced by the flow, as illustrated in Fig. 15a. The first force is due to the velocity gradient of the fluid in the gap between the wall and the rotating body; it pushes the body in the "rolling" direction. The



Fig. 15: (*left*) Flow and pressure fields, which develop for a rotating body near a planar or curved wall. (*right*) Definition of the Navier length b for simple shear flow near a solid wall.

second force is due to a pressure different between the two sides of the rotating body, which arises from the compression of the fluid into the gap on one side, and the decompression on the other side; it pushes the body against the rolling direction. The actual motion depends on the balance between these two forces. For a no-slip wall, they exactly cancel for an infinitely long cylinder, while the shear force dominates for a sphere (and a short spherocylinder) and generates a motion in the rolling direction. This implies that – in addition to the propulsive forward force – the bacterial body experiences a force to the right, the flagellar bundle to the left (due to the opposite directions of rotation), which results in a CW circling motion close to the wall.

Fluid slip on a surface at z = 0 is characterized by the Navier slip length b, which is defined by the boundary condition

$$\mathbf{v}_s = b \frac{\partial \mathbf{v}}{\partial z} \tag{35}$$

for the fluid flow field $\mathbf{v}(\mathbf{r}, t)$, where \mathbf{v}_s is the fluid velocity at the surface, as illustrated in Fig. 15b. A finite slip length b implies that the fluid velocity gradient in the gap between the rotating body and the wall is reduced, so that the corresponding force contribution is reduced compared to the no-slip case. For a perfect-slip wall, with $b = \infty$, the pressure contribution dominates, and the bacterium performs a CCW motion.

The dependence of the motion on the slip length can be studied numerically by a bacterium model displayed in Fig. 16a. This model consists of a spherocylindrical body and several helical flagella, which are driven by a motor torque at their anchoring points in the cell wall [92, 93]. The bacterium model is then embedded in a particle-based mesoscopic solvent (see Chapter A10 by M. Ripoll) to describe hydrodynamic flows and interactions. Simulations with this model reproduce the two limiting cases of slip lengths very well, see Fig. 16b,d. Obviously, there should be an intermediate value of the slip length, where the bacterium swims on a straight line, which is predicted from the simulations for *E. coli* to occur for b = 30 nm. This is an interesting result, because it shows that the bacterium motion is sensitive to changes in slip length on the tens-of-nanometer scale. Swimming bacteria could therefore act as slip-length sensor on this scale of slip lengths.



Fig. 16: Swimming bacteria sense the slip of its nearby surface. (a) The model bacterium of length ℓ consists of a spherocylindrical body of length ℓ_b and diameter d and four helical flagella each turned by a motor torque. The bacterial geometry and flagellar properties are in agreement with experiments of E. coli. The body and the flagellar bundle counter rotate. h is the gap width between the body and the surface. (b) CW, (c) noisy straight, and (d) CCW trajectories, obtained from hydrodynamic simulations of a bacterium swimming near homogeneous surfaces with different slip lengths b, as indicated. From Ref. [92].

5 Hydrodynamic Synchronization: Metachronal Waves in Cilia Carpets

In nature, the density of flagella, cilia, and various kinds of microswimmers can sometimes be very high. For example, in mammalian reproduction, the average number of *sperm* per ejaculate is tens to hundreds of millions, so that the average distance between sperm is on the scale of ten micrometers — comparable to the length of their flagellum — so that interactions between them are not negligible. In recent years, experiments [94–98] have revealed an interesting cooperative behavior of sperm at high concentration, e.g. the distinctive aggregations or 'trains' of hundreds of wood-mouse sperm [96,97], or the vortex arrays of swimming sea urchin sperm on a substrate [98]. Vortex arrays of curved flagella have also been obtained in simulations [99]. *Cilia* densely cover the surfaces of the protozoa Paramecium [100] and Opalina [101, 102] and of the green algae *Volvox* [45, 103], and line the airways in the human body. Lophotrichous and peritrichous *bacteria* have multiple flagella located at the same spot and located randomly, respectively, on the bacteria surface. Their run-and-tumble motion requires the synchronization, bundling and unbundling of these flagella [104–108]. In all these cases, hydrodynamic interactions play an important role.

Motile cilia on the surface of a cell or microorganism perform an active whip-like motion, which propels the fluid along the surface of cells and tissues. As explained in Sec. 1, this beat consists of a fast power stroke in which the cilium has an elongated shape, and a slower recovery stroke in which the cilium is curved and closer to the cell surface. The beat of different cilia is usually not random, but strongly synchronized in a wave-like pattern which is called a *metachronal wave* (MCW) [101].

Theoretical approaches to investigate hydrodynamic interactions between cilia and the formation of metachronal waves fall in three categories: (i) highly simplified model systems, designed to elucidate the mechanism of hydrodynamic synchronization of many active agents [11, 103, 109–112], (ii) models of an actively driven semi-flexible filament, which mimic the beat of a real cilium [113–117], and (iii) models of a filament, with a beat shape obtained from



Fig. 17: Simulation snapshot of cilia conformations during the formation of metachronal waves. From Ref. [117].

maximizing the pumping efficiency [118, 119].

The second class of models consists of semi-flexible filaments, which are deformed actively by internal forces to reproduce the power and recovery strokes of real cilia, and can react to the flow field generated by their neighbors [113–116]. The studies of such models have been restricted so far mainly to effectively one-dimensional chains of cilia [113–116], or to two-dimensional arrays of a small number of cilia [120]. They provide an indication of metachronal coordination, but the considered systems were too small or the time evolution too short to allow any prediction of MCW properties. This limitation has been overcome very recently by a large-scale, mesoscopic hydrodynamics simulation of arrays of up to 60×60 cilia, which are modeled as active filaments with geometric triggers for the switching between power and recovery strokes [117], see Fig. 17. This model provides clear evidence for the emergence of metachronal waves, and allows a detailed study of wave and transport properties [117].

Simulation results for the beat period, the fluid velocity above the cilia array, and the transport efficiency are compared in Fig. 18 for cilia arrays with either metachronal coordination or a completely synchronous beat [117]. The first surprising result is that the time required for a single beat actually increases in the presence of metachronal waves (MCWs), see Fig. 18(a). This can be understood be the larger hydrodynamic resistance when cilia are partially beating against each other in a MCW. However, the faster beating in synchronous beating is to little avail, since the fluid above is just pushed back and forth, which results in a much smaller net fluid transport velocity compared to MCWs, see Fig. 18(b). This is not only inefficient in the sense of generating only slow fluid transport, but also from an energetic point of view, because also the energy efficiency

$$\epsilon = \frac{16}{3} \frac{\eta v^2 d_c^2}{L_c P_c} \tag{36}$$

(where η is the fluid viscosity, v the average fluid velocity, d_c and L_c the cilium separation and length, respectively, and P_c the energy consumption per unit time) is much lower for synchronous beating (because also oscillating a viscous fluid generates dissipation) than for MCWs. This effect can be quite large, reaching almost an order of magnitude for the optimal cilia spacing.



Fig. 18: Mesoscale hydrodynamics simulations show that metachronal coordination greatly increases swimming velocity. (a) Beat period τ_b versus cilia spacing d_c . (b) Average fluid velocity v. Solid lines are fits to power laws $d_c^{-\nu}$ (using data for $d_c \ge 5$ with 20×20 cilia). In all figures, results are shown for MCWs in arrays of 20×20 cilia (red bullets) and in arrays of 60×60 cilia (red circles), as well as for synchronously beating cilia (blue squares). Error bars denote standard deviation. From Ref. [117].

From the experimental side, a detailed quantitative investigation of metachronal waves and their properties has been performed recently in Ref. [103] for *Volvox carteri*, whose large size and ease of visualization make it an ideal model organism for such studies. In particular, beat frequencies, decay times, and wave numbers have been determined. Interestingly, the decay time is found to be just a few times the beat period.

6 Summary and Conclusions

The swimming and crawling motility of biological cells and micro-organisms is an exciting field of current research. Interesting aspects of motility are (i) the prediction of individual particle motion from the underlying propulsion mechanism, (ii) the elucidation of the dynamics of two self-propelled particles which interact with each other via static (conservative) and hydrodynamic forces, (iii) the determination of the swimming or crawling behavior n complex environments, near walls and in confined geometry, and (iv) the understanding of the collective behavior of many interacting cells and micro-organisms.

The understanding of the dynamical behavior of many "classical" biological microswimmers (such as sperm, bacteria and algae) has made much progress in recent years. This opens up the way for addressing more intricate questions. What is universal about the swimming behavior, and what is specific for a certain class of swimmers (like the synchronization of different wave forms of flagella)? How can artificial microswimmers be designed to react to external stimuli and find their targets? How can motility be controlled? How do motile cells affect each other in collective migration? How can the biological complexity be incorporated in theoretical models? The investigations of these and other questions make cell motility an exciting research topic also in the future.

Appendices

A Oseen Tensor of Low-Reynolds-Number Hydrodynamics

The Stokes equation

$$-\nabla p(\mathbf{r}) + \eta \nabla^2 \mathbf{u}(\mathbf{r}) + \mathbf{f}(\mathbf{r}) = 0$$
(37)

for the hydrodynamic velocity field $\mathbf{u}(\mathbf{r})$ and the pressure field $p(\mathbf{r})$ of an incompressible fluid [with $\nabla \cdot \mathbf{u}(\mathbf{r}) = 0$] can be solved explicitly for a body-force field $\mathbf{f}(\mathbf{r})$ by Fourier transformation. With

$$p(\mathbf{k}) = \int d^3 r e^{i\mathbf{k}\cdot\mathbf{r}} p(\mathbf{r})$$
(38)

and similarly u(k), etc., the Stokes equation and the incompressibility condition become

$$-i\mathbf{k}p(\mathbf{k}) - \eta k^2 \mathbf{u}(\mathbf{k}) + \mathbf{f}(\mathbf{k}) = 0$$
(39)

$$\mathbf{k} \cdot \mathbf{u}(\mathbf{k}) = 0. \tag{40}$$

The multiplication of Eq. (39) with k then implies

$$p(\mathbf{k}) = -i\mathbf{k} \cdot \mathbf{f}(\mathbf{k})/k^2 \,. \tag{41}$$

Insertion of this result in Eq. (39) gives

$$-\mathbf{k}(\mathbf{k} \cdot \mathbf{f}(\mathbf{k}))/k^2 - \eta k^2 \mathbf{u}(\mathbf{k}) + \mathbf{f}(\mathbf{k}) = 0$$
(42)

so that

$$u_{\alpha}(\mathbf{k}) = \frac{1}{\eta k^2} \sum_{\beta} \left[\delta_{\alpha\beta} - \frac{k_{\alpha}k_{\beta}}{k^2} \right] f_{\beta}(\mathbf{k}) .$$
(43)

Fourier transformation back to real space then yields (convolution theorem)

$$\mathbf{u}(\mathbf{r}) = \int d^3 r' \mathbf{H}(\mathbf{r} - \mathbf{r}') \cdot \mathbf{f}(\mathbf{r}')$$
(44)

with the Oseen tensor

$$H_{\alpha\beta}(\mathbf{r}) = \int \frac{d^3k}{(2\pi)^3} e^{-i\mathbf{k}\cdot\mathbf{r}} \frac{1}{\eta k^2} \left[\delta_{\alpha\beta} - \frac{k_{\alpha}k_{\beta}}{k^2} \right] \,. \tag{45}$$

Because the tensor H(r) only depends on the vector r, it can be written in the form

$$H_{\alpha\beta}(\mathbf{r}) = A(r)\delta_{\alpha\beta} + B(r)r_{\alpha}r_{\beta}/r^2 ; \qquad (46)$$

Thus,

$$\sum_{\alpha} H_{\alpha\alpha}(\mathbf{r}) = 3A(r) + B(r) \tag{47}$$

$$\sum_{\alpha\beta} H_{\alpha\beta}(\mathbf{r}) r_{\alpha} r_{\beta} / r^2 = A(r) + B(r) .$$
(48)



Fig. 19: Flow lines of hydrodynamic monopole and dipole, oriented in the horizontal direction. (top) Far-field of the monopole, as given by Eq. (51). (bottom, left) Near-field of the dipole, as given by Eq. (53). The two force centers are marked by (red) bullets. (bottom, right) Far-field of the dipole, as given by Eq. (54). The separatrix between the inflow and outflow regions is shown by thick red lines.

An explicit calculation using Eq. (45) provides

$$\sum_{\alpha} H_{\alpha\alpha}(\mathbf{r}) = \frac{1}{2\pi\eta r} \quad ; \quad \sum_{\alpha\beta} H_{\alpha\beta}(\mathbf{r}) r_{\alpha} r_{\beta} / r^2 = \frac{1}{4\pi\eta r} \,. \tag{49}$$

Thus, the Oseen tensor finally reads

$$H_{\alpha\beta}(\mathbf{r}) = \frac{1}{8\pi\eta r} \left[\delta_{\alpha\beta} + \frac{r_{\alpha}r_{\beta}}{r^2} \right] \,. \tag{50}$$

The Oseen tensor shows that the hydrodynamic interaction is very long ranged, with a 1/r decay like the Coulomb potential. Equations (44) and (50) imply that the fluid velocity field for a point force at the origin in direction $\hat{\mathbf{e}}$, $\mathbf{f}(\mathbf{r}) = \hat{\mathbf{e}}\delta(\mathbf{r})$, is given by

$$\mathbf{u}(\mathbf{r}) = \frac{1}{8\pi\eta r} \left[\hat{\mathbf{e}} + \frac{(\mathbf{r} \cdot \hat{\mathbf{e}})\mathbf{r}}{r^2} \right] \,. \tag{51}$$

The streamlines of a force monopole is shown in Fig. 19(a)

B Hydrodynamic Force Dipoles

The flow field of a hydrodynamic force dipole, with opposite forces of equal magnitude at $\mathbf{r} = \pm \mathbf{r}_0$ and direction $\hat{\mathbf{e}} = \mathbf{r}_0/r_0$,

$$\mathbf{f}_1(\mathbf{r}) = f_0 \hat{\mathbf{e}} \delta(\mathbf{r} - \mathbf{r}_0) \quad ; \quad \mathbf{f}_2(\mathbf{r}) = -f_0 \hat{\mathbf{e}} \delta(\mathbf{r} + \mathbf{r}_0) \tag{52}$$

is obtained from Eq. (51) to be

$$\mathbf{u}(\mathbf{r}) = \frac{1}{8\pi\eta\sqrt{(\mathbf{r}-\mathbf{r}_0)^2}} \left[\hat{\mathbf{e}} + \frac{((\mathbf{r}-\mathbf{r}_0)\cdot\hat{\mathbf{e}})(\mathbf{r}-\mathbf{r}_0)}{(\mathbf{r}-\mathbf{r}_0)^2} \right] - \frac{1}{8\pi\eta\sqrt{(\mathbf{r}+\mathbf{r}_0)^2}} \left[\hat{\mathbf{e}} + \frac{((\mathbf{r}+\mathbf{r}_0)\cdot\hat{\mathbf{e}})(\mathbf{r}+\mathbf{r}_0)}{(\mathbf{r}+\mathbf{r}_0)^2} \right]$$
(53)

The flow lines in the near-field of this force dipole are shown in Fig. 19(b). An expansion to leading order in $\mathbf{r}_0/|\mathbf{r}|$ yields

$$\mathbf{u}(\mathbf{r}) = \frac{P}{8\pi\eta r^3} \left[-1 + 3\frac{(\mathbf{r}\cdot\hat{\mathbf{e}})^2}{r^2} \right] \mathbf{r}$$
(54)

where $P = 2f_0r_0$ is the dipole strength. Note that the flow field of a force dipole decays as $1/r^2$ from the center of the dipole, faster than the force monopole of Eq. (51). The flow lines of a hydrodynamic dipole are shown in Fig. 19(c). There are two inflow and two outflow regions in the *xy*-projection, which are determined by the separatrix $y = \pm \sqrt{2}x$. In three dimensions, the outflow region is of course a cone.

C Hydrodynamics near Planar Surfaces

An exact solution of the Stokes equation near a planar surface with no-slip boundary conditions (i.e. $\mathbf{u}(\mathbf{r}) = 0$ at the surface) is possible. The solution is called the Blake tensor [121]. Instead of the exact solution, we consider here a simple approximation for a dipole near a planar wall, oriented parallel to the wall, by employing the method of image charges known from electrostatics, i.e. we write for a wall at z = 0,

$$\mathbf{u}_{wall}(\mathbf{r} - \mathbf{r}_0) = \mathbf{u}_{dipole}(\mathbf{r} - \mathbf{r}_0; \hat{\mathbf{e}}) + \mathbf{u}_{dipole}(\mathbf{r} - \mathbf{r}_1; \hat{\mathbf{e}}')$$
(55)

where $\mathbf{r}_0 = (x_0, y_0, z_0)$ and $\mathbf{r}_1 = (x_0, y_0, -z_0)$, with $z_0 > 0$, and $\hat{\mathbf{e}}'$ is the mirror image of $\hat{\mathbf{e}}$ with respect to the z = 0 plane. This implies that at z = 0 the velocity field u_z perpendicular to the surface vanishes identically, $u_z \equiv 0$, but the no-slip boundary condition is *not* satisfied. The dipole experiences a force near the surface, which is determined by the hydrodynamic interaction between the dipole and the image charge. It is given by the z-component of the flow field of the image charge at the location of the dipole, so that

$$v_z(z_0) = -\frac{P}{32\pi\eta z_0^2} \left[1 - 3(\hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_z)^2 \right] \,.$$
(56)

because $(\hat{\mathbf{e}}' \cdot \hat{\mathbf{e}}_z)^2 = (\hat{\mathbf{e}} \cdot \hat{\mathbf{e}}_z)^2$. This result shows that the hydrodynamic force is attractive to the wall, and that it decays as the dipole flow field with the squared distance from the wall. Eq. (56) is in agreement with the exact result (10) for no-slip boundary conditions, except that the numerical prefactor in Eq. (56) is smaller by a factor 2/3.

References

- [1] E. M. Purcell, Am. J. Phys. 45, 3 (1977).
- [2] J. Toner, Y. Tu, and S. Ramaswamy, Ann. Phys. (N.Y.) 318, 170 (2005).
- [3] S. Ramaswamy, Annu. Rev. Condensed Matter Phys. 1(1), 323 (2010).
- [4] T. Vicsek and A. Zafeiris, Phys. Rep. 517, 71 (2012).
- [5] M. C. Marchetti, J. F. Joanny, S. Ramaswamy, T. B. Liverpool, J. Prost, M. Rao, and R. A. Simha, Rev. Mod. Phys. 85, 1143 (2013).
- [6] D. Saintillan and M. J. Shelley, C. R. Physique 14, 497 (2013).
- [7] J. Elgeti, R. G. Winkler, and G. Gompper, Rep. Prog. Phys. 78, 056601 (2015).
- [8] E. Lauga and T. R. Powers, Rep. Prog. Phys. 72, 096601 (2009).
- [9] T. Ishikawa, J. Royal Soc. Interface 6, 815 (2009).
- [10] D. L. Koch and G. Subramanian, Annu. Rev. Fluid Mech. 43, 637 (2011).
- [11] R. Golestanian, J. M. Yeomans, and N. Uchida, Soft Matter 7, 3074 (2011).
- [12] J. Elgeti and G. Gompper, Eur. Phys. J. Special Topics 225, 2333 (2016).
- [13] R. M. Harshey, Annu. Rev. Microbiol. 57, 249 (2003).
- [14] M. E. Cates, Rep. Prog. Phys. 75, 042601 (2012).
- [15] E. Lauga, Annu. Rev. Fluid Mech. 48, 105 (2016).
- [16] E. A. Gaffney, H. Gadelha, D. J. Smith, J. R. Blake, and J. C. Kirkman-Brown, Annu. Rev. Fluid Mech. 43, 501 (2011).
- [17] L. Alvarez, B. M. Friedrich, G. Gompper, and U. B. Kaupp, Trends Cell Biol. 24, 198 (2014).
- [18] A. Zöttl and H. Stark, J. Phys.: Condens. Matter 28, 253001 (2016).
- [19] C. Bechinger, R. Di Leonardo, H. Löwen, C. Reichhardt, G. Volpe, and G. Volpe, Rev. Mod. Phys. 88, i045006 (2016).
- [20] P. Romanczuk, M. Bär, W. Ebeling, B. Lindner, and L. Schimansky-Geier, Eur. Phys. J. Special Topics 1, 202 (2012).
- [21] M. E. Cates and J. Tailleur, Annu. Rev. Condens. Matter Phys. 6, 219 (2015).
- [22] G. A. Ozin, I. Manners, S. Fournier-Bidoz, and A. Arsenault, Adv. Mater. 17, 3011 (2005).
- [23] S. Sengupta, M. E. Ibele, and A. Sen, Angew. Chem. Int. Ed. 51, 8434 (2012).
- [24] P. Friedl and D. Gilmour, Nat. Rev. Mol. Cell Biol. 10, 445 (2009).
- [25] P. Sens and J. Plastino, J. Phys.: Condens. Matter 27, 273103 (2015).
- [26] I. S. Aronson, Physical Models of Cell Motility (Springer, New York, 2016).
- [27] C. Brennen and H. Winet, Annu. Rev. Fluid Mech. 9, 339 (1977).
- [28] D. Jones, New Scientist 197, 40 (2008).
- [29] H. C. Berg, E. coli in Motion (Springer, New York, 2004).
- [30] H. C. Berg, Annu. Rev. Biochem. 72, 19 (2003).
- [31] N. C. Darnton, L. Turner, S. Rojevsky, and H. C. Berg, J. Bacteriol. 189, 1756 (2007).
- [32] Http://www.tutorvista.com/topic/cilia-and-flagella.
- [33] B. M. Friedrich, I. H. Riedel-Kruse, J. Howard, and F. Jülicher, J. Exp. Biol. 213, 1226 (2010).
- [34] S. L. Tamm and G. A. Horridge, Proc. R. Soc. Lond. B 175, 219 (1970).
- [35] D. R. Brumley, K. Y. Wan, M. Polin, and R. E. Goldstein, eLife 3, e02750 (2014).
- [36] C. B. Lindemann and K. A. Lesich, J. Cell Sci. 123, 519 (2010).
- [37] C. B. Lindemann, Biophys. J. 107, 1487 (2014).
- [38] F. Jülicher and J. Prost, Phys. Rev. Lett. 78, 4510 (1997).

- [39] S. Camalet, F. Jülicher, and J. Prost, Phys. Rev. Lett. 82, 1590 (1999).
- [40] Http://en.wikipedia.org/wiki/Chlamydomonas_reinhardtii.
- [41] K. Drescher, R. E. Goldstein, N. Michel, M. Polin, and I. Tuval, Phys. Rev. Lett. 105, 168101 (2010).
- [42] J. S. Guasto, K. A. Johnson, and J. P. Gollub, Phys. Rev. Lett. 105, 168102 (2010).
- [43] Http://remf.dartmouth.edu/images/algaeSEM.
- [44] Http://en.wikipedia.org/wiki/Volvox.
- [45] K. Drescher, K. C. Leptos, I. Tuval, T. Ishikawa, T. J. Pedley, and R. E. Goldstein, Phys. Rev. Lett. 102, 168101 (2009).
- [46] J. R. Howse, R. A. L. Jones, A. J. Ryan, T. Gough, R. Vafabakhsh, and R. Golestanian, Phys. Rev. Lett. 99, 048102 (2007).
- [47] R. Dreyfus, J. Baudry, M. L. Roper, M. Fermigier, H. A. Stone, and J. Bibette, Nature 437, 862 (2005).
- [48] B. J. Williams, S. V. Anand, J. Rajagopalan, and M. T. A. Saif, Nat. Commun. 5, 3081 (2014).
- [49] Y. W. Kim and R. R. Netz, Phys. Rev. Lett. 96, 158101 (2006).
- [50] E. M. Gauger, M. T. Downton, and H. Stark, Eur. Phys. J. E 28, 231 (2009).
- [51] T. Sanchez, D. Welch, D. Nicastro, and Z. Dogic, Science 333, 456 (2011).
- [52] S. Sareh, J. Rossiter, A. Conn, K. Drescher, and R. E. Goldstein, J. Royal Soc. Interface 10, 20120666 (2012).
- [53] A. Keißner and C. Brücker, Soft Matter 8, 5342 (2012).
- [54] J. Branscomb and A. Alexeev, Soft Matter **6**, 4066 (2010).
- [55] M. Vilfan, A. Potočnik, B. Kavčič, N. Osterman, I. Poberaj, A. Vilfan, and D. Babič, Proc. Natl. Acad. Sci. USA 107, 1844 (2010).
- [56] D. A. Fletcher and R. D. Mullins, Nature 463, 485 (2010).
- [57] R. Ananthakrishnan and A. Ehrlicher, Int. J. Biol. Sci. 3, 303 (2007).
- [58] A. Mogilner and G. Oster, Biophys. J. 71, 3030 (1996).
- [59] J. Weichsel and U. S. Schwarz, Proc. Natl. Acad. Sci. USA 107, 6304 (2010).
- [60] A. Mogilner and K. Keren, Curr. Biol. 19, R762 (2009).
- [61] D. Shao, H. Levine, and W.-J. Rappel, Proc. Natl. Acad. Sci. USA 109, 6851 (2012).
- [62] R. Mayor and C. Carmona-Fontaine, Trends Cell Biol. 20, 319 (2010).
- [63] B. A. Camley, Y. Zhao, B. Li, H. Levine, and W.-J. Rappel, Phys. Rev. Lett. 111, 158102 (2013).
- [64] F. Ziebert and I. S. Aranson, npj Comput. Mater. 2, 16019 (2016).
- [65] M. Paoluzzi, R. D. Leonardo, M. C. Marchetti, and L. Angelani, Sci. Rep. 6, 34146 (2016).
- [66] C. A. Velasco, S. D. Ghahnaviyeh, H. N. Pishkenari, T. Auth, and G. Gompper, Soft Matter 13, 5865 (2017).
- [67] M. Skruzny, T. Brach, R. Ciuffa, S. Rybina, M. Wachsmuth, and M. Kaksonen, Proc. Natl. Acad. Sci. USA 109, E2533 (2012).
- [68] K. Drescher, J. Dunkel, L. H. Cisneros, S. Ganguly, and R. E. Goldstein, Proc. Natl. Acad. Sci. USA 108, 10940 (2011).
- [69] A. P. Berke, L. Turner, H. C. Berg, and E. Lauga, Phys. Rev. Lett. 101, 038102 (2008).
- [70] I. O. Götze and G. Gompper, Phys. Rev. E 82, 041921 (2010).
- [71] M. Doi and S. F. Edwards, *The Theory of Polymer Dynamics* (Clarendon Press, Oxford, 1986).
- [72] J. K. G. Dhont, An Introduction to Dynamics of Colloids (Elsevier, Amsterdam, 1996).

- [73] J. Gray and G. J. Hancock, J. Exp. Biol. 32, 802 (1955).
- [74] J. J. L. Higdon, J. Fluid Mech. 90, 685 (1979).
- [75] G. I. Taylor, Proc. Roy. Soc. Lond. A 209, 447 (1951).
- [76] H. C. Crenshaw, Biophys. J. 56, 1029 (1989).
- [77] J. Elgeti, U. B. Kaupp, and G. Gompper, Biophys. J. 99, 1018 (2010).
- [78] U. B. Kaupp, J. Solzin, E. Hildebrand, J. E. Brown, A. Helbig, V. Hagen, M. Beyermann, F. Pampaloni, and I. Weyand, Nat. Cell Biol. 5, 109 (2003).
- [79] J. F. Jikeli, L. Alvarez, B. M. Friedrich, L. G. Wilson, R. Pascal, R. Colin, M. Pichlo, A. Rennhack, C. Brenker, and U. B. Kaupp, Nat. Comm. 6, 7985 (2015).
- [80] J. J. L. Higdon, J. Fluid Mech. 94, 331 (1979).
- [81] H. Gadelha, E. A. Gaffney, D. J. Smith, and J. C. Kirkman-Brown, J. Royal Soc. Interface 7, 1689 (2010).
- [82] A. Bukatin, I. Kukhtevich, N. Stoop, J. Dunkel, and V. Kantsler, Proc. Natl. Acad. Sci. USA 112, 15904 (2015).
- [83] V. F. Geyer, P. Sartori, B. M. Friedrich, F. Jülicher, and J. Howard, Curr. Biol. 26, 1098 (2016).
- [84] G. Saggiorato, L. Alvarez, J. F. Jikeli, U. B. Kaupp, G. Gompper, and J. Elgeti, Nat. Commun. 8, 1415 (2017).
- [85] E. Lauga, W. R. DiLuzio, G. M. Whitesides, and H. A. Stone, Biophys. J. 90, 400 (2006).
- [86] L. Lemelle, J.-F. Palierne, E. Chatre, C. Vaillant, and C. Place, Soft Matter 9, 9759 (2013).
- [87] L. Lemelle, J. Palierne, E. Chatre, and C. Place, J. Bacteriol. 192, 6307 (2010).
- [88] R. Di Leonardo, D. Dell'Arciprete, L. Angelani, and V. Iebba, Phys. Rev. Lett. 106, 038101 (2011).
- [89] L. A. Pratt and R. R. Kolter, Mol. Microbiol. 30, 285 (1998).
- [90] W. R. DiLuzio, L. Turner, M. Mayer, P. Garstecki, D. B. Weibel, H. C. Berg, and G. M. Whitesides, Nature 435, 1271 (2005).
- [91] Podcast "Welt der Physik, Folge 189 Mikroschwimmer im Modell" (09. Juli 2015); http://www.weltderphysik.de/mediathek/podcast/.
- [92] J. Hu, A. Wysocki, R. G. Winkler, and G. Gompper, Sci. Rep. 5, 9586 (2015).
- [93] J. Hu, M. Yang, G. Gompper, and R. G. Winkler, Soft Matter 11, 7867 (2015).
- [94] H. D. M. Moore and D. A. Taggart, Biol. Reprod. 52, 947 (1995).
- [95] F. Hayashi, Funct. Ecol. 12, 347 (1998).
- [96] S. Immler, H. D. M. Moore, W. G. Breed, and T. R. Birkhead, PLoS ONE 2, e170 (2007).
- [97] H. D. M. Moore, K. Dvoráková, N. Jenkins, and W. G. Breed, Nature 418, 174 (2002).
- [98] I. H. Riedel, K. Kruse, and J. Howard, Science 309, 300 (2005).
- [99] Y. Yang, F. Qiu, and G. Gompper, Phys. Rev. E 89, 012720 (2014).
- [100] H. Machemer, J. Exp. Biol. 57, 239 (1972).
- [101] M. Sleigh, 2. The Biology of Cilia and Flagella (Pergamon Press, Oxford, 1962), chap. The Structure of Cilia, pp. 11–75.
- [102] R. Rikmenspoel, Biophys. J. 16, 445 (1976).
- [103] D. R. Brumley, M. Polin, T. J. Pedley, and R. E. Goldstein, Phys. Rev. Lett. 109, 268102 (2012).
- [104] M. J. Kim, J. C. Bird, A. J. Van Parys, K. S. Breuer, and T. R. Powers, Proc. Natl. Acad. Sci. USA 100, 15481 (2003).
- [105] M. Reichert and H. Stark, Eur. Phys. J. E 17, 493 (2005).
- [106] B. Qian, H. Jiang, D. A. Gagnon, K. S. Breuer, and T. R. Powers, Phys. Rev. E 80, 061919

(2009).

- [107] S. Y. Reigh, R. G. Winkler, and G. Gompper, Soft Matter 8, 4363 (2012).
- [108] S. Y. Reigh, R. G. Winkler, and G. Gompper, PLoS ONE 8, e70868 (2013).
- [109] A. Vilfan and F. Jülicher, Phys. Rev. Lett. 96, 058102 (2006).
- [110] T. Niedermayer, B. Eckhardt, and P. Lenz, Chaos 18(3), 037128 (2008).
- [111] N. Uchida and R. Golestanian, Phys. Rev. Lett. 104, 178103 (2010).
- [112] C. Wollin and H. Stark, Eur. Phys. J. E 34, 42 (2011).
- [113] S. Gueron, K. Levit-Gurevich, N. Liron, and J. J. Blum, Proc. Natl. Acad. Sci. USA 94(12), 6001 (1997).
- [114] S. Gueron and K. Levit-Gurevich, Proc. Nat. Acad. Sci. USA 96(22), 12240 (1999).
- [115] B. Guirao and J. Joanny, Biophys. J. 92, 1900 (2007).
- [116] Y. Yang, J. Elgeti, and G. Gompper, Phys. Rev. E 78, 061903 (2008).
- [117] J. Elgeti and G. Gompper, Proc. Natl. Acad. Sci. USA 110, 4470 (2013).
- [118] N. Osterman and A. Vilfan, Proc. Natl. Acad. Sci. USA 108, 15727 (2011).
- [119] C. Eloy and E. Lauga, Phys. Rev. Lett. 109, 038101 (2012).
- [120] S. Gueron and K. Levit-Gurevich, Proc. Biol. Sci 268(1467), 599 (2001).
- [121] J. R. Blake, Proc. Camb. Philos. Soc. 70, 303 (1971).

E 5 Bacterial biofilms – Physical determinants of microbial architecture

Berenike Maier Experimental Biophysics Institute of Theoretical Physics Universität zu Köln

Contents

1	Intr	oduction	2	
2	A m	olecular toolbox for controlling bacterial adhesion	2	
	2.1	Initial attachment to surfaces	2	
	2.2	Bacterial appendages are involved in surface sensing	4	
3	Interactions between bacteria and local order		4	
	3.1	Liquid-like order of spherical bacteria	4	
	3.2	Fluid crystalline-like order of rod-shaped bacteria	5	
	3.3	The role of the extracellular matrix on biofilm structure	6	
4	Physical mechanisms of cellular sorting		6	
	4.1	Passive spatial segregation in biofilms	6	
	4.2	Sorting with respect to differential interaction force	6	
	4.3	Sorting with respect to cell shape	8	
	4.4	Effects of the extracellular matrix on segregation and sorting		
5	Bacterial fitness relates to biofilm structure		8	
	5.1	Differential adhesion to surfaces affects bacterial fitness	9	
	5.2	Differential bacterial interactions affect bacterial fitness	9	
	5.3	Cell shape affects bacterial fitness	9	
6	Con	Conclusion		
A	For	Force spectroscopy for characterizing bacterial interaction forces 11		

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Biofilms are structured communities formed by a single or multiple microbial species. Within biofilms, bacteria are embedded into extracellular matrix allowing them to build macroscopic objects. It is becoming clear that life in biofilms is major mode of microbial life [1]. Since many persistent and untreatable bacterial infections are associated with biofilm formation, studying biofilms is of high societal relevance. Many bacterial species can switch from life as an individual to life in biofilms and vice versa. Within biofilms, bacteria differentiate creating spatial heterogeneity. Importantly, bacteria can adjust the structure of the biofilm in response to changes in the environment. For example, some bacterial species generate large biofilm structures in the presence of antibiotics (Fig. 1). It is well established that the biofilm structure responds to environmental changes, but it is also clear that bacteria can only tune their interactions at the level of individual cells. Currently, we know very little about the link between physical interactions at the level of single cells and biofilm structure. In this lecture, we will discuss recent progress in quantifying cell-cell and cell-surface interactions at the molecular level. Subsequently, we will introduce bacterial systems where single cell properties were correlated with the structure of biofilms. Finally, we will relate structure to bacterial fitness. In this chapter all structured bacterial communities including colonies grown in liquid or on agar plates are denoted as biofilms.



Fig. 1: Biofilm formation in response to external stress. Staphylococcus aureus was grown in the absence (left) or in the presence (right) of the antibiotic methicillin. Biofilm formation can be seen as a pellet in the right test tube (image adapted from [2]).

2 A molecular toolbox for controlling bacterial adhesion

2.1 Initial attachment to surfaces

Bacteria have evolved numerous cell surface appendages and extracellular matrix components that allow controlling the materials properties of biofilms. Many of these extracellular structures are polymers. Intuitively, it is clear that by tuning interaction forces with surfaces and other bacteria, bacteria can control the structure of biofilms. Still, we know very little about the correlation between these interaction forces and biofilm structure. Using force spectroscopy (Appendix A), a few research groups have started characterizing the impact of surface appendages on interaction forces in a systematic way. In the following, we will describe mechanical properties of various cell appendages.

Pili are ubiquitous filamentous surface appendages that serve multiple functions. Here, we will shortly introduce the best-characterized pilus types, namely type I pili (T1P), type IV pili (T4P), and P-pili [3]. They all serve the function of mediating bacterial attachment to the surface of other cells and to abiotic surfaces. All these types of pili are polymerized from pilin subunits that are secreted from the cytoplasm to the exterior and they are anchored to the bacterial cell envelope via membrane spanning complexes. The pilin subunits are specific to the pilus type but within a specific type they are well conserved between different bacterial species and even archea [3]. There are no known homologs in eucaryotes. During life at surfaces, bacteria face a fundamental problem. On the one hand, they need to attach firmly in order to avoid being flushed away by flow. On the other hand, during an immotile lifestyle, they rapidly deplete nutrients and they are incapable of colonizing the surface. Thus, bacteria have invented various mechanical 'tricks' that endow them with motility while being bound to the surface.

Type IV pili (T4P) mediate surface motility enabling bacteria to colonize surfaces while remaining attached [4]. When a pilus polymer is attached to a surface, it can retract by depolymerization (Fig. 2a) [5]. During retraction, it exerts force onto the cell body like a grappling hook (Fig. S1a) [6,7]. Interestingly, for *Neisseria gonorrhoeae*, the forces generated during pilus retraction exceed 100 pN making the T4P one of the strongest molecular motors characterized so far [8]. Some bacterial species can actively modify the interaction strength between the T4P and the surface by secreting exopolysaccharides [9]. As a consequence, bacteria tend to follow their own trail on the surface [10]. Trail formation impacts on the structure of the biofilm.



Fig. 2: Extracellular polymers induce surface attachement. a) The type IV pilus (T4P) mediates attachment to the surface. As the pilus retracts by depolymerization, the cell body approaches the surface. By cycles of attachment, retraction, and elongation, bacteria move over surfaces. b) The type I pilus (T1P) shows optimal properties for mediating adhesion in fluid flow. In the absence of flow, the probability that the ligand at the tip binds to a specific receptor is low. Under flow, the ligand undergoes a conformational change increasing bond stability. Simultaneously, the shaft of the T1P uncoils modulating the force acting on the tip to increase bond stability.

Type I pili (T1P) of *Escherichia coli* have evolved mechanical properties to optimize the stability of the tip adhesive under flow (Fig. 2b) [11]. The tip adhesive forms a catch-bond with its receptor, i.e. the interaction is reinforced by mechanical stress. As a consequence, the binding period of *E. coli* is maximal when forces in the range between (30 - 60) pN are applied. Interestingly, within the same force regime, the T1P uncoils increasing the length of the pilus seven-fold and acting as a "buffer" to keep the force in that regime [11] (Fig. S1a,c). Other pili including T4P pili and P-pili show similar elongation transitions of the pilus shaft [12,13].

Flagella are helical surface appendages that are connected to a rotary motor located within the cell envelope [14]. Flagella rotation enables bacteria to swim through acqueous environment and to swarm at humid surfaces. Only recently it was discovered that flagella play an important role during the first steps of biofilm formation at surfaces as described in 2.2. Once bacteria are committed to biofilm formation, flagella rotation tends to stop [15].

2.2 Bacterial appendages are involved in surface sensing

When bacteria switch from life as individual cells to life in biofilms, they change their patterns of gene expression. For this, bacteria must sense that they are in close contact to a surface. Recently, it was shown that T4P retraction is important for surface sensing in *Pseudomonas aeruginosa* [16,17]. The strength of the gene expression response was dependent of the stiffness of the surface. For a different species, *Caulobacter crescentus*, two different mechanisms of surface sensing were reported. One mechanism relies on T4P, in particular cells responded to blocking T4P retraction [18]. Another group found that the flagella motor was required for surface detection [19]. Even though the details of surface sensing remain far from being well understood, it becomes clear that motor properties of bacterial surface appendices play a crucial role in surface sensing.

3 Interactions between bacteria and local order

Recent advances in fluorescence microscopy and image analysis have enabled imaging of bacterial colonies at the single cell level. Thus, it became possible to address the question whether bacteria show order within biofilms. Bacterial cells are often in close contact, making mechanical interactions between adjacent bacteria particularly significant. If we think of bacteria as hard spheres or ellipsoids whose attractive interactions are controlled by the surface appendages described in chapter 2, then we expect that they behave like colloidal systems assuming crystalline, fluid crystalline, or liquid-like order. In contrast to colloids, however, bacteria reproduce. Moreover, some of their surface structures (including type IV pili and flagella) act as molecular motors as well as adhesins.

3.1 Liquid-like order of spherical bacteria

We have addressed the question whether biofilms formed by spherical bacteria (cocci) show local order. For the species *Neisseria gonorrhoeae*, T4P are the dominating surface structure in terms of mediating attractive interactions between bacteria [20]. We found that *N. gonorrhoeae* show short-range liquid-like order (Fig. 3a) [21]. In particular, the radial g(r) showed two maxima in early biofilms and became more prounounced in 24 h biofilms. g(r) is

defined so that $N/V g(r)r^2 dr$ is the probability that the center of a bacterium will be found in the range dr at a distance r from the center of another bacterium, where N/V is the number density of bacteria. Interestingly, fact that T4P are molecular motors is important for local order and mesoscopic structure of the biofilm. Strains that generate retraction-deficient T4P form biofilms without short-range order and lower density. We propose that bacteria continuously pull on each other in a tug-of-war like fashion to tune the structure of the biofilm.



Fig. 3: Confocal reconstructions of biofilms imaged with single cell resolution. a) Spherical cells (N. gonorrhoeae). The reconstruction has been cut in the center of the spherical colony. n denotes the number of nearest neighbors [21]. b) Rod-like cells (V. cholerae) (image adapted from [22] with permission).

3.2 Fluid crystalline-like order of rod-shaped bacteria

Many bacterial species have a rod-like (ellipsoid) cell body. Recent studies have suggested that rod-like bacteria can drive collective behaviors in microbial groups bacause of their tendency to aligh their orientations with adjacent bacteria [23]. The assymetric cell shape is therefore likely to influence the structure of biofilms. Drescher et al investigated local ordering in biofilms formed by *Vibrio cholerae* [22]. Indeed, they found that bacteria showed nematic order at the center of the biofilm. The degree of ordering as defined by the nematic order parameter $S(r) = \langle 3/2 (n_i n_j)^2 - 1/2 \rangle$ (with n_i , n_j depicting the relative cellular orientations) depends on the size of the colonies and shows a minimum for colonies containing 100 cells. Nematic ordering of the first layer can be explained by a combination of surface compression due to adhesion to and sterical hindrance which causes upright orientation of cells at the surface (Fig. 3b) [24]. Subsequent directional growth causes nematic order at the center of the colonies. Reducing the strength of surface adhesion leads to disoriented growth.

3.3 The role of the extracellular matrix on biofilm structure

The extracellular matrix is crucial for controlling biofilm structure. It can be thought of as an extracellular polymer network that simulatanously acts as a glue, as a barrier to invasion, and as a "sponge" for nutrients. For many species extracellular DNA is the most abundant matrix component. It facilitates biofilm formation [25], renders microcolony formation irreversible [20], and governs the rheological properties of the biofilm. In addition to DNA, exopolysaccharides are crucial for determining the structure of the biofilm [1].

In the example described in 3.2 deletion of the protein mediating attachment between mother and daughter cells causes an increase in the average cell-to-cell distance because growth of the colony is dominated by the expansion of the extracellular matrix. The biofilm matrix also enables *V. cholerae* to establish an osmotic pressure difference between the biofilm and the external environment. This pressure difference promotes biofilm expansion by physically swelling the colony, which enhances nutrient uptake [26].

4 Physical mechanisms of cellular sorting

Biofilms formed by a single species develop spatial heterogeneity in terms of gene expression [1]. Moreover, biofilms can be formed by different species. Thus, bacteria inside the biofilm are likely to express different types or densities of surface appendages. In the following, we summarize work characterizing the effect of physical parameters (including interaction forces, cell shape, entropic effects of extracellular matrix) on cell sorting and phase separation.

4.1 Passive spatial segregation in biofilms

Within biofilms bacterial motilty is strongly reduced or fully inhibited. Therefore, clonal clusters form without active aggregation. Cells grow and divide and in the absence of motility, the offspring of a single mother cell forms a sector [27]. Even when different strains are initially well-mixed, populations that grow mostly at the surface (due to nutrient and space limitation at the center) often experience strong spatial bottlenecks. This process is refered to as spatial genetic drift and induces subdivision into monoclonal sectors.

4.2 Sorting with respect to differential interaction force

Postioning within biofilms can be controlled by interaction forces between bacteria. We have demonstrated this by using genetically modifying type IV pili (T4P) in *N. gonorrhoeae* [28]. The underlying differential adhesion hypothesis has been formulated theoretically by Albert Harris when he tried to understand cell sorting during embyonic development [29]. The idea is that bacteria contiously extert force on each other by pilus-pilus attachment, retraction, and detachment. This generates a tug-of-war whereby the bacteria tend to move in the direction where the rupture force are highest. These rupture force were tuned by genetically modifying the pili. The rupture forces were characterized using laser tweezers (Fig. 4a, Fig. S1b) [28]. We found that the morphology of the mixed colony agreed remarkably well with the prediction (Fig. 4a), indicating that bacteria sort with respect to differential interaction forces.



Fig. 4: Cellular segregation and sorting in biofilms can be triggered by various physical factors. a) Differential interaction forces. Depending on the bacterial interaction forces F (experimentally determined by rupture forces between cells), bacteria sort to the center or to the surface of mixed bacterial colonies [28]. b) Cell shape. In a mixtured biofilm with spherical and rod-like bacteria, spherical bacteria (red) tend to sort to the top of biofilms (image adapted from [30] with permission). c) Depletion forces. In the presence of non-interacting extracellular matrix (yellow), rod-like cells self-assemble into clusters driven by depletion forces (image adapted from [31] with permission).

4.3 Sorting with respect to cell shape

Another mechanism for controlling cellular position within biofilms is tuning cell shape. Smith et al performed individual-based simulations of growth in a biofilm containing both spherical and rod-like bacteria [30]. They found that spherical cells rose to the upper surface of the biofilms (Fig. 4b). They tested their predictions experimentally using strains of *E. coli* with different mutations in the actin homolog *mreB*. Using these different mutants, they tuned the aspect ration of the rod-shaped *E. coli* cells. In agreement with their predictions, bacteria sorted with respect to their shape, with spherical cells at the top and rod-like cells dominating the basal surface (Fig. 4b) [30]. This work demonstrates that by tuning cell shape, bacteria can influence positioning within biofilms.

4.4 Effects of the extracellular matrix on segregation and sorting

The extracellular matrix affects the degree of cellular clustering and can most likely trigger cell sorting. So far, we have exclusively discussed attractive interactions between bacteria. Gosh et al address the effect of repulsive forces generated by the extracellular matrix [31]. They used a particle-based modeling approach to study self-organization of rod-like bacteria in the presence of self-produced matrix. They predict that the presence of non-absorbing matrix can lead to spontaneaous aggregation of bacteria by a depletion attraction (Fig. 4c) [31]. In a different study, preliminary evidence is provided that production of extracellular matrix by a subpopulation of bacteria induces sorting to the suface of a biofilm formed by *Pseudomonas fluorescens* [32]. In this study, mutations in a regulatory protein occur repeatedly and these mutations confer positioning to the top. One of the regulatory targets of the mutated protein (RsmE) has a function secreting matrix components.

5 Bacterial fitness relates to biofilm structure

So far, we have discussed physical parameters determining adhesion, biofilm structure, and bacterial positioning within biofilms. Here, we will focus on their effects on bacterial fitness. Positioning has important effects on bacterial fitness (measured by their growth rate) in structured environments. In the absence of external stresses, surface residing bacteria would benefit the most from unlimited space, high oxygen, and nutrient concentration. In other cases positioning could help bacteria further inside the biofilm to evade stresses such as shear flow, antibiotics or immune cells. It is also plausible that relocating bacteria to the surface can act as a dispersal mechanism to spread cells to new biofilm attachment sites. Alternatively, dynamic positioning of cells within a biofilm could lead to cooperation of cells and help to avoid a static growth-arresting biofilm structure. It is therefore conceivable that bacteria have evolved mechanisms for manipulating physical interactions and thus for tuning their positions. A beautiful example illustrating this effect is an experiment where S. aureus were grown either in the presence or in the absence of the antibiotic methicillin (Fig. 1)[2]. In the presence of methicillin, bacteria formed macroscopic biofilm-structures while the bacteria remained in suspension without methicillin. Finding out which physical parameters are tuned in response to stress will be a major challenge for the future. In the following, we will describe initial steps in that direction.

5.1 Differential adhesion to surfaces affects bacterial fitness

Adhesive force govern attachment and trigger biofilm formation. But how benficial is strong attachment at later stages of biofilm development? Surface-attached bacteria have reduced access to nutrient compared to bacteria located at the top of the biofilm. Schluter et al addressed the simple but fundamental question how the magnitude to adhesion force to a surface affects the fitness of bacteria [33]. In particular, they ran individual-based simulations and experiments with *V. cholerae* with two bacterial subpopulations whose only difference was strength of adhesion. The relative fitness was dependent on the geometry of the system, in particular on the location of fluid flow providing bacteria with nutrients. When grown on a semi-permeable membrane with nutrient supply from below the membrane, the more strongly adhesive bacteria [33]. In summary, the adhesive strength of bacteria affects bacterial fitness within biofilms.

5.2 Differential bacterial interactions affect bacterial fitness

Differential adhesion forces between bacteria can also affect fitness [34]. Using *N. gonorrhoeae* and its T4P as a model, we found that loss of T4P confers a benefit because nonpiliated bacteria segregate to the front of growing colonies [34]. At the molecular level, loss of piliation strongly reduces the attractive force between bacteria generated by pili (4.2) [28]. Loss of pili occurs repeatedly through various mutations that have particularly high mutation rates. These highly mutable regions can be interpreted as switches for controlling the physical properties of T4P. Another example of a reversible mutational target associated with positioning has been found in *P. fluorescens* [32]. The combined physical properties of these two mutants allow them to spread faster compared to the individual mutants; one strain provides a wetting polymer at the colony edge and the other strain pushes both along [32]. Interestingly, both *P. fluorescens* and *N. gonorrhoeae* modulated local physical interactions within the biofilm-model of expanding colonies through mutations in specific genes. Therefore, it is tempting to speculate that altering the local physical interactions within biofilms through reversible mutations is a common strategy in biofilm adaptation.

5.3 Cell shape affects bacterial fitness

Microbes can actively change their shape in response to environmental stress including exposure to antibiotics or predation. Moreover, phylogenetic studies indicate that specific cell shapes have evolved independently multiple times [35]. However, understanding when and why particular cell shapes offer a competitive advantage remains an unresolved question. Bacteria sort with respect to their cell shape with spherical bacteria accumulating at the top of the biofilm [30]. Positioning at the top can confer a selective advantage and indeed depending on the location of the nutrient source, spherical or rod-like bacteria grow faster in an experiment using *E. coli* with different cell shapes [30].

6 Conclusion

Relying on a few bacterial systems, the direct links between physical interactions, biofilm structure, and bacterial fitness have been characterized. Among the major physical determinants are attractive interactions mediated by cell appendages, depletion forces and osmotic pressure generated by the extracellular matrix, and cell shape. In the microbial world, the huge diversity of species hampers the search for general molecular principles for controlling biofilm structure. Threfore, searching for common physical principles which are likely implemented by different molecule in different species will be a major challenge for the future.

Appendices

A Force spectroscopy for characterizing bacterial interaction forces

Laser tweezers (laser traps, optical tweezers) are very useful tools for force spectroscopy and for characterizing the properties of molecular motors. Laser tweezers are created by focussing laser light through a microscope objective into a diffraction limited spot. A dielectric particle (such as a glass bead or a bacterium) experiences force, pulling the particle into the trap center. Displacements down to the sub-nanometer range (10^{-10} m) and forces in the range of hundreds of piconewtons (10^{-12} N) can be measured under optimal conditions [36]. Laser traps have been used for measuring the attractive force of bacteria to abiotic surfaces (Fig. S1a) and to other bacteria (Fig. S1b) [20,28] and for characterizing motor properties of surface appendages [8,37]. Atomic force microscopes (AFM) are particularly useful for force spectroscopy. Its force spectrum extends to the nanonewton (10^{-9} N) range. AFMs have been used for measuring rupture forces of bacterial appendages from surfaces (Fig. S1c) and other bacteria (Fig. S1d) [38] and for characterizing the elastic properties of cell appendages [11,12].



Fig. S1: Methods for measuring interaction forces. Laser tweezers were used for characterizing the attractive force of cell appendages to a) a surface and b) between cells. *AFM allow for measuring attractive forces of cell appendages to c) a surface and d) between cells.* Furthermore, the setup shown in c) is useful for characterizing elastic properties of cell appendages.

References

- [1] C. D. Nadell, K. Drescher, and K. R. Foster, Nat Rev Microbiol 14, 589 (2016).
- [2] J. B. Kaplan, E. A. Izano, P. Gopal, M. T. Karwacki, S. Kim, J. L. Bose, K. W. Bayles, and A. R. Horswill, MBio 3, e00198 (2012).
- [3] M. K. Hospenthal, T. R. D. Costa, and G. Waksman, Nat Rev Microbiol 15, 365 (2017).
- [4] B. Maier and G. C. Wong, Trends Microbiol 23, 775 (2015).
- [5] A. J. Merz, M. So, and M. P. Sheetz, Nature 407, 98 (2000).
- [6] R. Marathe et al., Nat Commun 5, 3759 (2014).
- [7] C. Holz, D. Opitz, L. Greune, R. Kurre, M. Koomey, M. A. Schmidt, and B. Maier, Phys Rev Lett **104**, 178104 (2010).
- [8] B. Maier, L. Potter, M. So, C. D. Long, H. S. Seifert, and M. P. Sheetz, Proc Natl Acad Sci U S A 99, 16012 (2002).
- [9] J. Ribbe, A. E. Baker, S. Euler, G. A. O'Toole, and B. Maier, J Bacteriol 199 (2017).
- [10] K. Zhao, B. S. Tseng, B. Beckerman, F. Jin, M. L. Gibiansky, J. J. Harrison, E. Luijten, M. R. Parsek, and G. C. L. Wong, Nature **497**, 388 (2013).
- [11] M. Forero, O. Yakovenko, E. V. Sokurenko, W. E. Thomas, and V. Vogel, PLoS Biol 4, e298 (2006).
- [12] N. Biais, D. L. Higashi, J. Brujic, M. So, and M. P. Sheetz, Proc Natl Acad Sci U S A **107**, 11358 (2010).
- [13] M. Andersson, E. Fallman, B. E. Uhlin, and O. Axner, Biophys J 91, 2717 (2006).
- [14] H. C. Berg, Annu Rev Biochem 72, 19 (2003).
- [15] A. Boehm *et al.*, Cell **141**, 107 (2010).
- [16] A. Persat, Y. F. Inclan, J. N. Engel, H. A. Stone, and Z. Gitai, Proc Natl Acad Sci U S A **112**, 7563 (2015).
- [17] Y. Luo, K. Zhao, A. E. Baker, S. L. Kuchma, K. A. Coggan, M. C. Wolfgang, G. C. Wong, and G. A. O'Toole, MBio 6 (2015).
- [18] C. K. Ellison *et al.*, Science **358**, 535 (2017).
- [19] I. Hug, S. Deshpande, K. S. Sprecher, T. Pfohl, and U. Jenal, Science **358**, 531 (2017).
- [20] L. Dewenter, T. E. Volkmann, and B. Maier, Integr Biol (Camb) 7, 1161 (2015).
- [21] A. Welker, T. Cronenberg, R. Zöllner, C. Meel, K. Siewering, E. R. Oldewurtel, and B. Maier, (submitted).
- [22] K. Drescher, J. Dunkel, C. D. Nadell, S. van Teeffelen, I. Grnja, N. S. Wingreen, H. A. Stone, and B. L. Bassler, Proc Natl Acad Sci U S A **113**, E2066 (2016).
- [23] D. Volfson, S. Cookson, J. Hasty, and L. S. Tsimring, Proc Natl Acad Sci U S A **105**, 15346 (2008).
- [24] J. Yan, A. G. Sharo, H. A. Stone, N. S. Wingreen, and B. L. Bassler, Proc Natl Acad Sci U S A **113**, E5337 (2016).
- [25] M. Zweig, S. Schork, A. Koerdt, K. Siewering, C. Sternberg, K. Thormann, S. V. Albers, S. Molin, and C. van der Does, Environ Microbiol **16**, 1040 (2014).
- [26] J. Yan, C. D. Nadell, H. A. Stone, N. S. Wingreen, and B. L. Bassler, Nat Commun 8, 327 (2017).
- [27] O. Hallatschek, P. Hersen, S. Ramanathan, and D. R. Nelson, Proc Natl Acad Sci U S A **104**, 19926 (2007).
- [28] E. R. Oldewurtel, N. Kouzel, L. Dewenter, K. Henseler, and B. Maier, Elife 4 (2015).
- [29] A. K. Harris, J Theor Biol **61**, 267 (1976).

- [30] W. P. Smith, Y. Davit, J. M. Osborne, W. Kim, K. R. Foster, and J. M. Pitt-Francis, Proc Natl Acad Sci U S A **114**, E280 (2017).
- [31] P. Ghosh, J. Mondal, E. Ben-Jacob, and H. Levine, Proc Natl Acad Sci U S A 112, E2166 (2015).
- [32] W. Kim, F. Racimo, J. Schluter, S. B. Levy, and K. R. Foster, Proc Natl Acad Sci U S A **111**, E1639 (2014).
- [33] J. Schluter, C. D. Nadell, B. L. Bassler, and K. R. Foster, Isme J 9, 139 (2015).
- [34] R. Zollner, E. R. Oldewurtel, N. Kouzel, and B. Maier, Sci Rep 7, 12151 (2017).
- [35] D. T. Kysela, A. M. Randich, P. D. Caccamo, and Y. V. Brun, PLoS Biol 14, e1002565 (2016).
- [36] J. R. Moffitt, Y. R. Chemla, S. B. Smith, and C. Bustamante, Annu Rev Biochem 77, 205 (2008).
- [37] T. L. Min, P. J. Mears, L. M. Chubiz, C. V. Rao, I. Golding, and Y. R. Chemla, Nat Methods 6, 831 (2009).
- [38] Y. F. Dufrene, Trends Microbiol 23, 376 (2015).
E 6 Tissue Growth

Jens Elgeti Theoretical Soft Matter and Biophysics Institute of Complex Systems Forschungszentrum Jülich GmbH

Contents

1	Introduction	2
2	Homeostatic pressure	4
3	Particle based Modeling	5
4	Competition	7
5	Fluid or Solid	10
6	Spheroids	12
7	Conclusions	15

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.



Fig. 1: *Examples of tissue mechanics* – (**a-e**) *Dorsal Closure in the development of Drosophila.* The inner tissue (amnioserosa) contracts and the outer ephithelium closes the opening. The mechanical force for closure comes from tension in the amnioserosa (increased by apoptosis) and from a super-cellular pursestring contracting at the leading edge (from [1]). (**F-H**) *Growth of the tissue making up the entestine against an elastic background causes an elastic buckling instability resulting in the typical wrinkled patterns (from [2])*

1 Introduction

Materials have classically been studied either in equilibrium or close to equilibrium. Often, the driving force comes from outside, either in form of moving boundaries as in a shear plate rheometer, or from an external field like gravity.

Active materials are intrinsically out of equilibrium. They consume energy internally on the scale of the microscopic constituents. Materials can be active in many ways. For one, internal forces can create stresses that create spontaneous flows, as in collections of filaments and motors (C5 and C6) or swimming micro-organisms (E4). See chapter E3 vor an overview on active matter.

A different route to active materials that we want to address here is breaking mass conservation, a material that locally generates and destroys itself. Examples are polymerizing gels or cellular tissues. In this lecture, we will focus on cellular tissues, but many aspects are more generic and can be applied to any material that has a stress dependent self generation.

From a biological point of view, the notion has spread that "Mechanosensitivity in one form or another appears to be a property shared by all cells of the body and by all phyla from mammals to bacteria."[4] (See also chapter D5, D6, and D7). Examples of mechanical feedback range across all organisms found on earth: In the development of an embryo tissues grow and large scale motions of tissues reorganize the growing embryo[5]. Dorsal closure in drosophila



The stack air-cushion applied to the breast ; front and back river.—(a.) The spring passing over the shoulder, and fixed in front to the centre point of the shield (b), and in the back to the back-pad g.—(b.) The shield, receiving on a screw at its centre point of the spring (a), where also is attached a double buckle for straps; round the border of the rin appears an edge (b) of the air-cushion.—(c.) The stop-cock for inflating the air-cushion.—(c.) The stop-cock for inflating the air-cushion.—(d.) A belt fixed to the centre point of the pad (g) behind, passing round underneath the healthy breast, and fixed in front by leather straps to the buckle attached to the shield, and also to studies at two opposite points of the circumference of the latter.—(e.) A soft moreable pad attached to the inner surface of that part of the bet (d.), which its bedon the healthy breast.—(f.) A second bit passing much the arm of the unaffected side by a shoulder strap, carried behind, and fixed to the buckle in front of the frame.—(g.) The counter pressure pad.—(h.) Border of the air-cushion, as before a bell substitued for the spring (a), and a spring faxed to the pad (g), passed round the vasist in the direction of, and lying outside the bett (d).]

Fig. 2: Soft mechanical pressure as a possible route for treatment of cancer has been tried already in the early 19th century. This device from Walshe in 1846 [3] exerts a "constant, equal, and uniform pressure" on the tumor.

development is mediated via tensile stresses in the amnioserosa, gastrulation is caused by a spontaneous curvature of a certain group of cells[6]. Also the human adult body undergoes a constant cell turnover. For example our intestine is under constant renewal. Again, mechanics are speculated to be the cause of the observed patterns [7, 2]. Finally, cancer, the "emperor of all maladies" [8] is a disease of growth. A notion that mechanics may play a role and could be used for treatment has appeared early on [3] (See Fig. 2), but only recently received much more attention [9, 10, 11].

Thus, understanding the mechanics of growth can not only help us understand the development of an organism, it may also contribute to understand and treat diseases or even design artificial tissues. Indeed tissue engineering and organ printing are becoming thriving fields of research. We begin to understand that mechanical properties of the environment are key players when growing an artificial organ.

In this lecture we will adress some simple key concepts of modeling the mechanics of tissue growth. We begin with a simple continuum model in the spirit of hydrodynamics, followed by a particle based model similarly inspired by particle based hydrodynamics. Subsequently we study growth and competition and mechanical response with these models. Finally, we address an experimental setup that can be well compared with the models.



Fig. 3: Gedankenexperiment to the homeostatic pressure. (A) A tissue inside a finite compartment permeable to nutrients, but impermeable to cells, will grow to fill the available space. If one of the walls is a movable piston connected to a spring, the tissue will press on the spring until the restoration force is strong enough to hinder further growth. The pressure in this steady state is than termed the homeostatic pressure. (B) If the spring of the second compartment is replaced by a tissue of different homeostatic pressure, the pressure in the chamber evolves to an intermediate pressure. The weaker cells die, while the stronger tissue grows. Eventually, the stronger tissue completely takes over the compartment.

2 Homeostatic pressure

When considering the mechanics of growing tissues, growth is the key difference to the otherwise classic mechanics of liquids or solids $[12, 13]^1$. To account for growth and death, a source term is added to the continuity equation

$$\partial_t \rho + div(\vec{J}) = (k_d - k_a)\rho \tag{1}$$

with k_d and k_a the division and apoptosis rate respectively and \vec{J} the flux of cells. Here, we will study the effect of mechanical forces on growth. Certainly, growth depends on the biochemical environment, but we will assume this to be constant across the system at hand. Instead, we will ask ourselves which effect mechanical stress has on growth. A simple assumption is the existence of a homeostatic state, where division and apoptosis are balanced. Homeostasis is the tissue dynamics equivalent of a steady state. We will assume that the homeostatic state is characterized by a homeostatic density or a homeostatic pressure, the two being coupled to each other via the tissues equation of state. Pressure in this context is the force felt by the cells hindering growth, not to be confused with the hydrostatic pressure. In terms of thermodynamics, it is the conjugate force to cell volume. The homeostatic pressure (or stress) $P_H = -\sigma_H$ and density ρ_H is the pressure (density) at which division is balanced by apoptosis.

We can then expand the growth rates around this homeostatic state [14]

$$\sigma = \chi^{-1}(\rho - \rho_H) + \sigma_H \tag{2}$$

$$(k_d - k_a) = \kappa'(\rho - \rho_H) \tag{3}$$

$$(k_d - k_a) = \kappa (\sigma - \sigma_H). \tag{4}$$

This implies that a tissue, left in a finite compartment grows to its homeostatic density, exerting its homeostatic pressure on the surrounding (See Fig. 3).



Fig. 4: *Tissues are modeled as collection of sticky soft colloids. Two particles representing a cell expand actively. Once a critical size is reached, the cell divides, creating new particles.*

3 Particle based Modeling

Particle based modeling has been used since a long time to model materials in general and fluids in particular. One key concept is to introduce simplified particles and interactions to capture the right dynamics on more mesoscopic length and time scales. For example both multi particle collision dynamics and dissipative particle dynamics (See A10 by M. Ripoll) use very simple, but also very different interactions. However, because they both satisfy momentum balance locally, they result in the Navier Stokes equation on larger length-scales. In the same spirit, we introduce a simplified particle based model for tissue growth. One cell is represented by two point particles *i* and *j* at positions $\vec{r_i}$, which repel each other with a growth force

$$\vec{F}_{ij}^g = \frac{B}{(r_{ij} + r_0)^2} \hat{r}_{ij},$$

with a growth strength B. $\vec{r_{ij}} = \vec{r_j} - \vec{r_j}$ is shorthand for the vector connecting particles *i* and *j*. This growth force drives the particles constituting one cell apart. Motivated by the "size checkpoint" of real cells, virtual cells divide, once the particles reach a critical distance R_c . For division, two new particles are inserted in the immediate surrounding of the constituents of the mother cell to form the daughter cells. This is the key active mechanism where mechanics feed back on growth. If the surrounding medium provides additional (mechanical) resistance to growth (due to pressure for example), the process slows down. The second active process is cell apoptosis or death. In the spirit of minimal modeling, apoptosis is introduced by a constant rate of cell removal. In principle this rate could also change through a mechanic feedback, however, this would introduce further parameters.

The rest of the particle based growth model is "classic" soft colloids. A soft repulsive potential of strength f_0 prevents overlap of particles belonging to different cells

$$\vec{F}_{ij}^v = f_0 \left(\frac{R_{pp}^5}{r_{ij}^5} - 1\right) \hat{r}_{ij},$$

where all interactions are set to zero beyond a cut off distance R_{pp} . Adhesion between cells is modeled by a simple constant attractive force

$$\vec{F}^a_{ij} = -f_1 \hat{r}_{ij}.$$

¹Active stresses by cells, and stresses generated by divisions and apoptosis events, can also significantly contribute, and even change the appropriate material laws. We will consider some of these later in the lecture.



Fig. 5: Net growth rate $k_b = k_d - k_a$ as a function of the difference between the imposed pressure P^i and the homeostatic pressure P^b_H . Close to the homeostatic pressure, the growth rate is indeed a linear function of the pressure differences. The slope across different parameters is also very similar. However, for larger pressure differences, nonlinearities appear.

Finally, as in all simulations of active systems, energy has to be dissipated. Here, a Dissipative Particle Dynamics (DPD) -Thermostat (See B3) offers simultaneously momentum conserving dissipation \vec{F}_{ij}^d , and random fluctuations \vec{F}_{ij}^r ,

$$\vec{F}_{ij}^r = \sigma \omega^R(r_{ij})\xi_{ij}\hat{r}_{ij} \tag{5}$$

$$\vec{F}_{ij}^d = -\gamma \omega^D(r_{ij})(\vec{v}_{ij} \cdot \hat{r}_{ij})\hat{r}_{ij}.$$
(6)

where σ and γ are the amplitudes of fluctuations and dissipation, $\xi_{ij} = \xi_{ji}$ a symmetric gaussian random variable with zero mean and unit variance and a weight functions $\omega(r_{ij})$ as in DPD. Additionally, dissipation by friction with the background can be modeled by $\vec{F}^b = \gamma_b \vec{v}$. In total, the force between particles reads

$$\vec{F}_{i} = \underbrace{\vec{F}_{ik}^{g} + \vec{F}_{ik}^{d} + \vec{F}_{ik}^{r}}_{\text{intracellular forces}} + \underbrace{\sum_{j} \left(\vec{F}_{ij}^{v} + \vec{F}_{ij}^{a} + \vec{F}_{ij}^{d} + \vec{F}_{ij}^{r} \right)}_{\text{intercellular forces}} + \vec{F}_{j}^{b} \tag{7}$$

On large enough time and length-scales, this model reproduces the concept of homeostatic pressure. A tissue confined by a piston grows to a steady state pressure. When the tissue is compressed, the division rate goes down and apoptosis reduces the number of cells, until the pressure is back at its homeostatic value. Close to the homeostatic pressure, growth rates indeed follow the linear expansion postulated in eq. 2 (See Fig. 5). However, they also indicate strong nonlinearities as the pressure further deviates from the homeostatic one.

Finally, a note on parameters and interpretation. Similar to fluid models, where particles are not representing water molecules, the analogy of two particles describing a cell is just an interpretation to simplify the explanation. When interpreting data from particle based simulations, a simulation cell could represent many real cells, or vice versa. It is important to match the right mesoscopic quantities when interpreting simulation results in light of experiments.



Fig. 6: Two tissues in a finite compartment face competition for space. **(top)** In isolation, tissue A grows faster, but tissue B reaches a higher homeostatic pressure. **(bottom)** In competition, the compartment reaches a pressure intermediate of the two homeostatic pressures. This causes tissue A to shrink, and tissue B to grow, until the compartment is completely taken over by tissue B.

4 Competition

If two tissues are present in a finite compartment, they enter competition. One of the two tissues will take over the compartment eliminating the other.²

Using the theory of homeostatic pressure, we can understand the general dynamics. Assuming the tissue as incompressible and neglecting elastic and viscous forces, competition can be expressed by simple number balance. The division rates of both tissues are determined by the stress

$$k_c = \kappa (\sigma - \sigma_{Hci}),\tag{8}$$

where we abbreviated the net growth rate of tissue c as $k_c = (k_{dc} - k_{ac})$. Due to incompressibility the compartment has space for a total number $N_{tot} = N_A + N_B$ cells of tissue A and B. The total number of cells will be constant over time, i.e.

$$0 = \partial_t N_{tot} = k_A N_A + k_B N_B. \tag{9}$$

This is easily solved for the stress

$$\sigma = \frac{N_A \sigma_{HA} + N_2 \sigma_{HB}}{N_A + N_B} \tag{10}$$

and the growth-rates

$$N_A k_A = \kappa \frac{N_A N_B 2(\sigma_{HA} - \sigma_{HB})}{N_A + N_B} = -N_B k_B.$$
 (11)

²Ongoing research actually suggests the existens of stable steady states of mixed tissues. This will be explored in the future.

If we define $\phi = N_A/N_{tot}$ as the fraction of A-cells and the difference in homeostatic pressure as $\Delta \sigma = (\sigma_{HA} - \sigma_{HB})$, this simplifies to

$$\partial_t \phi = \kappa \Delta \sigma \phi (1 - \phi) \tag{12}$$

i.e. a simple logistic growth on time scales $\kappa \Delta \sigma$.

In a particle based simulation, we can study the more complex behavior of a three dimensional competition and more generic parameters. In particular, we can compare tissues that have a higher devision rate if unconstrained, versus tissues with a lower division rate, but a higher homeostatic pressure. As predicted by the analytic model, even in this extreme example, it is the tissue with the higher homeostatic pressure that wins the competition (see Fig. 6). Even though these tissues are neither incompressible nor do they have the same κ , the shape of the curve is still well reproduced by a logistic growth predicted by the simple number balance (Eq. 12).

If the competition instead happens on a substrate, stress is no longer constant[15]. The friction with the substrate results in an effective decay length, over which the pressure decays to the homeostatic one. Thus cell number balance (eq. 12) is now local, and has a convective derivative:

$$\partial_t \phi_A + v \partial_x \phi_A = (k^A - k^B) \phi_A (1 - \phi_A).$$
(13)

where we have neglected for simplicity terms corresponding to cell diffusion. Now cellnumberfraction ϕ is a local quantity, with

$$\phi_A(x,t) = \frac{n_A(x,t)}{n_A(x,t) + n_B(x,t)}$$
(14)

the number fraction of A type cells, where $n_c(x,t)$ is the total number density of cells of type $c = \{A, B\}$ at time t and position x. The velocity field v is determined by a generalized incompressibility condition

$$\partial_x v = k^A \phi_A + k^B (1 - \phi_A), \tag{15}$$

which relates the velocity gradient to the growth rates k^c . As above, the growth rates k^c , depend on the local stress σ

$$k^{c}(\sigma) = -\kappa^{c}(\sigma_{H}^{c} - \sigma), \tag{16}$$

Again, the growth coefficients κ^c and homeostatic stresses σ_H^c are tissue properties and the only model parameters. Eq. 13 is solved in a comoving coordinate system $s = x - v_0 t$ by $\phi_A(s) = \Theta(s)$, the Heaviside step function describing a sharp interface, where v_0 is the interface velocity. Force balance implies

$$\partial_s \sigma(s) = 2\rho_b \gamma_{bg} v(s) \tag{17}$$

with the cell density ρ_b and the substrate friction coefficient γ_{bg} . The velocity v(s) is determined by Eq. 15, which then results in the stress profile

$$\sigma(s) = \begin{cases} \sigma_H^B + \lambda_B \frac{\sigma_H^A - \sigma_H^B}{\lambda_A + \lambda_B} \exp\left(\frac{s}{\lambda_B}\right) & \text{for } s \le 0\\ \sigma_H^A + \lambda_A \frac{\sigma_H^B - \sigma_H^A}{\lambda_A + \lambda_B} \exp\left(-\frac{s}{\lambda_A}\right) & \text{for } s > 0. \end{cases}$$
(18)



Fig. 7: Properties of a quasi one-dimensional competition a) Simulation snapshots at early and late times for relatively thin channels. The image shown is zoomed in into the interface region. The stronger tissue (red) invades into the weaker tissue (blue). The initially flat interface roughens over time, while propagating to the right. Note that the scale is the same as for c) and e). b) Interface height h (average position of the interface) as a function of time t for three different homeostatic stress differences. c) Velocity profile $v_x(s)$ as a function of position s. Solid blue line represents analytical prediction, where the parameters are obtained by separate independent simulations. d) Interface width w as a function of time t for three different homeostatic stress differences $\Delta \sigma_H^*$. e) Stress as a function of position s. The dashed lines correspond to the homeostatic stresses of the two tissues. Solid blue line as in c). f) Interface velocity v_0 as a function of the homeostatic stress difference $\Delta \sigma^*$ (solid blue line as in c)). c), e) correspond to time averages of one simulation with a run length of $Tk_a = 80$ and a system width L/R = 20. Each point and line in b), d), f) correspond to a single simulation. From Ref. [15]



Fig. 8: Rheology of three dimensional tissues in the particle based model. Results from an in-silico rheology experiment. In a tissue one layer 0 < z < 1 is pulled in the x-direction, a second layer is pulled in the opposite direction. (left) Without cell division and apoptosis the tissue behaves like a yield stress solid. Only above a critical stress does the tissue begin to flow. (right) With cell division and apoptosis, the tissue flows for any given force.

The characteristic length scale $\lambda_c^2 = (\kappa^c 2\rho_b \gamma_{bg})^{-1}$ is fixed by the friction γ_{bg} , the bulk density ρ_b , and the growth coefficients κ^c . Note that $\sigma(0) = (\sigma_H^A \lambda_B + \sigma_H^B \lambda_A)/(\lambda_A + \lambda_B)$ is the weighted average of the homeostatic stresses, which simplifies to $\sigma(0) = (\sigma_H^A + \sigma_H^B)/2$ for $\lambda_A \approx \lambda_B$. The finite length λ_c leads to interface motion with a constant interface propagation velocity

$$v_0 = \frac{1}{2\rho\gamma_{bq}} \partial_s \sigma|_{s=0} = \frac{\Delta\sigma_H}{2\rho\gamma_{bq}(\lambda_A + \lambda_B)},\tag{19}$$

which depends linearly on the difference in the homeostatic stresses, $\Delta \sigma_H = \sigma_H^A - \sigma_H^B$. We find the prediction for the stress profile to be in very good agreement with the simulation

results (see Fig. 7). The velocity profile, which can be obtained analytically from Eqs. 17 and 18, also agrees very well with the simulations (see Fig. 7b) without adjusting any parameters. Furthermore, the linear dependence of the interface velocity v_0 on the homeostatic pressure difference $\Delta \sigma_H$ in the simulations can be understood by the analytical calculation (see Fig. 7c) [15]. When the chamber is no longer two dimensional, a finite interface width develops. Similar to classical interfacial growth, the interface width scales with a powerlaw with the interfacial length [15].

5 Fluid or Solid

Real tissues have a very complex viscoelastic behavior. On timescales of seconds or minutes, the rheological properties can be measured in a relative straight forward extension of classical rheological experiments. However, growth happens on very long timescales, which are difficult to access experimentally. Extrapolating from short timescales to the long term behavior is difficult, but we can exploit the long timescales to predict some generic features. The timescale of growth is the longest timescale in the system, much longer than any other cellular timescale like attachment and detachment of cellular adhesion proteins. One can thus already guess that on the timescales of growth, the tissue behaves as a viscous fluid. Even if the tissue behaves like an elastic solid on short timescales, division and apoptosis relax stress, leading to fluid



Fig. 9: Rheology of three dimensional tissues in the particle based model. (left) Viscosity as determined by a virtual shear plate experiment for different tissue parameters. For cell turnover-rate k larger than the shear rate, the predicted 1/k behavior is well reproduced. For smaller cell turnover rates, the yielding behavior leads to a finite plateau of the viscosity (right) The viscosity can be measured in many different ensembles, which agree with each other when the shear rate is small compared to the cell turnover rate. Fitting a Maxwell model to oscillatory shear flow experiments, yields an long timescale small stress viscosity. This agrees best with the expected 1/k behavior.

behavior on long timescales. Combining these simple arguments with a simple Maxwell model for viscoelastic behavior results in a rather good prediction for the tissue viscosity η . In the Maxwell model, the product of the elastic modulus E and the relaxation time τ give the viscosity

$$\eta = E\tau. \tag{20}$$

Arguing that the relaxation time is proportional to the cell turnover time, we get

$$\eta = E\tau \approx E/k,\tag{21}$$

with cell turnover $k = k_a = k_d$.

Using the particle based simulations, these arguments can be tested with a very different approach. First, consider a tissue in periodic boundary conditions, one layer is pulled up, and another downwards. A liquid material would display a linear velocity gradient between the two layers, while an elastic material would show no velocity, but a strain gradient. Indeed, if apoptosis is disabled by setting the apoptosis rate $k_a = 0$, division stops automatically when the finite compartment is densely packed with cells. In this situation, the particles behave like a colloidal glass. Below a critical shear stress, no continuos motion can be observed. Above a critical stress, the tissue yields, and starts to flow with rather small velocities (See Fig. 8). With apoptosis, the tissue is in a homeostatic state with continuos cell turnover $k = k_a = k_d$. Each division and apoptosis event locally relax some stress, leading to a viscous behavior.

To really understand the rheological behavior different types of rheology experiments have to be performed. The simplest one is the shear plate experiment: The tissue is confined between two parallel plates and the top one is moved by a constant velocity. This imposes a shear on the tissue. If the shear rate is larger than the cell turnover rate, the rheology is dominated by the yielding behavior described above for no cell turnover. The effective viscosity is dominated by



Fig. 10: Growth of tissue spheroids in dialysis bags. Cells from a colon carcinoma cell line (CT26) are placed inside a dialysis bag. Large molecular weight dextran is added to the growth medium to exert an osmotic pressure via the bag on the spheroids. After 12 days most spheroids are sacrificed for cryosections. Two spheroids are allowed to continue to grow. They grow to the same final volume after pressure is released. From [16]

the yield-stress. If the cell turnover-rate k is larger than the shear rate, the stress-relaxation via cell-turnover becomes important, and the effective viscosity is well predicted by Eq.(21). On a more quantitative level, oscillatory shear experiments can be performed. A simple way to implement these, is to impose an oscillating (in space and time) force density. From these, the full complex viscosity η_C can be extracted for different frequencies.

6 Spheroids

Experimentally measuring the mechanical properties of tissues on long timescales is very challenging indeed. The tissue has to experience the same mechanic stimulus over the course of many days or weeks, while keeping the biochemical environment constant. One such experiment is growing tissue spheroids under isotropic osmotic stress. Cells are placed in a dialysis bag. The bag is put in a standard growth medium, with an added high molecular weight dextran polymere. The dialysis membrane is permeable to all the relevant nutrients, but not to the dextran. The resulting osmotic pressure of the dextran is transmitted over the membrane onto the spheroid. It turns out that indeed pressures as low as 500Pa are able to significantly slow growth.

Unfortunately this approach is very time consuming, because of difficulties in handeling spheroids in dialysis bags. A simpler approach, is to use the outer cell membrane as dialysis membrane. At the concentrations necessary to exert pressures of the order of many kPa dextran does not effect cells if grown on substrates. One can thus assume it has no poisonous effect. Also, it can not penetrate into the spheroid. Thus, the main effect of dextran in the growth medium is to exert a well controlled mechanic stress on the growing spheroid. These experiments can be performed on a much larger number of spheroids and significant statistics can be obtained. The results (see Fig. 11) confirm the growth reducing effect of pressure. However, some details



Fig. 11: Growth of tissue spheroids under isotropic pressure. Cells from a colon carcinoma cell line (CT26) are placed in wells and grown in standard growth medium with added dextran. The dextran exerts a well defined isotropic pressure on the spheroids. Volume is measured over the course of two weeks. From [16].

come at a surprise. First of all, the tissue never shrinks due to pressure. Equation 1 would predict an exponential growth or shrinkage linearly depending on pressure. However, the growth reducing effect seems to saturate above 5kPa. Furthermore the growth is clearly not exponential. The answer is, that while pressure is constant across the system, the biochemical, but also the biomechanical, environment is not. While the first certainly has some effects, it is the latter that we focus on here.

Particle based simulations can proof the existence of the biomechanical differences in the environment. With the right parameters, the growth under pressure can be well reproduced. Because biochemistry is assumed uniform in the simulation, there is a mechanical explanation that accounts for the observed growth phenomena.

In order to increase volume, a cell has to deform its surrounding, i.e. it builds up a strain dipole. Insertion of a strain dipole in an elastic medium is much easier close to a free surface, than it is in bulk. Thus divisions at the surface are favored mechanically over divisions in bulk.

Indeed both simulations and experiments exhibit increased division close to the surface (see Fig. 12). Simulations allow for a much more detailed analysis of the division profile. The simulation data suggests that division is enhanced in a thin layer close to the interface, but than saturates below (!) the apoptotic rate. This finally explains the growth curves. On average, cells die more than they divide ($k = (k_d - k_a) < 0$) but in a thin layer close to the surface, division is strongly increased (δk_s). Using simple cell number balance, we arrive at an equation for the growth of the spheroid

$$\partial_t N = kN + N_s \delta k_s,\tag{22}$$

with N_s the number of cells in the surface layer that is favored for division. Assuming a constant density and a constant surface layer thickness λ , we get

$$\partial_t V = kV + \lambda \delta k_s (V^{2/3} (36\pi)^{1/3}).$$
(23)



Fig. 12: Cryosections and virtual cryosections of growing spheroids. In experiments, spheroids are sacrificed, frozen and cut into thin slices. By labeling Ki67 cell divisions are marked in cyan. In simulations, recently divided cells are marked in cyan, others in grey. Independent of pressure, most divisions happen at the surface, while the division is strongly suppressed by pressure in the bulk. From [17]

Fitting Eq. 23 to the growth curves of simulations and experiments is able to reproduce the growth curves quite well and allows to extract how the different growth rates depend on pressure. It turns out that the surface rate $\lambda \delta k_s$ is largely unaffected, while k is negative for all applied pressures and decreases further with increasing pressure (see Fig. 13).

The negative growth rate in the bulk of the spheroid, and division at the surface naturally lead to a stable steady state where a flow of cells from the surface balances the net death in the bulk. For an incompressible fluid of cells this flow can be estimated analytically. The conservation equation reads

$$\nabla \vec{v}(r) = k(r),\tag{24}$$

where we assumed radial symmetry. As above, the growth rate k(r) is k in the bulk, and increased by δk_s within a range λ of the surface at position R(t), i.e.

$$k(r) = k + \delta k_s \Theta(R(t) - r - \lambda)$$

with Θ the Heavyside step function. With the boundary condition of zero velocity at the center, we can solve for v and get:

$$v_r = \begin{cases} \frac{1}{3}kr & \text{if } r < R(t) - \lambda\\ \frac{1}{3}(k+\delta k_s)r - \frac{1}{3}\delta k_s \frac{(R(t)-\lambda)^3}{r^2} & \text{if } r > R(t) - \lambda \end{cases}$$
(25)

The spheroid grows with the velocity at the surface $\partial_t R(t) = v_r(R)$. Note that expanding $\partial_t R$ to first order leads back to Eq. 23. Thus both equations predict the same steady state size

$$R_{\infty} = -3\delta k_s \lambda/k \tag{26}$$

The flow of cells predicted by Eq. 25 can be measured experimentally by marking cells at the surface and following them over the course of time[18]. This fit allows an independent measurement of the growth rates from the growth curves and leads to very similar values.



Fig. 13: Growth rates extracted from growth curves by fitting Eq 23. The growth rate at the surface δk_s is much less effected than the bulk growth rate k_d . From [17].

7 Conclusions

Growing materials display a set of phenomena completely unknown to conventional materials. Key aspects are a self developed homeostatic stress, viscous behavior due to material turnover and competition. To understand the dynamics of growing materials continuum mechanics very similar to classic viscoelastic materials can be used. The main difference comes from an additional source (or sink) term in the continuity equation. Mesoscopic particle based simulations can bridge the difficult gap between analytic theory and real experiments, but more importantly provide a third perspective on the problems at hand.

The field of growing materials is still young, and many aspects have not been addressed yet. For example, both, simulations and analytic model, have a source term in the continuity equation, while of course in reality matter remains conserved. The truth is that nutrients and fluid (the "interstitial fluid") pass between the cells and get turned into tissue material. Thus the growth (and death) rate is indeed a conversion rate from interstitial fluid to tissue material. See for example Ref. [19] for a further discussion.

References

- D. P. Kiehart, C. G. Galbraith, K. A. Edwards, W. L. Rickoll, and R. A. Montague, J. Cell Biol. 149(2), 471 (2000).
- [2] A. E. Shyer, T. Tallinen, N. L. Nerurkar, Z. Wei, E. S. Gil, D. L. Kaplan, C. J. Tabin, and L. Mahadevan, Science 342(6155), 212 (2013), ISSN 0036-8075.
- [3] W. Walshe, *The Nature and Treatment of Cancer*, The Nature and Treatment of Cancer (Taylor and Walton, 1846).
- [4] A. W. Orr, B. P. Helmke, B. R. Blackman, and M. A. Schwartz, Dev. Cell 10(1), 11 (2006).
- [5] M. A. Wozniak and C. S. Chen, Nat. Rev. Mol. Cell Biol. 10(1), 34 (2009).
- [6] C.-P. Heisenberg and Y. Bellaiche, Cell 153(5), 948 (2013), ISSN 0092-8674.
- [7] E. Hannezo, J. Prost, and J. . F. Joanny, Phys. Rev. Lett. 107(7), 078104 (2011).
- [8] S. Mukherjee, *The Emperor of All Maladies* (Scribner, 2010).
- [9] D. T. Butcher, T. Alliston, and V. M. Weaver, Nat. Rev. Cancer 9(2), 108 (2009).

- [10] A. Fritsch, M. Hockel, T. Kiessling, K. D. Nnetu, F. Wetzel, M. Zink, and J. A. Kas, Nat. Phys. 6(10), 730 (2010), ISSN 1745-2473.
- [11] O. J. T. McCarty, M. R. King, and P. A. Insel, American Journal of Physiology Cell Physiology 306(2), C77 (2014), ISSN 0363-6143.
- [12] L. Landau, E. Lifshitz, A. Kosevich, and L. Pitaevskiĭ, *Theory of Elasticity*, Course of theoretical physics (Butterworth-Heinemann, 1986), ISBN 9780750626330.
- [13] L. Landau and E. Lifshitz, *Fluid Mechanics*, no. Bd. 6 in Course of theoretical physics (Elsevier Science, 2013), ISBN 9781483140506.
- [14] M. Basan, T. Risler, J.-F. Joanny, X. Sastre-Garau, and J. Prost, HFSP J 3(4), 265 (2009).
- [15] N. Podewitz, F. Jlicher, G. Gompper, and J. Elgeti, New J. Phys. 18(8), 083020 (2016).
- [16] F. Montel, M. Delarue, J. Elgeti, L. Malaquin, M. Basan, T. Risler, B. Cabane, D. Vignjevic, J. Prost, G. Cappello, and J. F. Joanny, Phys. Rev. Lett. 107(18), 188102 (2011).
- [17] F. Montel, M. Delarue, J. Elgeti, D. Vignjevic, G. Cappello, and J. Prost, New J. Phys. 14(5), 055008 (2012).
- [18] M. Delarue, F. Montel, O. Caen, J. Elgeti, J.-M. Siaugue, D. Vignjevic, J. Prost, J.-F. Joanny, and G. Cappello, Phys. Rev. Lett. 110(13), 138103 (2013).
- [19] J. Ranft, J. Prost, F. Jlicher, and J.-F. Joanny, Eur Phys J E Soft Matter 35(6), 46 (2012).

E 7 Necessity and feasibility of large-scale neuronal network simulations

Maximilian Schmidt 1, Markus Diesmann 2,3,4 and Sacha J. van Albada 2

 ¹ Laboratory for Neural Circuit Theory, RIKEN Brain Science Institute (RIKEN BSI), Wako-Shi, Saitama, Japan
 ² Institute of Neuroscience and Medicine (INM-6), Institute for Advanced Simulation (IAS-6), JARA Institute Brain Structure Function Relationships (JBI 1/INM-10), Jülich Research Centre, Jülich, Germany
 ³ Department of Psychiatry, Psychotherapy and Psychosomatics, Medical Faculty, RWTH Aachen University, Aachen, Germany
 ⁴ Department of Physica Faculty 1, DWTU Aachen, University

⁴ Department of Physics, Faculty 1, RWTH Aachen University, Aachen, Germany

Contents

1	Neuronal network modeling	2
2	Minimal layered cortical network model	3
3	Necessity and feasibility of brain-scale simulations3.1Irreducibility of network dynamics3.2Simulation technology	7 8 8
4	Multi-area model of macaque visual cortex as stepping stone to the human brain	10
5	Conclusion	16

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Neuronal network modeling

Computational neuroscience employs mathematical and computational tools to understand the information processing properties of the brain. The fundamental challenge in this endeavor is to understand the relationships between the structure of the brain, its dynamics, and the emerging function. The dynamics in general feeds back to brain structure via activity-dependent plasticity of synaptic strengths and the rewiring of connections. Nerve cells, called neurons, are thought to be the main information carriers in the brain, and as such are the primary object of investigation in computational neuroscience. Neuronal network models describe sets of interconnected neurons via coupled differential equations to study how the connectivity structure and the dynamics of the individual elements determine collective network properties. Since models can be modified at will, for instance by systematically changing parameters, they are useful tools for creating and testing hypotheses.

Modeling the brain and its neural circuits is a formidable challenge for several reasons. First, there is actually no such thing as "the brain": brains differ substantially between species and between individuals of the same species. Second, brains feature an enormous complexity, the human brain being organized into hundreds of specialized regions each with a variety of cell types and gene expression patterns, which are connected and interact in intricate ways. Connectivity patterns are specific at the level of brain regions, cell types, individual cells, and even subregions of dendritic trees [1]. Furthermore, plasticity on temporal scales ranging from milliseconds to years reorganizes brain circuitry and changes its chemical composition. Third, the number of physical, environmental, and behavioral conditions in which animals can find themselves is enormous. Thus, we need to focus on a limited set of conditions to investigate, for instance the resting state in which the subject performs no particular task and receives no particular sensory stimuli.

When choosing a modeling approach, one needs to decide on the desired level of abstraction. Including all details is not feasible, but the details relevant to the question at hand need to be included. Also germane to the choice of method is what measurements are available to constrain the model. Different measurement techniques probe brain structure and activity at different spatial and temporal resolutions (Figure 1 on page 4). Brain connectivity can for example be investigated using axonal tracing, an invasive but precise and directed method visualizing one or a few projections at a time, or diffusion tensor imaging (DTI), a non-invasive but less precise and non-directed method visualizing whole-brain connectivity. Techniques for probing brain activity include microscopic methods such as optical imaging and intracellular recording, and meso- and macroscopic methods such as voltage-sensitive dye (VSD) imaging, local field potential (LFP) recording, functional magnetic resonance imaging (fMRI), electroencephalography (EEG), magnetoencephalography (MEG), and positron emission tomography (PET). Measurements are often not available for the living brain (in vivo), but are obtained from brain slices or cell cultures (in vitro), in which external inputs are absent and the conditions differ from those we primarily wish to understand, adding another level of complexity to the modeler's job.

Neuronal network models can be classified into two categories. Bottom-up models start from a description of the model components and their connectivity, and analyze the emerging dynamics. Since complete sets of experimental data necessary for defining neuronal network models are not available for individual brains, models combine data from multiple experiments, and thus in a sense represent an average across brains. This averaging removes much of the detailed structure of individual brains that is for instance created by experience-dependent plasticity. For this reason, bottom-up modeling, while it can provide insight into relationships between

structure and dynamics, is unlikely to reproduce the biological function of large neuronal networks. Top-down studies address this issue by defining a desired functionality of the model, e.g., recognition of certain visual features, and using intuition about possible mechanisms to propose network architectures or plasticity rules necessary to achieve the given function. The drawback of this research strategy is that there may be many different network structures implementing the same function. There is therefore a consensus in the neuroscience community that both approaches have to be followed simultaneously and continuously checked for compatibility and consistency.

One of the primary brain structures studied with neuronal network models is the cerebral cortex, the thin sheet of cells on the outside of mammalian and reptilian brains. A reason for the wide interest in cerebral cortex is its central role in higher functions such as planning, memory, attention, language, and conscious awareness. The cerebral cortex of mammals consists of allocortex which can be divided into fewer than six layers, and the evolutionarily more recent neocortex, with six layers. We are mainly interested in neocortex and also refer to it simply as 'cortex'. The cortex is a relatively homogeneous structure and its organization is fairly well preserved across species, which enables knowledge about individual cortical areas in particular species to be used to make informed guesses about cortex in general. Researchers hope to uncover a universality of cortex as, although from mouse to man its size increased by three orders of magnitude, its local structure has been well preserved in evolution, and it appears very similar under the microscope, independent of whether it is processing visual or auditory information, or is responsible for motor planning [2].

2 Minimal layered cortical network model

Cortical networks are characterized by a high density of neurons ($\sim 10^5$ neurons/mm²), sparse connectivity ($\sim 10\%$ average pairwise connection probability locally within areas), and a high average number of incoming synapses per neuron (the so-called 'in-degree') of around 10^4 . Taking these constraints together implies that a patch under 1mm² of cortical surface represents a minimal building block where a substantial fraction of the incoming synapses arise within the patch. In smaller networks one of the two constraints is necessarily not fulfilled: either there are too few synapses per neuron or the connection probability is too high. We refer to such a local cortical network as a 'microcircuit'. The interpretation of cortical microcircuits as generic building blocks supporting cortical function, in view of broad similarities between local cortical circuits across species and cortical areas, is expressed in the term 'canonical microcircuit'. The connectivity at this scale has been characterized in an increasingly refined manner (see, e.g., [3], and [4, 5] for reviews), and detailed wiring diagrams have been assembled.

Connection probabilities have been estimated using different methodologies: For instance, using neuron reconstructions along with simple assumptions on the distributions of synapses over dendrites [6], or from fractions of pairs of neurons in which stimulation of one neuron elicits a response in the other neuron [7]. In the microcircuit model described in [8], these different datasets are combined to derive a consistent connectivity map for the populations of excitatory and inhibitory neurons in each of the cortical layers 2/3, 4, 5, and 6. Here, layers 2 and 3 are grouped together because of extensive similarities, and layer 1 is disregarded because it contains only a scattering of neurons. From the computational perspective, this layered network is an extension of the classic balanced random network model [e.g., 9, 10]. In this model, an excitatory and an inhibitory population balance each other in a dynamic manner to produce a



Fig. 1: Organization and measurement of cerebral cortex. Cerebral cortex consists of areas that are defined based on functional, architectonic, molecular, connectivity, or topographic properties, or a combination thereof. Areas can have sizes ranging from a few hundred thousand to tens of millions of neurons, and consist in turn of highly connected local circuits of around a hundred thousand neurons. A variety of measurement modalities probe cortex at different scales; for instance, the electroencephalogram (EEG) has a spatial resolution on the order of several centimeters, while certain types of optical imaging can visualize microscopic brain circuits in detail. Figure by Susanne Kunkel.

network state with low synchrony that is referred to for simplicity as 'asynchronous irregular'. Such activity with low synchrony across neurons and irregular spiking of individual neurons is characteristic of cerebral cortex especially in the awake condition.

To emphasize connectivity rather than the intrinsic properties of neurons, the work of [8] represents single cells by a simple model called the leaky integrate-and-fire (LIF) model. The LIF neuron model approximates a biological neuron as an electrical circuit with a resistance and a capacitance connected in parallel (an RC circuit), where the resistance replaces the voltage-gated ion channels of the cell membrane and the capacitance mimics the membrane's ability to maintain a separation between electrically charged ions. Describing the current across the membrane due to single synaptic inputs as a jump followed by an exponential decay, the membrane potential V and the synaptic input I_s are given by two differential equations reading

$$\tau_m \frac{\mathbf{V}}{\mathbf{t}} = -(V - E_{\mathrm{L}}) + R_{\mathrm{m}} I_{\mathrm{s}}(t),$$

$$\tau_s \frac{\mathbf{I}_{\mathrm{s}}}{\mathbf{t}} = -I_{\mathrm{s}} + \tau_{\mathrm{s}} \sum_{i,j} J_i \delta(t - t_i^j),$$
 (1)

where τ_m is the membrane time constant, $R_m = \tau_m/C_m$ is the resistance of the neuron, C_m is the capacitance, E_L is the leak potential, τ_s is the synaptic time constant, J_i is the maximum



Fig. 2: Leaky integrate-and-fire neuron model and balanced random network model. A Postsynaptic potential (PSP) of a leaky integrate-and-fire (LIF) neuron in response to a single spike. The inset shows the RC circuit represented by the LIF neuron model. **B** Responses to constant (black) and fluctuating (red) input. Spikes shown as vertical lines above the plot. For constant input, the neuron is regularly driven to the threshold, while for fluctuating input, the spiking occurs irregularly. **C** Sketch of a balanced random network with two populations with $N_{\rm E}$ excitatory neurons and $N_{\rm I}$ inhibitory neurons, respectively. The populations are coupled both to themselves and to each other, and they receive an external input consisting of Poisson spike trains. **D** Synchronous, irregular regime with slow oscillations for weak external input. Each horizontal line shows the spike times of a given neuron. **E** Asynchronous, irregular regime for strong external input.

of the post-synaptic current induced by a single action potential ("spike") from neuron *i*, and *j* indexes the incoming spikes from the given neuron. The arrival of a spike evokes a postsynaptic potential (PSP) in the receiving neuron which has the form of the difference between a fast-rising exponential with time constant τ_s and a slowly decaying exponential with time constant τ_m (Figure 2 on page 5, panel A). When the membrane potential exceeds a fixed threshold θ , it is lowered to the reset potential and a spike is sent to all postsynaptic partners. Depending on the nature of the input, the neuron will display different types of behavior: For a constant input current, the time course of the sub-threshold membrane potential follows the charging curve of an RC circuit. The neuron will either settle in an equilibrium or, for sufficiently large input, it will emit spikes at regular intervals (Figure 2 on page 5, panel B). In such a case where spiking is mainly caused by the mean input current driving its membrane potential across threshold, the neuron is said to be in the "mean-driven regime". Even when the mean input does not cause the membrane potential to cross the threshold, sufficiently large fluctuations in the input can still cause spiking at irregular intervals. The neuron is then said to be in the "fluctuation-driven



Fig. 3: The microcircuit model of early sensory cortex. A Sketch of the microcircuit model of [8], consisting of four layers with one excitatory (triangles) and one inhibitory (circles) population each. Arrows indicate connections between populations. Each population is driven by spike trains drawn from a Poisson process representing inputs from parts of the brain that are not explicitly modeled. **B** Stationary state of the model. Left: Raster plot showing spikes of excitatory populations as black dots and spikes of inhibitory populations as gray dots. Right: Box plot [11] of the population-averaged firing rates (top), irregularity quantified as the coefficient of variation of the inter-spike interval distribution (CV ISI) (middle) and synchrony quantified as the variance of the spike count histogram divided by its mean (bottom). **C** Sketch of the propagation of transient thalamocortical input in the circuit. Figure adapted from [8] (with permission).

regime". This situation is particularly interesting for the use in cortical models because, as mentioned above, the cerebral cortex displays irregular spiking.

In the balanced random network model, an excitatory and an inhibitory population of neurons are coupled with random connectivity defined by a pairwise connection probability, which is typically on the order of 0.1. The excitatory population has outgoing synapses with positive weights, so that their spikes increase the membrane potentials of the target neurons, while synapses formed by the inhibitory population have negative weight and decrease the membrane potentials of the target neurons. Both populations are driven by stochastic external input, typically trains of spikes generated from a Poisson process with a time-independent probability of spike occurrence.

Depending on the relative strengths of the excitatory and inhibitory synaptic weights and the strength of the external input, this model can display different types of dynamics. For weak external input and strong inhibition, individual neurons emit spikes at irregular intervals, while

slow oscillations occur on the network level (Figure 2 on page 5, panel C). Increasing the external input retains the irregular spiking of neurons but removes the slow oscillation, leading to a nearly asynchronous state (Figure 2 on page 5, panel D). For weak inhibitory synapses, the network is no longer balanced and exhibits synchronous regular spiking.

The microcircuit of [8] uses this two-population network as a building block for each of the four layers (Figure 3 on page 6, panel B). The eight populations are connected to each other with population-specific connection probabilities derived from the aforementioned anatomical [6] and electrophysiological experiments (with the largest contribution from [12]). External inputs are modeled as Poisson spike trains with rates proportional to estimated numbers of synapses received from external sources.

The circuit reproduces several aspects of cortical dynamics. Neurons spike irregularly with low synchrony across cells and heterogeneous firing rates across populations (Figure 3 on page 6, panel B). Just like in cerebral cortex, inhibitory neurons exhibit higher rates than excitatory neurons, despite their identical parameters in the model, showing that the connectivity of the circuit itself can produce this feature. All populations exhibit fast oscillations of around 80 Hz. When a transient thalamic stimulus is fed into layer 4 and to a lesser extent into layer 6, activity propagates in a manner matching experimental findings, first to layer 2/3 and then to the lower ("infragranular") layers 5 and 6 (Figure 3 on page 6, panel C). Propagation of activity without ringing is achieved via a "handshaking" mechanism where an activated layer inhibits the layer from which it received its input, thereby terminating the activity in the sending layer, as if to signal that it has received the message.

3 Necessity and feasibility of brain-scale simulations

The 1 mm^2 cortical microcircuit captures only about half the synapses impinging on the neurons in the circuit. The remaining synapses come from surrounding cortex, other cortical areas, and subcortical regions. A self-consistent description of the circuit dynamics is therefore only possible when the recurrent interactions with these other parts of the brain are taken into account. Interactions between cortical areas are associated for instance with low-frequency neuronal population activity [13, 14] and express themselves in clustered correlation patterns across cortical areas as measured by functional magnetic resonance imaging (fMRI) [15].

Another reason why modeling a 1 mm^2 circuit does not suffice is that this only enables links with experiments at or below this scale, such as intracellular recordings or spike trains extracted from the high-frequency part of extracellular local field potential (LFP) recordings. However, as mentioned in Section 1, many brain activity data are obtained at larger scales: the spatial reach of the LFP including its low-frequency part can be several millimeters or more depending on the correlations in the underlying neuronal dynamics [16]; voltage-sensitive dye (VSD) imaging is similarly a mesoscopic method; fMRI can cover the whole brain with a resolution that ranges from about 1 to 5 mm; and electroencephalography (EEG) and magnetoencephalography (MEG) are macroscopic imaging methods with low spatial resolution. Both in order to aid the interpretation of such mesoscopic and macroscopic signals, and to enable constraining network models with these data modalities, we need to model larger parts of the brain.

3.1 Irreducibility of network dynamics

To simulate macroscopic neuronal networks, the corresponding models are commonly downscaled in terms of their numbers of neurons and synapses in order to fit the models into computer memory and to limit the simulation time. This raises the question to what extent such downscaled models preserve the dynamics of their full-scale counterparts. For cortex, because of its low-synchrony irregular activity, this question can be addressed with the help of mean-field theory in which the inputs to the neurons are approximated as Gaussian noise. This so-called diffusion approximation results from the central limit theorem for large numbers of (near-)independent inputs. Given a sufficiently simple neuron model, analytical equations for the first few moments of the population-averaged neuronal activity can be derived. We have studied the scalability question for networks of binary neurons with input-dependent stochastic transitions between two states, and for networks of leaky integrate-and-fire neurons. The networks can have any number of populations which are randomly connected with populationspecific connection probabilities, delayed interactions, and they receive external input in the form of a DC drive or Poisson spike trains. Assuming a stationary state, the mean-field theory for such networks shows that, if we do not allow the single-neuron properties to be changed, the mean activity and the population-averaged pairwise correlations between neurons cannot simultaneously be preserved beyond a certain level of downscaling [17]. More specifically, reducing the number of neurons increases the correlations, and reducing the number of synapses per neuron changes the shape of correlations as a function of time lag. These analytical results are confirmed by network simulations. Preserving the second-order statistics of neuronal spiking activity is relevant because mesoscopic measures of brain activity like the local-field potential (LFP) and the EEG are driven by the fluctuations of neuronal activity and the fluctuations are governed by the cross-correlations between neurons: if a model does not faithfully capture the correlation structure, also the mesoscopic measures will be inaccurate. Furthermore, synaptic plasticity, considered to be the main biophysical basis of learning, also depends on the temporal relationship between pre- and postsynaptic spiking. This does not mean that simulations have to be carried out at full scale, but once a suspiciously interesting result is obtained its validity should be checked with a full-scale simulation.

3.2 Simulation technology

Historically it has not been possible to simulate the microcircuit model at full scale. Despite the difficulties in interpreting the results from downscaled networks researchers therefore decided on a number of neurons and scaled down the number of synapses per neuron such that the connection probability remained in a reasonable range. As synapses outnumber neurons by far, the former dominate memory consumption. When using the aforementioned scaling procedure, the total number of synapses depends quadratically on the total number of neurons, conventionally termed network size. With a quadratic dependence of memory consumption on problem size, scaling up simulations to the full complement of neurons and synapses therefore seems rather hopeless. However, we need to consider that the number of synapses per neuron is finite and known. Once the scale of the microcircuit is reached where every neuron is supplied with its natural number of synapses, the number of synapses grows linearly with network size, and networks become necessarily more sparsely connected. In this view, simulating larger networks just requires increasing the computer size. In 2005 we broke through this barrier of

 10^5 neurons with a distributed simulation code optimized for spiking neuronal networks [18]. A relevant insight of the study is that if synapses are located on the compute node harboring the postsynaptic neuron, the presynaptic neuron only needs to send a single event, independent of the number of target neurons it has on this node. In addition, with this memory layout also network creation becomes an ideally parallel task. A second insight is that communication in intervals of the minimal delay in the system is sufficient to maintain causality [19] and decouples the computation step size of the solver from the duration of the communication interval. Initially the global pairwise exchange algorithm (CPEX) was used for communication but this was later replaced by a collective operation (MPI Allgather, [20]) sending all spike data to all compute nodes based on the realization that with a round-robin distribution of neurons across nodes and fewer compute nodes than synapses per neuron, each spike finds at least one target on all nodes. Today most simulation codes are based on the principles of [18], and NEST, originating from first steps made in the middle of the 1990s [21], is the most widely used code for large-scale spiking neuronal networks. With growing network size and the increasing complexity of computers the development of simulation technology has become a research area requiring substantial resources. A major further step in adapting NEST to modern computer architectures was a hybrid parallelization scheme combining distributed computing with multithreading [22, 23]. The initial commit to a repository with version control, using CVS at the time, was made in 1996. In 2001 the NEST Initiative (www.nest-initiative.org) was formed as a collaboration between several labs with two goals: (1) release of open source simulation software for neuronal networks, and (2) publication of underlying simulation technology in scientific journals. The first public release was made in 2004 at the European Advanced Course in Computational Neuroscience (ACCN) in Óbidos. The first GPL release came in 2012 and the NEST Initiative transformed into a registered society. Since 2015 an open development platform is operated with formal code reviews, continuous integration, and a fortnightly open developer video conference. A push towards brain-scale neuronal networks was given in 2008 by the start of the project in Japan leading to the construction of the K computer [24]. With the perspective of petascale supercomputers routinely available for neuroscience, first concepts for multi-area models at cellular resolution were developed. As a first step we developed a mathematical model of memory consumption showing which component of the software becomes the dominant consumer of memory at the different scales [25]. This model is continuously refined and guides the software design. On September 1, 2011, NEST ran on the K computer for the first time, establishing the usability of supercomputers for neuroscience [26]. Representing synapses only on the compute node where the postsynaptic neuron resides has the disadvantage that for all neurons in the network the compute node needs to maintain a list of the local synapses. With growing network size, however, connectivity becomes sparser. Consequently, the number of local target lists with only a single or no synapse increases. This sparseness was taken care of by a novel data structure for the target lists which dynamically modifies its representation depending on the number of elements, and by a specific sparse table optimized for the rejection of spikes without local targets [27]. The technology led to the largest general spiking neuronal network simulation carried out to date [28]. Demonstrating that neuroscience can orchestrate the full memory of a petascale computer in a simulation of a single neuronal network constituted a milestone in the discussion on the relevance of exascale computers for this field. At the same time it became clear that the idea of representing synapses solely on the compute node harboring the postsynaptic neurons reached a limit. For networks larger than one billion neurons, the sparse table alone would consume all the computer memory. Therefore, we initiated a project exploring a simulation scheme consisting of two phases. In the first phase compute nodes construct their part of the network as before but subsequently communicate to the other compute nodes from which presynaptic neurons they expect to receive spikes. This preserves parallel network construction but enables the compute node to decide where its outgoing spikes need to be sent. In the second phase the dynamical state of the network is propagated using a collective operation (MPI_Alltoall, [20]) that sends spikes only to compute nodes where they are needed. Consequently, the sparse table can be removed. Ironically, researchers come back here to some of the original ideas of [18]. While the new design aims at exascale computers, it already pays off at the petascale. Initial benchmarks show that simulation time is substantially reduced on existing supercomputers due to a reduced communication load and improved parallelization across the compute nodes [29]. At present it remains unclear whether traditional computer architectures will ultimately deliver the performance required in terms of time-to-solution and energy-to-solution for brain-scale simulations of processes like learning and development. Therefore, in parallel to the work on traditional supercomputers, intense research is carried out on neuromorphic computers [30].

4 Multi-area model of macaque visual cortex as stepping stone to the human brain

As an example of a large-scale cortical neuronal network, we study macaque visual cortex, and investigate interactions between areas and relationships between the structure and dynamics of cortex. Macaque visual cortex lends itself to investigation for several reasons: First, the macaque brain is relatively similar to the human brain [31, 32], and thus, macaque models may teach us something about how the human brain works. Second, a wide range of experimental data are available on the anatomy, physiology, and function of this system which partly cannot be obtained in the human. Third, visual cortex is organized in a systematic manner, which enables anatomical properties to be estimated where experimental data are missing [33]. The strategy to predict data from from mathematical regularities at parameter settings where no measurements have been obtained before is routine in physics. Nevertheless, in neuroscience, researchers found the idea so relevant that they termed it "predictive neuroinformatics" [34]. Hierarchically ascending connections and signals are called 'feedforward', those in the reverse direction 'feedback', and those between areas at the same hierarchical level 'lateral'. Our vision is to develop concepts and theoretical tools in the monkey for later transfer to the human.

The macaque visual cortex can be parcellated into a set of 32 cortical areas [35], which differ in their size, structure, and functional role in visual processing. As a first approximation, we represent each area by a 1 mm^2 patch of cortex. Averaged across all the modeled areas, this represents a reduction by a factor of 200 compared to the biological system. While Section 3 has shown that larger simulations are possible, this reduction increases the turnover rate of the simulations, enabling a more extensive parameter exploration. For each local circuit, we take the microcircuit model of [8] as a prototype. Each area thus consists of four layers, except area TH which lacks the granular layer 4, with an excitatory and an inhibitory population in each layer.

To perform dynamical simulations, we need to derive population sizes and synapse numbers for the connections between populations. The fundamental challenge is to find a way to integrate data from different sources into a common framework and estimate missing data based on regularities of cortical structure. The resulting network structure is therefore a combination of experimental data and heuristics based on experimental literature ([36], Figure 4 on page 12).

For instance, laminar cell densities are not available for all 32 areas of the model so that we have to estimate values for areas without available data. We use a categorization of the areas into architectural types that reflects laminar distinctiveness and the thickness of the granular layer 4 [37]. We then estimate the laminar cell densities for areas with missing information as the averages for areas of the same architectural type. Furthermore, measured total cortical thicknesses and relative laminar thicknesses are complemented with linear estimates based on the neuron densities. In the experimental data, we observe that the number of neurons decreases along the visual hierarchy from primary visual cortex V1 (197,936 neurons per mm²) to area TH (73,251 neurons per mm²). In conjunction with taking the layer-specific ratio of excitatory to inhibitory neurons to be the same as in cat V1 as given by [6], this leads us to a complete set of population sizes.

To derive the connectivity of the model, we conceptually distinguish four different types of synapses: The first two types are simulated in the model: local synapses originating within each 1mm² patch (type I), and cortico-cortical synapses formed by projections from other vision-related cortical areas (type II). Two other types are beyond the scope of the network under investigation: synapses from the same cortical area, but outside the 1mm² patch (type III), and synapses from other cortical non-visual and non-cortical areas (type IV). For simplicity, inputs from these connections are modeled as stochastic input.

We model local connectivity (type I) as spatially uniform. However, to determine the number of type I synapses, we assume the probability for two neurons to form one or more synapses to decay with the distance between them according to a Gaussian profile. The radius of the modeled patch ($R = \sqrt{1/\pi}$ mm) then determines the cutoff of the Gaussian and therefore the number of synapses within the modeled patch. We take the 8×8 connectivity matrix of the microcircuit model as a prototype and adapt it to the area-specific population sizes in the following manner: in order to keep the population-averaged spike rates approximately equal to those in the microcircuit model, we assume the relative numbers of synapses per target neuron across projections to be preserved. These heuristics enable us to derive the local connectivity for each area.

The connectivity between areas is mostly based on the results of axonal tracing experiments. In these experiments, a tracing substance is injected and taken up by a small number of cells. The tracer then propagates from the axon of a neuron to its cell body (retrograde tracing) or from the cell body forward along the axon (anterograde tracing). This technique enables labeling source neurons (retrograde tracing) and axonal ramifications on the target side (anterograde tracing).

The model combines such tracing data with statistical regularities to fill in missing values. Its inter-area connectivity is defined in four steps: First, a directed connection between two areas is established if it is included in the CoCoMac database [35, 38, 39, 44, 45, 46] or the retrograde tracing dataset of [40]. Second, the total number of synapses between two areas is based on fractions of labeled neurons (FLN) as given by [40] or estimated using the exponential fall-off of FLN with inter-areal distance [41]. Third, we determine laminar patterns of connections based on fractions of supragranular labeled neurons (SLN), both measured [42] and estimated using a sigmoidal fit against differences in neuronal densities, and layer-specific data on target [35, 45, 46, 47, 48, 49, 50, 51, 52] and source [35, 46, 49, 51, 53, 54] patterns from CoCoMac. Finally, we account for the possibly different laminar positions of synapses and cell bodies in the target areas by computing the conditional probability that a cortico-cortical synapse in a layer ν is formed on a dendrite of a specific cell type, using the morphological reconstructions of [6]. This derivation assumes that the population-averaged number of synapses formed on a dendritic tree is proportional to the length of the dendrite, a version of so-called Peters's rule [55]. In this way, we systematically derive the probability of two neurons being connected



Fig. 4: Construction principles of the multi-area model. Top: Determination of the population sizes. The size of a population is computed as the product of the layer-specific cell density and its volume computed as thickness times the fixed surface area of 1mm². Population sizes decrease along the gradient of architectural types. The ratio between excitatory and inhibitory neurons is taken from layer-specific data of [6] leading to an average proportion of $\sim 80\%$ excitatory and $\sim 20\%$ inhibitory cells. Figure adapted from [8], with permission. Bottom: Construction of the model connectivity. Local connectivity within areas is based on the microcircuit model of [8]. Cortico-cortical connections are first determined on the area level from binary data from the CoCoMac database [38, 39] and quantitative tracing data from [40], which is completed using the exponential fall-off of connection density with inter-areal distance [41]. Synapses between each pair of cortical areas are then distributed over source and target layers based on layer-specific tracing data from [42] and CoCoMac. Synapses in the receiving area are subsequently assigned to cells according to layer- and cell-type-specific dendritic densities from [6]. These derivations result in distinct laminar patterns between feedforward, feedback, and lateral connections. Based on a theoretical method using mean-field theory [43], the resulting connectivity matrix is refined to improve the phase space of the network. Panels adapted from [36], which contains the detailed derivation.

in any pair of populations in the network. The resulting network comprises about 4 million neurons interconnected via approximately $2 \cdot 10^{10}$ synapses.

Simulations of this network on the supercomputer JUQUEEN reveal a bistable phase space, with a low-activity attractor with reasonable activity except for near-silent excitatory populations in layers 5 and 6, and an attractor with unrealistically high activity across areas and populations. To improve the dynamical state of the network, which we take to be the low-activity state, we make use of mean-field theory to predict the stationary population-averaged firing rates [43]. We furthermore derive equations describing the influence of single connections on the network activity, which allow us to make targeted modifications of the network connectivity within the uncertainty of the data. This enables us to include constraints from experimental activity measurements, namely that cortical populations exhibit spiking rates within an approximate range of [0.05, 30] spikes/s, into the network definition.

The resulting network reproduces experimental findings on cortical activity on multiple scales. The neurons spike irregularly with an overall average spike rate of 14.6 spikes/s. By increasing the strengths of synapses in cortico-cortical connections to a moderate level through multiplication by a factor χ , with stronger synapses onto inhibitory than excitatory neurons, multiplying cortico-cortical synaptic strengths onto inhibitory neurons by an additional factor $\chi_{\mathcal{I}}$, we move the network into a regime of metastable dynamics. In this regime, the system is poised just below the edge of stability of the low-activity fixed point, with transient excursions closer to the edge, corresponding to elevated activity.

Figure 5 on page 14, panels A–C show the spiking patterns of all populations in three example areas. The area-averaged activities exhibit fluctuations on various time scales (Figure 5 on page 14, panel D), in contrast to the activity of the isolated microcircuit model of [8], which is concentrated at high frequencies. Thus, interactions between the areas cause slow fluctuations to emerge when the model is in the metastable regime. For primary visual cortex (V1), recordings of the spiking activity of 140 neurons across all cortical layers during rest are available [56, 57], of which Figure 5 on page 14, panels E and F show two segments. To compare the power spectra between simulation and experiment, it is important to match the number of neurons, since this influences the relative contributions of autocorrelations and cross-correlations to the overall spectrum. The reason is that the number of contributions from autocorrelations equals the number of neurons. Taking this matching into account, the simulated spectrum of V1 spiking activity is close to the spectrum computed from the experimental recordings (Figure 5 on page 14, panel G). Furthermore, the distribution of firing rates across neurons in V1 in the model corresponds closely to that of the experimental data (Figure 5 on page 14, panel H).

Another type of experiment with which we can compare the simulated activity is so-called resting-state functional magnetic resonance imaging (fMRI). Here, 'resting-state' implies that the animal does not receive any particular sensory stimuli (especially visual, auditory, olfactory, and gustatory stimuli are absent; somatosensory input is difficult to control and is therefore present), and is not performing a task. In the resting state, the activity measured by fMRI, namely the blood oxygenation level-dependent (BOLD) signal, waxes and wanes together in clusters of areas at a rate below 0.1 s^{-1} [15, 62]. The activity within clusters is positively correlated, and that between areas in different clusters is less positively or even negatively correlated. These correlations between the activity traces of different areas are also referred to as 'functional connectivity', to distinguish them from structural connectivity and to emphasize that the anatomical structure can support different dynamical patterns. The BOLD signal has been shown to reflect neuronal activity, and in particular the synaptic inputs to neurons [63]. Thus,



Fig. 5: Spiking activity of the network. A–C Raster plots of spiking activity of 3% of the neurons in primary visual cortex (V1) (A), secondary visual cortex (V2) (B), and the frontal eye field (FEF) (C). Blue dots: spike times of excitatory neurons, red dots: spike times of inhibitory neurons. (D) Area-averaged firing rates, shown as raw binned spike histograms with 1ms bin width (gray) and convolved histograms, with a Gaussian kernel (black) of optimal width [58]. E–H Spike recordings from 140 neurons in primary visual cortex of macaque monkey [56, 57]. E–F Raster plot of experimentally recorded spiking activity for $t \in [5s, 8s]$ (E) and $t \in [392s, 395s]$ (F). Neurons are sorted according to depth of the recording electrodes with neurons closest to the surface of the brain at the top. G Power spectra of spike histograms for experimental recording (yellow), and simulated activity of 140 neurons (gray) and all neurons (black) across all populations of area V1. Inset: enlargement for frequencies up to 5Hz. H Distribution of spike rates across single cells in experimental data (yellow) and of simulated spike trains across all layers and populations of V1 (black).



Fig. 6: Inter-area interactions. A Simulated functional connectivity (FC) for cortico-cortical synaptic strength scaling factor $\chi = 1.9$ measured as the zero-time-lag correlation coefficient of synaptic input currents. **B** FC of macaque resting-state fMRI. The matrix elements are sorted according to simulated clusters determined with the Louvain algorithm [59]. C Pearson correlation coefficient of simulated FC vs. experimentally measured FC as a function of scaling factor χ for cortico-cortical synaptic weights with $\chi_{\tau} = 2$ (dots) and $\chi_{\tau} = 1$ (triangle). Dashed line, Pearson correlation coefficient of structural connectivity vs. experimentally measured FC. **D** Clusters in the connectivity graph, indicated by the color of the nodes: The hierarchically lowest visual areas (green), areas in the so-called 'dorsal stream', which processes moving stimuli (red), polysensory areas in the temporal lobe (light red), mixed cluster (light blue), areas in the so-called 'ventral stream', which performs object recognition (dark blue), and areas in the frontal lobe (purple). Black, connections within clusters; gray, connections between clusters. Line thickness encodes logarithmized numbers of outgoing connections per neuron (outdegrees). Only edges with outdegrees representing > 0.1% of all outgoing connections from a given area are shown. For visual clarity, clusters are spatially segregated, and inside clusters, areas are positioned using a force-directed algorithm [60]. E Alluvial diagram [61] showing the differences in the clusters for the structural connectivity and the simulated FC (colors distinguish clusters).

we can compare functional connectivity from resting-state fMRI with correlations between the traces of summed synaptic inputs to different areas in our model. A structured interaction pattern emerges in the network dynamics with a cluster of strongly correlated areas and one that is less correlated both internally and with the other cluster (Figure 6 on page 15, panel A). At these moderate cortico-cortical synaptic strengths where the model exhibits metastable dynamics, the functional connectivity of the simulation agrees well with resting-state fMRI data from macaque monkeys (Figure 6 on page 15, panels B, C). Importantly, the agreement between simulation and experiment exceeds the correlation between the structural connectivity and the experimental functional connectivity, indicating the enhanced explanatory power of the dynamical model.

We compare these groups of dynamically correlated areas with 7 clusters of strongly connected areas that we identify in the structural connectivity using an unsupervised clustering algorithm [64]. These clusters of the structural connectivity correspond with known groupings of areas based on functional properties (Figure 6 on page 15, panel D). Areas within structural clusters mostly stay together in the dynamical clusters, but there is also substantial mixing (Figure 6 on page 15, panel E).

5 Conclusion

In recent years we have derived a systematic way to combine the growing body of experimental knowledge on the structure of cortex into a consistent network description that enables us to study the dynamics of mammalian cortex using numerical simulations at the single-cell level. The resulting model reproduces fundamental aspects of cortical dynamics at different spatial and temporal scales. There is, however, much room for improvement, and the model should be regarded as the first prototype in a series of models. For instance, the restriction to random connectivity below the level of neuronal populations should be lifted by incorporating spatial dependence of connectivity within and between areas as well as higher-order connectivity statistics found in experimental studies [65, 66]. This would be a step towards studying stimulus-driven dynamics and functional aspects of such cortical networks.

Next we will make an executable formal description of the multi-area model of visual cortex publicly available so that it can help develop refined and extended models, akin to the microcircuit model of [8], which already serves as a building block and example network, both for own collaborative work [16, 30, 36, 67, 68, 69] and in the wider community [70, 71, 72]. Adding further brain structures like the cerebellum, hippocampus, thalamus, and basal ganglia will require a larger variety of biophysical mechanisms to be included in the simulation code. With neuromodulators [73] and gap junctions [74], first steps have been made. It may also be required to turn to multi-scale models in the sense that not all components of the model are described at the same level of abstraction. Many questions remain open, but recent work has opened a pathway for the combination of population and spiking models [75]. Ultimately multi-simulator approaches may be required [76]. Our hope is that with the progress of simulation technology we will reach a situation where experts on subsystems publish their codes in common formats such that they can be combined into more complete brain models.

Acknowledgments

We are grateful to our colleagues in the NEST developer community for continuous collaboration. Use of the JUQUEEN supercomputer in Jülich was made possible by VSR computation time grant JINB33. Partly supported by Helmholtz Portfolio Supercomputing and Modeling for the Human Brain (SMHB), the European Union Seventh Framework Programme under grant agreement no. 604102 (Human Brain Project, HBP), the European Union's Horizon 2020 research and innovation programme under grant agreement no. 720270 (HBP SGA1), and Priority Programme "Computational Connectomics" (SPP 2041) of the German Research Foundation (DFG). This research used resources of K computer at the RIKEN Advanced Institute for Computational Science. Supported by the project Exploratory Challenge on Post-K Computer (Understanding the neural mechanisms of thoughts and its applications to AI) of the Ministry of Education, Culture, Sports, Science and Technology (MEXT). All network simulations carried out with NEST (http://www.nest-simulator.org).

References

- [1] M. F. Iacaruso, I. T. Gasler, and S. B. Hofer, Nature 547(7664), 449 (2017).
- [2] V. Braitenberg and A. Schüz, *Cortex: Statistics and Geometry of Neuronal Connectivity* (Springer-Verlag-Verlag, Berlin, 1998), 2nd ed., ISBN 3-540-63816-4.
- [3] H. Markram, E. Muller, S. Ramaswamy, M. W. Reimann, M. Abdellah, C. A. Sanchez, A. Ailamaki, L. Alonso-Nanclares, N. Antille, S. Arsever, *et al.*, Cell 163(2), 456 (2015).
- [4] R. J. Douglas and K. A. C. Martin, Neuron 56, 226 (2007).
- [5] R. J. Douglas and K. A. C. Martin, Curr. Biol. 17(13), 496 (2007).
- [6] T. Binzegger, R. J. Douglas, and K. A. C. Martin, J. Neurosci. 39(24), 8441 (2004).
- [7] A. M. Thomson, D. C. West, Y. Wang, and A. Bannister, Cereb. Cortex 12, 936 (2002).
- [8] T. C. Potjans and M. Diesmann, Cereb. Cortex 24(3), 785 (2014).
- [9] C. van Vreeswijk and H. Sompolinsky, Science 274, 1724 (1996).
- [10] N. Brunel, J. Comput. Neurosci. 8(3), 183 (2000).
- [11] J. W. Tukey, Exploratory data analysis (Addison-Wesley, 1977).
- [12] A. M. Thomson and C. Lamy, Front. Neurosci. 1(1), 19 (2007).
- [13] A. Von Stein and J. Sarnthein, Int. J. Psychophysiol. 38(3), 301 (2000).
- [14] M. Siegel, T. H. Donner, and A. K. Engel, Nat. Rev. Neurosci. 13(2), 121 (2012).
- [15] M. P. Van Den Heuvel and H. E. H. Pol, Eur. Neuropsychopharmacol. 20(8), 519 (2010).
- [16] H. Lindén, T. Tetzlaff, T. C. Potjans, K. H. Pettersen, S. Grün, M. Diesmann, and G. T. Einevoll, Neuron 72(5), 859 (2011).
- [17] S. J. van Albada, M. Helias, and M. Diesmann, PLOS Comput. Biol. 11(9), e1004490 (2015).
- [18] A. Morrison, C. Mehring, T. Geisel, A. Aertsen, and M. Diesmann, Neural Comput. 17(8), 1776 (2005).
- [19] A. Morrison and M. Diesmann, in *Lectures in Supercomputational Neurosciences: Dynamics in Complex Brain Networks*, edited by P. b. Graben, C. Zhou, M. Thiel, and J. Kurths (Springer, Berlin, Heidelberg, 2008), pp. 267–278, ISBN 978-3-540-73159-7, URL http://dx.doi.org/10.1007/978-3-540-73159-7_10.
- [20] Message Passing Interface Forum, *MPI: A Message-Passing Interface Standard*, Tech. Rep. UT-CS-94-230 (1994).

- [21] M. Diesmann, M.-O. Gewaltig, and A. Aertsen, SYNOD: an Environment for Neural Systems Simulations. Language Interface and Tutorial, Tech. Rep. GC-AA-/95-3, Weizmann Institute of Science, The Grodetsky Center for Research of Higher Brain Functions, Israel (1995).
- [22] H. E. Plesser, J. M. Eppler, A. Morrison, M. Diesmann, and M.-O. Gewaltig, in *Euro-Par* 2007: Parallel Processing, edited by A.-M. Kermarrec, L. Bougé, and T. Priol (Springer-Verlag, Berlin, 2007), vol. 4641 of *Lecture Notes in Computer Science*, pp. 672–681.
- [23] T. Ippen, J. M. Eppler, H. E. Plesser, and M. Diesmann, Frontiers in Neuroinformatics 11, 30 (2017).
- [24] M. Diesmann, BioSupercomputing Newsletter 8, 8 (2013).
- [25] S. Kunkel, T. C. Potjans, J. M. Eppler, H. E. Plesser, A. Morrison, and M. Diesmann, Front. Neuroinformatics **5**, 35 (2012).
- [26] M. Helias, S. Kunkel, G. Masumoto, J. Igarashi, J. M. Eppler, S. Ishii, T. Fukai, A. Morrison, and M. Diesmann, Front. Neuroinformatics 6, 26 (2012).
- [27] S. Kunkel, M. Schmidt, J. M. Eppler, G. Masumoto, J. Igarashi, S. Ishii, T. Fukai, A. Morrison, M. Diesmann, and M. Helias, Front. Neuroinformatics 8, 78 (2014), ISSN 1662-5196.
- [28] RIKEN BSI, *Largest neuronal network simulation achieved using K computer*, Press release (2013).
- [29] J. Jordan, T. Ippen, M. Helias, I. Kitayama, S. Mitsuhisa, J. Igarashi, M. Diesmann, and S. Kunkel, Front. Neuroinformatics (under review) (2017).
- [30] S. J. van Albada, A. G. Rowley, J. Senk, M. Hopkins, M. Schmidt, A. B. Stokes, D. R. Lester, M. Diesmann, and S. B. Furber, (under review) (2017).
- [31] M. I. Sereno and R. B. Tootell, Curr. Opin. Neurobiol. 15(2), 135 (2005).
- [32] A. Goulas, M. Bastiani, G. Bezgin, H. B. Uylings, A. Roebroeck, and P. Stiers, PLOS Comput. Biol. 10(3), e1003529 (2014).
- [33] M. Hinne, A. Meijers, R. Bakker, P. H. Tiesinga, M. Mørup, and M. A. van Gerven, PLOS Comput. Biol. 13(1), e1005374 (2017).
- [34] P. Tiesinga, R. Bakker, S. Hill, and J. G. Bjaalie, Curr. Opin. Neurobiol. 32, 107 (2015).
- [35] D. J. Felleman and D. C. Van Essen, Cereb. Cortex 1, 1 (1991).
- [36] M. Schmidt, R. Bakker, C. C. Hilgetag, M. Diesmann, and S. J. van Albada, Brain Struct. Func. (2017).
- [37] S. Dombrowski, C. Hilgetag, and H. Barbas, Cereb. Cortex 11(10), 975 (2001).
- [38] K. Stephan, L. Kamper, A. Bozkurt, G. Burns, M. Young, and R. Kötter, Phil. Trans. R. Soc. B 356, 1159 (2001).
- [39] R. Bakker, W. Thomas, and M. Diesmann, Front. Neuroinformatics 6, 30 (2012).
- [40] N. T. Markov, M. M. Ercsey-Ravasz, A. R. Ribeiro Gomes, C. Lamy, L. Magrou, J. Vezoli, P. Misery, A. Falchier, R. Quilodran, M. A. Gariel, J. Sallet, R. Gamanut, *et al.*, Cereb. Cortex 24(1), 17 (2014).
- [41] M. Ercsey-Ravasz, N. T. Markov, C. Lamy, D. C. V. Essen, K. Knoblauch, Z. Toroczkai, and H. Kennedy, Neuron 80(1), 184 (2013), ISSN 0896-6273.
- [42] N. T. Markov, J. Vezoli, P. Chameau, A. Falchier, R. Quilodran, C. Huissoud, C. Lamy, P. Misery, P. Giroud, S. Ullman, P. Barone, C. Dehay, *et al.*, J. Comp. Neurol. **522**(1), 225 (2014), ISSN 1096-9861.
- [43] J. Schuecker, M. Schmidt, S. J. van Albada, M. Diesmann, and M. Helias, PLOS Comput. Biol. 13(2), 1 (2017).
- [44] W. A. Suzuki and D. G. Amaral, J. Neurosci. 14(3), 1856 (1994).
- [45] K. S. Rockland and D. N. Pandya, Brain Res. 179, 3 (1979).
- [46] C. L. Barnes and D. N. Pandya, J. Comp. Neurol. 318(2), 222 (1992).
- [47] C. Distler, D. Boussaoud, R. Desimone, and L. G. Ungerleider, J. Comp. Neurol. 334(1), 125 (1993).
- [48] E. Jones, J. Coulter, and S. Hendry, J. Comp. Neurol. 181(2), 291 (1978).
- [49] A. Morel and J. Bullier, Vis. Neurosci. 4(06), 555 (1990).
- [50] M. Webster, L. Ungerleider, and J. Bachevalier, J. Neurosci. 11(4), 1095 (1991).
- [51] W. L. Suzuki and D. G. Amaral, J. Comp. Neurol. 350(4), 497 (1994).
- [52] M. J. Webster, J. Bachevalier, and L. G. Ungerleider, Cereb. Cortex 4(5), 470 (1994).
- [53] D. J. Perkel, J. Bullier, and H. Kennedy, J. Comp. Neurol. 253(3), 374 (1986).
- [54] B. Seltzer and D. N. Pandya, J. Comp. Neurol. 343(3), 445 (1994).
- [55] C. L. Rees, K. Moradi, and G. A. Ascoli, Trends Neurosci. (2016).
- [56] C. C. J. Chu, P. F. Chien, and C. P. Hung, Vision Res. 96, 113 (2014).
- [57] C. C. J. Chu, P. F. Chien, and C. P. Hung, CRCNS.org (2014), URL http://dx.doi. org/10.6080/K0J1012K.
- [58] H. Shimazaki and S. Shinomoto, J. Comput. Neurosci. 29(1), 171 (2010), ISSN 1573-6873.
- [59] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, J. Stat. Mech. Theory Exp. 2008(10), P10008 (2008).
- [60] T. Kamada and S. Kawai, Inf. Process. Lett. **31**(1), 7 (1989).
- [61] M. Rosvall and C. T. Bergstrom, PLOS ONE 5(1) (2010).
- [62] B. Biswal, F. Zerrin Yetkin, V. M. Haughton, and J. S. Hyde, Magn. Res. Med. 34(4), 537 (1995).
- [63] M. L. Schölvinck, A. Maier, Q. Y. Frank, J. H. Duyn, and D. A. Leopold, Proc. Nat. Acad. Sci. USA 107(22), 10238 (2010).
- [64] M. Rosvall, D. Axelsson, and C. T. Bergstrom, Eur. Phys. J. Spec. Top. 178(1), 13 (2009).
- [65] S. Song, P. Sjöström, M. Reigl, S. Nelson, and D. Chklovskii, PLOS Biol. 3(3), e68 (2005).
- [66] R. Perin, T. K. Berger, and H. Markram, Proc. Nat. Acad. Sci. USA 108(13), 5419 (2011).
- [67] M. Schmidt, R. Bakker, K. Shen, G. Bezgin, C.-C. Hilgetag, M. Diesmann, and S. J. van Albada, arXiv preprint arXiv:1511.09364v4 (2016).
- [68] N. Wagatsuma, T. C. Potjans, M. Diesmann, and T. Fukai, Front. Comput. Neurosci. 5, 31 (2011).
- [69] E. Hagen, D. Dahmen, M. L. Stavrinou, H. Lindén, T. Tetzlaff, S. J. van Albada, S. Grün, M. Diesmann, and G. T. Einevoll, Cereb. Cortex (2016).
- [70] N. Cain, R. Iyer, C. Koch, and S. Mihalas, PLOS Comput. Biol. 12(9), e1005045 (2016).
- [71] J. H. Lee, C. Koch, and S. Mihalas, Front. Comput. Neurosci. 11(28) (2017).
- [72] T. Schwalger, M. Deger, and W. Gerstner, PLOS Comput. Biol. 13(4), e1005507 (2017).
- [73] W. Potjans, A. Morrison, and M. Diesmann, Front. Comput. Neurosci. 4(141) (2010).
- [74] J. Hahne, M. Helias, S. Kunkel, J. Igarashi, M. Bolten, A. Frommer, and M. Diesmann, Front. Neuroinformatics **9**(22) (2015).
- [75] J. Hahne, D. Dahmen, J. Schuecker, A. Frommer, M. Bolten, M. Helias, and M. Diesmann, Front. Neuroinformatics 11, 34 (2017), ISSN 1662-5196.
- [76] M. Djurfeldt, J. Hjorth, J. M. Eppler, N. Dudani, M. Helias, T. C. Potjans, U. S. Bhalla, M. Diesmann, J. Hellgren Kotaleski, and O. Ekeberg, Neuroinformatics 8, 43 (2010).

F 1 How Physics Shaped our Senses

U. B. Kaupp and L. Alvarez Molecular Sensory Systems Center of Advanced European Studies and Research Bonn, Germany

Contents

1	Introd	luction	2	
2	Photo	reception - Detecting Single Photons		
	2.1	Morphology		
	2.2	High quantum yield of light absorption by rhodopsin	4	
	2.3	Highly amplified cascade of biochemical reactions	5	
	2.4	Thermal stability of rhodopsins	6	
	2.5	Uniform electrical response	7	
	2.6	Fluctuations of hydrolysis rate and electrical noise		
3	Chem	oreception – Detecting Single Molecules	9	
	3.1	The precision of concentration measurements by cells		
	3.2	Single-molecule sensitivity of sperm	11	
	3.3	The sensory and visual signalling pathways: a comparison	11	
R	eferenc	es	12	

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Our senses command over exquisite sensitivity and precision. Visual cells in the eye can detect single photons; sensory neurons in the antennae of insects can detect single pheromone molecules; sperm can detect single molecules of attractants that guide sperm to the egg for fertilization; some animals can detect temperature changes as small as about 0.01 0 C, and hair cells in the inner ear detect vibrations at the limit of Brownian motion. Many of these sensory feats happen in specialized cellular compartments, called cilia or flagella: long, slender filaments that emanate from the cell body (Figure 1). These compartments serve as antennae that register chemical of physical cues and transform the stimulus into an electrical signal that eventually evokes a neuronal and behavioural response. The detectability of weak stimuli is limited by noise from extrinsic and intrinsic sources. The extrinsic noise results from the stochastic nature of quantal stimuli – photons or molecules – impinging onto sensory cells. The intrinsic noise results from thermal fluctuations within cellular signalling pathways that translate or transduce the stimulus into an electrical signal. The impact of noise on the accuracy of sensing has been discussed in general terms from different perspectives [1-3]. The ultimate goal is to reveal what single cells actually do, which requires identifying and characterizing every probabilistic chemical step in precise quantitative terms. In my lecture, I will outline the physical principles and cellular mechanisms underlying such ultra-sensitivity for two sensory modalities: photoreception and chemoreception. Compared to other sensory modalities, the molecules involved in the visual and chemical senses have been identified and the mechanisms have been disclosed with great clarity.



Fig. 1: Scheme of the morphology and size scale of four different sensory cells. The cell body hosts all the housekeeping organelles like the nucleus. A specialized cell compartment, called cilium or flagellum, serves as a sensory antenna to detect a chemical or physical stimulus.

2 Photoreception - Detecting Single Photons

Rod photoreceptors in the retina, across a wide variety of species, are the most sensitive light detectors in nature. Rods can register single photons and translate the absorption event into a minuscule voltage pulse of about 1 mV [4]. Thus, rod photoreceptors operate at the ultimate physical limit. Their exquisite sensitivity rests on a number of structural and mechanistic features described below.

2.1 Morphology

Key among these features is a peculiar cellular morphology (Figure 1). The outer segment of mammalian rods hosts about 1000 flattened membrane sacs, called disks that accommodate about 10^8 molecules (mouse) of the visual pigment rhodopsin (Rh). The density is about 25,000 Rh/µm²; at such density, more than 50% of the disk surface is covered by the light receptor. The optical density $E(\lambda)$ along the long rod axis, according to Lambert-Beer's law:

$$E(\lambda) = -\ln\left(\frac{I}{I_0}\right) = c \cdot \varepsilon(\lambda) \cdot l \tag{1}$$

is almost unity; thus a photoreceptor cell captures most incident photons. The stack of disks serves two additional functions: it organizes and compartmentalizes the cellular signalling events.

First, light triggers a multi-stage sequence of biochemical reactions ranging from transformation of Rh into its active form (Rh*) to degradation of the intracellular messenger cGMP (Figure 2). The signalling molecules involved in this enzyme cascade are tethered to the flat disk surface and can meet by 2D rotational and translational diffusion on the disk surface, a mechanism called *collision coupling*. Furthermore, Rh forms long threads of dimers; the threads may form special pairs - called tracks (Figure 3) [5, 6]. The signalling components downstream of Rh* may assemble on these tracks and molecules could meet for activation by hopping along these tracks in 1D. Reducing a 3D to a 2D or 1D process enhances the probability that molecules meet in the correct geometric arrangement for successful activation

Second, the slender rod outer segment represents a small compartment where the biochemical reactions happen. The stack of disks further compartmentalizes the reaction volume. The radial diffusion of cellular messengers like cGMP or Ca²⁺ is probably unrestricted, whereas the longitudinal diffusion along the length of the outer segment is severely restricted due to the baffling effect of disks. Therefore, the apparent or effective diffusion coefficient of cGMP is about 20 times smaller than in free solution [7]. Consequently, the longitudinal spread of biochemical and active electrical excitation is restricted to about 3 μ m [7] (compare with the length of the mouse outer segment of 24 μ m).





(a) Sensory signalling pathways in (i) rod photoreceptors (ii) olfactory neurons, and Fig. 2: (iii) sea urchin sperm. The signalling pathways begin with a receptor to register the respective stimulus. In rods and olfactory neurons members of the large family of G protein-coupled receptors (GPCR) serve as light and odorant receptors, respectively. In sea urchin sperm a guanylate cyclase serves both as chemoreceptor that binds the ligand and as enzyme that synthesizes the cellular messenger cGMP. The GPCRs activate a G protein (called G_t in rods and G_{olf} in olfactory cells). In turn, the G_t activates a phosphodiesterase (PDE) that removes cGMP by hydrolysis. As cGMP concentration declines, cvclic nucleotide-gated ion channels (CNG) close and the cell hyperpolarizes. The G_{olf} in olfactory cells activates an adenylate cyclase (AC) that synthesizes cAMP; cAMP-gated ion channels (CNG) open and Ca^{2+} flows into the cell. Finally, the Ca^{2+} influx opens Cl channels (Cl) and Cl ions leave the cell producing a large depolarizing response. In sperm, the cGMP synthesis opens a K^+ selective CNG channel (CNGK); the cell hyperpolarizes. During recovery from hyperpolarization, Ca^{2+} channels (CatSper) open and the intracellular Ca^{2+} concentration rises. This rise in Ca^{2+} concentration results in changes of the flagellar beat and motility. Cvclic AMP, cGMP, and the sperm chemoattractant are depicted as vellow, red, and black circles, respectively.

(b) Mechanism of electrical excitation in rod photoreceptors. In the dark, a sustained current enters the outer segment via CNG channels and leaves the inner segment via a K^+ channel. This circulating current, called dark current, is suppressed by light due to closure of CNG channels. Rh in disks is represented by a red dot and light by a green flash.

2.2 High quantum yield of light absorption by rhodopsin

Visual pigments quite generally consist of two parts: a small organic molecule, *11-cis* retinal, and a protein. The *11-cis* retinal is covalently attached to and embedded into the protein opsin, thus forming rhodopsin. Absorption of light triggers the photo-isomerization from the *11-cis* to *all-trans*. The quantum yield of the photochemical reaction is exquisitely high (0.75) and the isomerization is completed within 200 femtoseconds [8, 9]. By comparison, the quantum yield of the free retinal in solution is low (about 0.05, depending on the solvent, viscosity, and wavelength of excitation) and the photo-isomerization proceeds about 10,000-fold slower

(nanosecond range). Thus, by wrapping up retinal with the protein opsin, a small organic molecule becomes an excellent chromophore and a passive dye turns into a reliable trigger of a photochemical reaction. In summary, the high probability of photon absorption in combination with a high quantum yield ensures that almost every photon that hits a rod photoreceptor elicits a physiological response.



Fig. 3: Supramolecular organization of rhodopsin in the disk membrane. Rh (red spheres) forms dimers that form elongated tracks. Some evidence suggests that components of the signalling pathway assemble on the Rh tracks (green and yellow spheres).

2.3 Highly amplified cascade of biochemical reactions

The initial light absorption is highly amplified by an enzyme cascade encompassing several stages (Figure 2a). Light-activated mammalian Rh* can activate about 10-20 G proteins. One G^* protein activates an enzyme called phosphodiesterase (PDE^{*}) that rapidly destroys about 100 cGMP molecules by hydrolysis, thus, capture of a single photon initiates the removal of about 2000 cGMP molecules [10, 11]. Ion channels in the plasma membrane, called cyclic nucleotide-gated (CNG) channels, are kept open when bound to cGMP in the dark; when the cGMP concentration is lowered by light-driven hydrolysis, about 300 channels close and the membrane potential becomes more negative - called hyperpolarization (Figure 2b). A set of different biochemical reactions inactivates Rh*, G*, and the PDE*. Finally, the cGMP pool is replenished by guanylate cyclase (GC), an enzyme that synthesizes cGMP from GTP (Figure 4). The GC activity is inhibited by Ca^{2+} . In the dark GC is inactive because a small fraction of the current passing through CNG channels is carried by Ca²⁺ ions and Ca²⁺ is elevated. When CNG channels close in light, the Ca²⁺ concentration drops and GC starts cGMP synthesis. Finally, CNG channels open again, Ca^{2+} flows into the cell, and cGMP synthesis is terminated. Thus, the two cellular messengers control each other by a perfect negative feedback mechanism. The Ca²⁺ feedback onto cGMP synthesis is important as this mechanism contributes to the uniform waveform of the photo-response and, thereby, allows photoreceptors to discriminate single-photon events.

2.4 Thermal stability of rhodopsins

In complete darkness, spontaneous discrete current fluctuations (*dark events*) can be recorded from rod photoreceptors [12, 13]. The spontaneous fluctuations are indistinguishable from those evoked by single photons, suggesting that they result from thermal activation of visual pigments. The spontaneous activity can be considered as dark noise that interferes with light detection and limits sensitivity. This phenomenon was recognized 150 years ago by Theodor Fechner and by Hermann von Helmholtz and referred to as *Eigengrau* or *Eigenlicht* [14]. Depending on species, the thermally evoked responses occur every other 30 - 120 s/per rod cell.



Fig. 4: Feedback cycle between light-driven hydrolysis and Ca^{2+} -regulated synthesis of cGMP. GCAP is a small Ca^{2+} -binding protein that regulates cGMP synthesis by guanylate cyclase. A $Na^+/Ca^{2+}/K^+$ exchanger extrudes Ca^{2+} from the cell. In the dark, Ca^{2+} influx via CNG channels is in balance with Ca^{2+} extrusion by the exchanger. When CNG channels close, Ca^{2+} drops and cGMP synthesis commences.

Considering the high Rh density in retinal rods, a rhodopsin molecule is activated about once every 300-1000 years. Thus, rhodopsin is thermally extremely stable. Other members of the family of G protein-coupled receptors (GPCR) are not as stable, but also by far not as abundant as rhodopsin [15].

The activation energies E_a of visual pigments by light and heat have been estimated to be 40-50 kcal/mol and 20-25 kcal/mol, respectively [16]. The discrepancy suggested that the two events follow two different molecular routes and activation by light and heat may not be identical. The apparent discrepancy of activation energies is not merely of academic interest. Barlow [17] proposed that long-wavelength absorbing pigments are noisier. If the groundstate energy barriers for isomerization by light and heat are different, the relation between the peak absorption wavelength λ_{max} and thermal isomerization would remain enigmatic. The thermal activation energy was initially derived from the temperature dependence of the rate of dark events (Arrhenius analysis) according to a simple Boltzmann distribution of energies. However, for complex molecules with many vibrational modes, a Boltzmann energy distribution is inappropriate [16, 18] and must be replaced by a Hinshelwood distribution that considers the contribution of vibrational modes to the free internal energy. A molecule's probability of possessing a thermal energy $\geq E_a^T$ and thus being able to become activated thermally is:

$$f = e^{\frac{-E_a^T}{RT}} \sum_{i=1}^m \frac{1}{(i-1)!} \left(\frac{E_a^T}{RT}\right)^{i-1}$$
(2)

wherein R is the universal gas constant, T is the absolute temperature, and m is the number of molecular vibrational modes that contribute to the thermal energy for activation [18, 19]. Using the Hinshelwood distribution, the discrepancy disappears, suggesting that heat and photons activate the same pathway along the reaction coordinate and firmly establish the relationship between λ_{max} and E_a .

The mammalian retina is furnished with two different types of photoreceptors - rods and cones – that are specialized for low and high light levels prevailing at night and during the day, respectively. Compared to rhodopsin in rods (500 nm), visual pigments in cone photoreceptors maximally absorb light at higher wavelengths (up to 620 nm). The spontaneous activity of cone pigments and, therefore, the dark noise level is much higher compared to rods. Therefore, these photoreceptors are not suitable for detecting and counting single photons. The relative thermal instability probably is the reason why far-red-or infrared-detecting animals like snakes use special thermosensors (TRP ion channels) rather than far-red-absorbing pigments [18].

2.5 Uniform electrical response

A remarkable feature of the single-molecule response is its reproducibility and uniformity (Figure 5a). The response amplitude varies over a surprisingly small range [4]. Why was this finding intuitively unexpected? The number of cGMP molecules that become removed by hydrolysis depends on the lifetime of active Rh*. During its lifetime, Rh* activates several Gt molecules that in turn activate several PDE molecules. The longer the Rh* lifetime, the more cGMP molecules are removed and the larger the response amplitude. If the probability of Rh* inactivation would be constant in time, it would be expected that the lifetime is exponentially distributed (like radioactive decay times of isotopes), and the coefficient-of-variation (CV = s.d./mean) of response amplitudes would become unity. Contrary to this expectation, the CV of the light response is 0.25 - 0.35. The small CV lowers photon noise and allows rod photoreceptors to discern individual absorption events. Up to this date, the mechanisms underlying this remarkable feat, are not entirely clear and a matter of debate [7]. Three different mechanisms have been entertained: negative Ca2+ feedback that blunts the photoresponse, saturation, and control of Rh* lifetime [20]. The spatio-temporal dynamics of cGMP also control Ca²⁺ dynamics that terminates the photoresponse and initiates recovery. Simulation of messenger dynamics suggests that this feedback contributes to the uniformity of single-photon responses. Alternatively, two different saturation mechanisms could render responses uniform. First, the cGMP pool in the compartment between two disks is completely exhausted and all CNG channels in this and a few neighbouring compartments close. Alternatively, a response unit is defined by the length of rhodopsin tracks and the number of G_t proteins that are associated with that track (see Figure 3). Given a Rh : G_t ratio of 10 : 1, a 100 nm track built from about 100 Rh molecules could host 10 Gt proteins. If Rh* can activate only G_t proteins hopping on its own track, but not G proteins from neighbouring tracks, the uniformity of the photoresponse is determined by the distribution of G_t proteins among tracks [5].

Finally, a particular interesting mechanism for generating uniform responses is control of Rh* lifetime. If Rh* inactivates stochastically - similar to radioactive decay - the lifetime is exponentially distributed with a CV = 1. Inactivation in multiple steps *n* can narrow the width of the lifetime distribution by $1/\sqrt{n}$ (Figure 5b). In fact, an enzyme, called *rhodopsin kinase*, can attach up to six phosphate groups to serine or threonine amino-acid residues on Rh* ("phosphorylation"). Phosphorylated Rh* is recognized by another molecule, called *arrestin*, that caps phospho-Rh* and removes it from the reaction cycle. Manipulating the number of phosphorylated sites, in fact, affects uniformity of the photoresponse [21]. Activation and inactivation of proteins by phosphorylation/dephosphorylation is a recurrent cellular motif; therefore, defeating shot noise by multiple phosphorylation (or other biochemical modifications) might be a ubiquitous mechanism.



Fig. 5: (A) single-photon responses recorded with a suction pipette from the outer segment of a rod photoreceptor (upper). Superposition of several photoresponses illustrates the uniformity of discrete events (lower). (B) left: Three different mechanisms of rhodopsin shut-off. right: Distribution of integrated rhodopsin activity for the respective shut-off mechanism.

2.6 Fluctuations of hydrolysis rate and electrical noise

Apart from the discrete fluctuations produced by thermal Rh activation, there is continuous noise of the dark current produced by thermal activation (and deactivation) of PDE [13, 22, 23]. This continuous noise compromises the signal-to-noise ratio of the single-photon response.

There is also noise due to the stochastic opening and closing of CNG channels that mediate the light response. Compared to the other two noise sources it is small, and its frequency range is much higher [22]. However, a back-of-the-envelope calculation illustrates that nature takes many measures to keep all noise sources at bay. The dark current flowing from the outer to the inner segment is about 40 pA (Figure 6). Currents carried by single prototypical Na⁺ or K⁺ channels are of the order of one pA. If such channels carry the dark current, on average N = 40 open channels would be required. Channels rapidly fluctuate between open and closed state, giving rise to current noise $\delta I = i \sqrt{N} = 6.3 \text{ pA}$ (*i* is the single-channel current). A single photon suppresses the dark current by about 3% or 1 pA; this minuscule change would be buried in the 6-pA noise. CNG channels are exceptional; their single-channel current is only 4 fA, i.e. about 200-fold smaller than that of a prototypical channel. To produce 40 pA of dark current, on average 10,000 channels must be in the open state; the current noise due to channel fluctuations then becomes $\delta I = i \sqrt{N} = 4 \text{ fAx} \sqrt{10,000} = 0.4 \text{ pA}$.



Fig. 6: Current noise produced by stochastic opening and closing events of ion channels. The single-channel current i is assumed to be 1 pA. The current carried by Na^+ ions is blocked by Ca^{2+} ions; therefore, in the presence of external Ca^{2+} , i is estimated to be about 4 fA.

3 Chemoreception – Detecting Single Molecules

The detection of chemical cues in the environment - which provide information on food, mates, danger, predators, and pathogens - is essential for the survival of most animals. Organisms employ specialized chemosensory cells furnished with intricate cellular pathways to map the chemical landscapes in their environment. Some chemosensory cells can register ligands at extremely low concentrations: Sensory neurons in the vomeronasal organ of the mouse respond to pheromone concentrations of <1 pM [24, 25]. Sperm of the sea urchin *Arbacia punctulata* respond to chemoattractant molecules at femtomolar concentrations [26, 27]. Finally, sensillae in the antennae of the male silkmoth *Bombyx mori* can also register femtomolar concentrations of the female pheromone [28]. These chemosensory cells are equipped with a thin cilium or flagellum that is even thinner than the outer segment of rod photoreceptors (diameter ca. 0.25 μ m vs. 1-8 μ m). The volume of a flagellum (diameter 0.25 μ m, length 45 μ m) is about 2ft. To illustrate the consequences of such a tiny volume, it's worth pointing out that one molecule is equivalent to a concentration of ca. 1 nM. Thus, like in photoreceptors, the cellular signalling events evoked by ligands also happen in a small

reaction compartment.

The evidence for single-molecule detection is not as direct and obvious as for single-photon sensitivity. Nonetheless, for sperm, three lines of evidence argue for single-molecule sensitivity. In experiments with suspensions of sea urchin sperm, the delivery on average of m = 0.2 chemoattractant molecules/sperm cell evokes voltage and Ca²⁺ responses [26, 27]. According to the Poisson distribution, the probability that a cell receives two or more hits is negligible (<1.6 %). Thus, the macroscopic responses probably represent the ensemble average of single-molecule events. Moreover, for mean values m < 1, the shape of the Ca²⁺ and voltage responses becomes independent of the chemoattractant concentrations, i.e. responses scale [27]. Finally, the number of maximal binding events can be estimated from the number of molecules impinging on a cilium. The shape of a cilium can be approximated by an ellipsoid of revolution (*prolate*). The rate of ligand molecules hitting the prolate surface is:

$$K = \frac{4\pi a N_{\rm A} D}{\ln(2a/b)} \tag{3}$$

Wherein *D* is the diffusion constant of the ligand, and *a* and *b* are the long and short semiaxes, respectively, of the prolate, and N_A is Avogadro's number. For a concentration of 1 pM and $a = 25 \mu m$, $b = 0.15 \mu m$, about seven molecules/s hit the prolate surface [36, 38]. Thus, apart from geometric factors that account for different length scales and cell shapes, responses evoked by sub-picomolar ligand concentrations probably reflect single-molecule events.

3.1 The precision of concentration measurements by cells

Organisms and cells use chemical gradients for directed navigation. In a landmark paper [29], derived from first physical principles the limit of precision by which cells can register molecules. Chemoattractants randomly impinging on the cell surface cause "molecule noise" that limits the precision of chemical sensing. A measure of the precision of sensing is the uncertainty $\Delta c/c$, i.e., the incremental increase Δc of a concentration c; only if $\Delta c/c$ is sufficiently high, a new binding event can be distinguished from molecule noise. The uncertainty is given by:

$$\Delta c/c = 1/\sqrt{N_R s D \tau c} \tag{4}$$

wherein N_R is the number of the receptors, *s* is the effective radius of the binding site, *D* is the diffusion constant of the chemoattractant, and τ is the sampling time [29]. In the past 40 years, eqn. (4) and the underlying assumptions has been scrutinized (see [3] for a comprehensive and excellent discussion). Berg & Purcell (1977) and most subsequent studies assume that cells integrate binding events over time. However, maximum-likelihood estimates have also been considered [30]. All studies - within a factor of two or so - arrive at the Berg-Purcell limit of chemical sensing.

Eqn. (4) provides important insights. To lower the sensing error, cells must increase the number of measurements, either by increasing the number of receptors $N_{\rm R}$ or the measuring (integration) time τ . The limit has been experimentally tested for two model systems bacteria and the slime mould *Dictyostelium discoideum*. Bacteria respond to a 0.2-0.3% change in receptor occupancy by glutamate (an amino acid) [31, 32]. The number of receptors that bind

the ligand is about 15,000 [33]. Thus, the binding (or dissociation) of about 30-50 molecules evokes a behavioural response. The response "threshold" is about 10 nM glutamate, corresponding to about 10 bound molecules.

In *D. discoideum*, for the shallowest gradient of the chemoattractant cAMP that allows directed movement $(3.3 \cdot 10^{-3} \text{ nM/}\mu\text{m})$, the difference in receptor occupancy between the front and the back of the cell is about 10 ligand molecules ($N_{\rm R}$ is about 25,000 in front and back of a cell 10 µm in length) [34, 35]. In conclusion, bacteria and slime mould operate in the low nanomolar range of chemoattractants using about 10,000-25,000 receptors/cell; they can detect changes in receptor occupancy of 0.1-0.5%, but not single-molecule events.

3.2 Single-molecule sensitivity of sperm

By contrast to bacteria and slime mould, some sperm species operate in the femto- to picomolar range. According to eqn. (4), sperm must be endowed with a much higher receptor density to achieve a tolerable sensing error in the single-molecule regime. This expectation is born out by experiment: *A. punctulata* sperm host at least 300,000 receptors/cell; the receptors are located on a sperm cell's sensory antenna, the flagellum. The sensing error $\Delta c/c$ using the following parameters: $N_{\rm R} = 3 \times 10^5$ receptors [36], $D_{\rm L} = 2.4 \times 10^6$ cm²/s [37], and $\tau = 0.5$ s [27, 37]. The dimension *s* of the binding site is not known; upper and lower limits are the radius of the extracellular GC domain (2.65 nm) or the radius of the chemoattractant peptide (0.65 nm), respectively. At 1 pM resact, the uncertainty $\Delta c/c$ to count molecules is 0.13–0.27 or 13–27%. For an ellipsoid-shaped flagellum (half-axis $a = 25 \ \mu m$ and $b = 0.15 \ \mu m$; [38], on average 3.5 ligand molecules are impinging during the sampling time of 0.5 s; thus, molecules can be counted with an error of about 0.5 molecules. In conclusion, the high receptor density - in principle - endows sperm operating at the physical limit with molecular precision.

3.3 The sensory and visual signalling pathways: a comparison

The chemosensory and visual signalling pathways share several functional features and signalling components (Figure 2). Both sensory cells employ cGMP as cellular messenger and use CNG channels to generate a hyperpolarization; single photons or molecules evoke voltage responses of about 1 mV and 3 mV, respectively, and a duration of 1-2 s. However, also large quantitative differences exist. The geometric volume of the rod outer segment of a mouse (36 μ m³) is 16-fold larger than the volume of the flagellum (2.2 μ m³). The number of receptor molecules is about 300-fold larger and the receptor concentration (molecules /volume) is about 20-fold larger in rods compared to sperm. However, the membrane density of Rh is only about 2-fold larger compared to sperm GC (25,000/µm² versus 10,000/µm²). Rhodopsin activation initiates the hydrolysis of about 2000 cGMP molecules, whereas GC stimulation initiates the synthesis of about 10 cGMP molecules [36]. Thus, the relative changes in cGMP concentration per unit volume differ by about 13-fold. Furthermore, during a single-photon response in rods about 300 CNG channels close, whereas in sperm <10 CNG channels are gated open. The cGMP concentration in sperm at rest is not known, but it must be as low as a few nM, equivalent to a few free cGMP molecules. In rod photoreceptors the basal cGMP concentration is about 1000-fold higher (1 μ M).

In order to achieve a molecular and mechanistic understanding of chemosensation that is as profound as that of photoreception, recording of single-molecule events from single cells are required. Despite efforts over the past 50 years, this has not yet been accomplished. It will be a frontier of biophysical research for the years to come.

References

- [1] W. Bialek, Annu. Rev. Biophys. Biophys. Chem. 16, 455-78 (1987)
- [2] I. Lestas et al., Nature 467, 174-8 (2010)
- [3] P.R. ten Wolde et al., J. Stat. Physics 162, 1395–1424 (2016)
- [4] D.A. Baylor et al., J. Physiol. 288, 613-634 (1979)
- [5] M. Gunkel et al., Structure 23, 628-38 (2015)
- [6] D. Fotiadis et al., Nature 421, 127-128 (2003)
- [7] O.P. Gross et al., Front. Mol. Neurosci. 8, 6 (2015)
- [8] R.W. Schoenleinet al., Science 254, 412-415 (1991)
- [9] A. Warshel, Angew. Chem. Int. Ed. Engl. 53, 10020-31 (2014)
- [10] M.E. Burns and E.N.Pugh, Jr. Physiology (Bethesda) 25, 72-84 (2010)
- [11] V.Y. Arshavsky and M.E. Burns, Cell Logist. 4, e29390 (2014)
- [12] K.-W. Yau et al., Nature 279, 806-807 (1979)
- [13] D.A. Baylor et al., J. Physiol. 309, 591-621 (1980)
- [14] H. Helmholtz, Zeitschrift für Psychologie und Physiologie der Sinnesorgane (1890)
- [15] A. Manglik and B. Kobilka, Curr. Opin. Cell Biol. 27, 136-43 (2014)
- [16] P. Ala-Laurila et al., Biophys. J. 86, 3653-62 (2004)
- [17] H.B. Barlow, Nature 179, 255-256 (1957)
- [18] D.G. Luo et al., Science 332, 1307-12 (2011)
- [19] R.C. St George, J. Gen. Physiol. 35, 495-517 (1952)
- [20] F. Rieke and D.A. Baylor, Rev. Mod. Phys. 70, 1027-1036 (1998)
- [21] T. Doan et al., Science 313, 530-533 (2006)
- [22] F. Rieke and D.A. Baylor, Biophys. J. 71, 2553-72 (1996)
- [23] J. Reingruber et al., Proc. Natl. Acad. Sci. USA 110, 19378-83 (2013)
- [24] T. Leinders-Zufall et al., Science 306, 1033-1037 (2004)
- [25] T. Leinders-Zufall et al., Nature 405, 792-796 (2000)
- [26] T. Strünker et al., Nat. Cell Biol. 8, 1149-1154 (2006)
- [27] U.B. Kaupp et al., Nat. Cell Biol. 5, 109-117 (2003)
- [28] K.E. Kaissling, Annu. Rev. Neurosci. 9, 121-45 (1986)
- [29] H.C. Berg and E.M. Purcell, Biophys. J. 20, 193-219 (1977)
- [30] R.G. Endres and N.S. Wingreen, Phys Rev Lett 103, 158101 (2009)
- [31] V. Sourjik and H.C. Berg, Proc. Natl. Acad. Sci. USA 99, 123-127 (2002)
- [32] J.E. Segall et al., Proc. Natl. Acad. Sci. USA 83, 8987-8991 (1986)
- [33] M. Li and G.L. Hazelbauer, J. Bacteriol. 186, 3687-94 (2004)
- [34] L. Song et al., Eur. J. Cell Biol. 85, 981-989 (2006)
- [35] D. Dormann et al., J. Cell Sci. 114, 2513-23 (2001)
- [36] M. Pichlo et al., J. Cell Biol. 206, 541-57 (2014)
- [37] N.D. Kashikar et al., J. Cell Biol. 198, 1075-91 (2012)
- [38] H.C. Berg, Random Walks in Biology, Princeton University Press (1993)

F 2 Population Genetics and Evolution

J. Krug

Institute for Theoretical Physics University of Cologne

Contents

1	Intr	oduction	2			
2	Mor	Moran model				
	2.1	Selection	4			
	2.2	Drift	5			
	2.3	Fixation	5			
	2.4	Parallel evolution	7			
3	Regimes of evolutionary dynamics 8					
	3.1	Molecular clock	8			
	3.2	Clonal interference	8			
	3.3	Speed of evolution	10			
	3.4	Muller's ratchet and the evolutionary benefits of sex	11			
	3.5	Spatial structure	12			
4	Epistasis and fitness landscapes 13					
	4.1	Pairwise interactions	13			
	4.2	Multiple loci	14			
	4.3	Mutational pathways	14			

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Evolution is all about processes that almost never happen. Daniel C. Dennett [1]

The mechanism of evolution by natural selection was proposed independently by Charles Darwin and Alfred Russell Wallace in 1858, and constitutes to this day the conceptual framework within which the entire mind-boggling diversity of biological phenomena can be organized and, at least in principle, explained. In the famous words of Theodosius Dobzhansky, "*Nothing in biology makes sense except in the light of evolution*" [2]. Unlike, for example, the mechanics of Newton or the electrodynamics of Maxwell, the theory of evolution was not originally phrased in mathematical terms. In fact, a mathematical formulation would not have been possible at the time of Darwin and Wallace, because the nature of heredity was not yet properly understood. Although Gregor Mendel published his laws of particulate inheritance already in 1865, they were largely ignored and rediscovered only around 1900.

At this point the discreteness of the carriers of genetic information was perceived to be in contradiction to the Darwinian view of evolutionary change being accumulated continuously in small steps over long periods of time. In a striking parallel to the controversy about the existence of atoms that overshadowed the early days of statistical mechanics [3], the proponents of Mendelian genetics were faced with the challenge of explaining how discrete random changes in the genes of individual organisms could give rise to a continuous and seemingly deterministic evolution of traits on the level of populations and species.

Not surprisingly, the reconciliation of the two viewpoints required a mathematical formulation of the basic processes through which the genetic composition of a population evolves. This development is known as the modern synthesis of evolutionary biology, which culminated around 1930 in the foundational works of Ronald A. Fisher, John Burdon Sanderson Haldane and Sewall Wright [4, 5, 6, 7]. Wright's summary of the key insight underlying mathematical population genetics makes it clear why this field may be aptly characterized as a "statistical mechanics of genes":

The difficulty seems to be the tendency to overlook the fact that the evolutionary process is concerned, not with individuals, but with the species, an intricate network of living matter, physically continuous in space-time, and with modes of response to external conditions which it appears can be related to the genetics of individuals only as statistical consequences of the latter. [6]

Our understanding of the molecular basis of genetic and biological processes has progressed greatly and undergone multiple revolutionary changes since the early days of population genetics. Nevertheless the theory remains fundamental to the interpretation of genomic data in laboratory experiments as well as in natural populations. In this lecture some key concepts of mathematical population genetics will be introduced within the most elementary setting, and a few recent applications addressing evolution experiments with bacteria and the evolution of antibiotic drug resistance will be described. Correspondingly, the focus will be on asexually reproducing populations; effects of genetic recombination that play a central role in traditional treatments of population genetics will not be covered. Pointers to the literature for further reading are provided throughout the text, and some of the derivations are left as exercises for the reader.



Fig. 1: Illustration of the Moran model for a population of size N = 5. Initially the population consists of three types. After 7 time steps, the blue type has gone extinct and the red type is close to fixation. Note that in the second time step the population does not change, because the same individual has been chosen for reproduction and removal.

2 Moran model

We consider a population consisting of a fixed number N of individuals labeled by an index i = 1, ..., N. Each individual carries a set of hereditary traits which will be collectively referred to as its *type*. Types can change through *mutations* which will however not be explicitly included in the present discussion. The most important trait governing the dynamics of type frequencies in the population is the *fitness* f_i , a real-valued number which quantifies (in a way to be specified shortly) the reproductive success of the individual. In the Moran model individuals reproduce asexually according to the following steps (see Fig. 1) [8, 9, 10]:

- An individual i is chosen for reproduction with a probability proportional to its fitness f_i .
- The chosen individual creates an offspring of the same type. At this point the population size is N + 1.
- To maintain the constraint of fixed population size, an individual is chosen randomly (without reference to its fitness) and killed. The killed individual could be the same as the one that was chosen for reproduction.

We say that N of these reproduction events make up one *generation*. The Moran model is one of two commonly used reproduction schemes in population genetics, the other being the Wright-Fisher model introduced by two of the pioneers mentioned in the Introduction [4, 6]. In the Wright-Fisher model the entire population is replaced by its offspring in a single step. This scheme is advantageous for numerical simulations [11], but the Moran model is more easily tractable by analytic means. For large N and under a suitable rescaling of time the two models are largely equivalent [12, 13].

We now specialize to a situation with only two types, denoted by A and B, to which we assign the fitness values $f_A = 1 + s$ and $f_B = 1$. The parameter s is called the *selection coefficient* of the A-type relative to the B-type. The state of the population is then fully characterized by the number n_A of A-individuals, which we henceforth denote by n; the number of B-individuals is correspondingly equal to N - n. In a single reproduction step the variable n can change by at most ± 1 . Mathematically speaking the time evolution is a Markov chain on the state space $\{0, 1, ..., N\}$ governed by the transition rates

$$T(n+1|n) = \frac{(1+s)n(N-n)}{N^2}, \quad T(n-1|n) = \frac{n(N-n)}{N^2}$$
(1)

and T(n|n) = 1 - T(n+1|n) - T(n-1|n), where T(n'|n) denotes the transition rate from n to n'. Markov chains of this kind are also known as birth-death processes [14]. For later reference we note that by construction T(0|0) = T(N|N) = 1, which implies that n = 0 and n = N are *absorbing states*.

2.1 Selection

We first ask how the average frequency x = n/N of the A-type changes over time. According to (1), the average number of A-individuals changes by $\Delta n = sx(1-x)$ in one reproduction step. Counting time in units of generations we thus arrive at

$$\frac{dx}{dt} = \dot{x} = sx(1-x). \tag{2}$$

This dynamical system has fixed points at x = 0 and x = 1. For s > 0 the fixed point at x = 1 (x = 0) is stable (unstable) and the stability is reversed for s < 0. Thus under the deterministic dynamics (2) the type with higher fitness dominates the population at long times, provided it was at all present initially. This can be seen as a simple mathematical representation of the Darwinian principle of the survival of the fittest.

It is instructive to rewrite Eq. (2) in terms of the mean population fitness $\bar{f} = f_A x + f_B(1-x) = 1 + sx$, which yields

$$\dot{x} = (f_A - \bar{f})x. \tag{3}$$

Thus the frequency of type A grows (shrinks) whenever the mean population fitness is smaller (larger) than f_A . Equation (3) can be naturally generalized to an arbitrary number K of types with fitnesses $f_1, ..., f_K$ and frequencies $x_1, ..., x_K$, where it takes the form

$$\dot{x_k} = (f_k - \bar{f})x_k, \quad \bar{f} = \sum_{k=1}^K f_k x_k.$$
 (4)

This is a simple example of a class of dynamical systems known as *replicator equations* [15].

Exercise: Find the general solution of (4) for arbitrary fitness values f_k . *Hint:* Look for a transformation that eliminates the nonlinearity $\bar{f}x_k$.

As an immediate consequence of (4), the mean fitness of the population evolves as

$$\dot{\bar{f}} = \sum_{k=1}^{K} f_k (f_k - \bar{f}) x_k = \sum_{k=1}^{K} f_k^2 x_k - (\bar{f})^2 = \operatorname{Var}[f] \ge 0,$$
(5)

where $\operatorname{Var}[f]$ is the population variance of the fitness. This is a simple version of a statement known as *Fisher's fundamental theorem* [4]: The mean population fitness increases under selection, and the rate of fitness increase is proportional to the amount of genetic variability in the population. Within the framework of the replicator equations (4) with constant fitness values, selection comes to an end when the fittest type takes over and $\operatorname{Var}[f] \to 0$.

2.2 Drift

If the two types A,B have the same fitness, s = 0, then the frequency of the A-type does not change on average. However, the size of the A-population still changes randomly according to the transition rates (1). These fluctuations, which ultimately arise from the sampling process in a finite population, are referred to as *genetic drift*. They are not visible in the deterministic selection equation (2), which is rigorously valid in the limit $N \to \infty$, because they occur on a different (longer) time scale.

To study drift we need to refine our analysis and consider the full distribution $P_n(t)$ of the random variable n, which evolves according to the master equation [14]

$$P_n(t) = T(n|n+1)P_{n+1}(t) + T(n|n-1)P_{n-1}(t) - [T(n+1|n) + T(n-1|n)]P_n(t).$$
 (6)

Replacing n by Nx and expanding the transition rates and the distribution function in powers of 1/N one arrives at an evolution equation for the distribution of the A-type frequency x which takes the form of a Fokker-Planck equation [14],

$$\frac{\partial}{\partial t}P(x,t) = -\frac{\partial}{\partial x}sx(1-x)P + \frac{1}{N}\frac{\partial^2}{\partial x^2}x(1-x)P$$
(7)

up to corrections of order N^{-2} . The fact that the diffusion term is of order N implies that *drift* phenomena occur on a time scale of N generations. Provided $s \gg 1/N$ the selection term dominates the evolution and one speaks of strong selection. Since selection coefficients are often small, the regime of weak selection where $s \sim 1/N$ is also of importance.

Exercise: Derive Eq. (7) from Eq. (6). Can you find stationary solutions of Eq. (7)? What boundary conditions should you impose at x = 0 and x = 1? Are the solutions normalizable?

2.3 Fixation

We return to the discrete Markov chain governed by the transition rates (1). Since the states n = 0 and n = N are absorbing, for long times the process will reach one of them and subsequently stay there. When this happens we say that *fixation* has occurred at the A-type (if n = N) or the B-type (if n = 0). The probability of fixation at either of the two types depends on the starting value of n. We denote by π_n the probability of fixation at the A-type when there were n A-individuals present initially. This quantity can be computed recursively subject to the obvious boundary conditions

$$\pi_0 = 0, \quad \pi_N = 1.$$
 (8)

To set up the recursion we consider the evolution of the process starting from n. After one time step there are three possibilities:

- The process has moved to n + 1 and subsequent fixation occurs with probability π_{n+1} .
- The process has moved to n-1 and subsequent fixation occurs with probability π_{n-1} .
- The process stays at n and subsequent fixation occurs with probability π_n .

Summing the three possibilities weighted with their respective transition probabilities we arrive at the relation

$$\pi_n = T(n+1|n)\pi_{n+1} + T(n-1|n)\pi_{n-1} + T(n|n)\pi_n \tag{9}$$

which defines a second order recursion relation for π_n . It is worth noting that the terms on the right hand side of (9) are subtly different from those on the right hand side of the master equation (6). This is because the fixation probability is an eigenvector of the *adjoint* of the time evolution operator of the process, which encodes the dynamics backwards in time. The solution of (9) which satisfies the boundary conditions (8) reads

$$\pi_n = \frac{1 - (1 + s)^{-n}}{1 - (1 + s)^{-N}}.$$
(10)

Exercise: Derive Eq. (10) from Eqs. (9) and (8).

We next examine some limiting cases of Eq. (10).

Neutral evolution. When $f_A = f_B$ all types have the same fitness and changes in type frequency arise purely by genetic drift. This situation is referred to as *neutral* evolution [16]. Taking the limit $s \to 0$ in (10) we obtain

$$\pi_n = \frac{n}{N},\tag{11}$$

a result that can be derived more intuitively from the following argument. Suppose that initially every individual is of a distinct type but all types have the same fitness. As time progresses more and more types go extinct until, after a time of the order of the drift time N, a single type remains. This implies that one individual in the initial population is destined to become the common ancestor of the entire population in the far future. Since all individuals are equivalent under neutral evolution, the probability for any given individual to acquire this role is 1/N. For the situation with two types it follows that the probability that the future population will be dominated by the A-type is equal to the fraction of individuals that are initially of type A, i.e. n/N.

Weak selection. Taking the joint limit $N \to \infty, n \to \infty, s \to 0$ in (10) at fixed values of x = n/N and $\bar{s} = sN$ we arrive at the expression

$$\pi(x) = \frac{1 - e^{-\bar{s}x}}{1 - e^{-\bar{s}}} \tag{12}$$

which can also be derived from the backwards (adjoint) version of the Fokker-Planck equation (7) [17].

Strong selection. We specialize to the case of a single initial A-individual which has arisen through a mutation, and take the limit $N \to \infty$ at fixed $s \neq 0$. Then (10) reduces to

$$\pi_1 = \max\left[\frac{s}{1+s}, 0\right] \approx \max[s, 0] \tag{13}$$

independent of N, where in the last step we have assumed that s > 0 is small in absolute terms, $0 < s \ll 1$. The biological significance of Eq. (13) is that newly arising *deleterious* mutations with s < 0 cannot fix in large populations, whereas *beneficial* mutations with s > 0 fix with a probability that is proportional to s. This important result was first derived by Haldane for the Wright-Fisher model, where $\pi_1 \approx 2s$ for $0 < s \ll 1$ [5]. Quite generally one expects that $\pi_1 \sim s$ in this limit with a prefactor that depends on the precise reproduction scheme. Since the selection coefficients of beneficial mutations are rarely larger than a few percent, this implies that, contrary to the scenario suggested by the deterministic selection equation (2), a large fraction of beneficial mutations is lost due to genetic drift.

For later reference it is also of interest to know the time until fixation, conditioned on it taking place. It turns out that, for the case of strong selection, the correct result for large N can be obtained from the deterministic dynamics (2). Starting the process with a single A-individual implies that the initial frequency is $x_{\text{initial}} = 1/N$, and we say that fixation has occurred when $x = x_{\text{final}} = 1 - 1/N$. Thus the time to fixation is given by

$$t_{\rm fix} = \int_{x_{\rm initial}}^{x_{\rm final}} \frac{dx}{sx(1-x)} = \frac{2\ln(N-1)}{s},$$
(14)

which agrees to leading order in N with the full stochastic calculation. The main contributions to the integral come from the boundary layers where the frequency x is close to 0 or 1. In this regime the stochastic dynamics is well approximated by a branching process, which behaves similar to the deterministic evolution when conditioned on survival [10].

2.4 Parallel evolution

As a simple application of the fundamental result (13) we ask for the probability that two populations evolving from the same starting configuration will fix the same mutation. We assume that the populations have access to a repertoire of m beneficial mutations with positive selection coefficients $s_1, s_2, ..., s_m$. Then it is plausible (and can be proved more formally [18]) that the probability q_k that the first mutation to fix is the one with label k is given by the fixation probability $\pi_1(s_k)$ normalized by the sum of all fixation probabilities of beneficial mutations,

$$q_k = \frac{\pi_1(s_k)}{\sum_{i=1}^m \pi_1(s_i)} \approx \frac{s_k}{\sum_{i=1}^m s_i}$$
(15)

under the approximation (13). The probability that the k'th mutation is fixed in two independent populations is q_k^2 . Correspondingly, the probability that any one of the available mutations fixes in both populations is obtained by summing over k, yielding the expression

$$P_2 = \sum_{k=1}^{m} q_k^2 = \sum_{k=1}^{m} \frac{s_k^2}{\left(\sum_{i=1}^{m} s_i\right)^2}$$
(16)

for the probability of parallel evolution [19].

Similar to entropy measures, P_2 quantifies the deviation of the discrete normalized probability distribution q_k from the equidistribution. If all selection coefficients are the same P_2 takes on its minimal value 1/m, and any variation in the selection coefficients increases the probability of parallel evolution. In order to actually compute P_2 one needs to invoke empirical data or make specific assumptions about the distribution of selection coefficients. Assuming that the tail of the fitness distribution conforms to the Gumbel class of extreme value theory [20], which comprises many common distributions like the normal or exponential distributions, Orr showed that P_2 takes the universal value 2/(m+1) [19]. An application to data obtained in the context of antibiotic resistance evolution is shown in Fig. 2. Here P_2 increases systematically with the concentration of antibiotics and considerably exceeds the prediction for the Gumbel class, indicating that the distribution of selection coefficients is heavy-tailed [21]. In addition to P_2 , the figure also shows the probability $P_{\text{max}} = \max_k q_k$ that the mutation of largest effect fixes first.



Fig. 2: Quantification of parallelism in the evolution of antibiotic resistance. Selection coefficients were estimated as a function of antibiotic concentration for a panel of 48 mutations in the enzyme TEM-1 β -lactamase that increase resistance against cefotaxime. Blue crosses show the probability of parallel evolution (16) and red squares the probability P_{max} that the mutation of largest effect size fixes. The horizontal dashed line shows the prediction for fitness values distributed according to the Gumbel class of extreme value theory. Cefotaxime concentration is measured in μ g/ml. Adapted from [21].

Exercise: Prove that $P_{\max} \ge P_2$.

3 Regimes of evolutionary dynamics

We make use of the results derived in the preceding section to describe some simple but biologically relevant regimes of evolutionary dynamics.

3.1 Molecular clock

Suppose that *neutral* mutations arise in the population at rate U_n per individual and generation. The average number of mutations arising in the whole population per generation is then U_nN , a fraction $\pi_1 = 1/N$ of which go to fixation according to Eq. (11). Thus the rate at which nucleotide changes occur in the genetic sequence, referred to as the (neutral) substitution rate

$$\nu_n = \pi_1 U_n N = U_n \tag{17}$$

is independent of population size. This simple observation underlies the molecular clock hypothesis, which posits that molecular evolution at the level of amino acid substitutions in proteins occurs at approximately constant rate across widely different species [22].

3.2 Clonal interference

When a population is not optimally adapted to its environment, a certain fraction of mutations are beneficial. Suppose that such mutations occur at rate U_b per individual and generation,

and that the associated typical selection coefficients satisfy $N^{-1} \ll s \ll 1$. Then the rate of substitution is $\nu_b = \pi_1 U_b N = s U_b N$ which increases linearly in N. Stated differently, the time between substitutions

$$t_{\rm sub} = \frac{1}{\nu_b} = \frac{1}{NsU_b} \tag{18}$$

decreases with increasing population size. On the other hand, the time (14) required for a single fixation event increases logarithmically in N. As a consequence, we need to distinguish two regimes of adaptive evolution that occur at small vs. large population sizes.



Fig. 3: Illustration of the periodic selection regime. Beneficial mutations that overcome genetic drift emerge at the substitution rate $\nu_b = 1/t_{sub}$ and fix in independent selective sweeps.

In the *periodic selection regime* ($t_{sub} \gg t_{fix}$) each mutation has time to fix before the next one appears, and fixation events are independent (Fig. 3). As N increases, it becomes increasingly likely that fixation events overlap, in the sense that a second beneficial mutation arises while the first is still on the way to fixation. Because the second mutation increases the mean fitness of the population, it reduces the selective advantage of the first mutation [compare to (4)]. Through this mechanism multiple beneficial clones that are present simultaneously in the population compete for fixation, a phenomenon known as *clonal interference* [11, 23, 24]. Clonal interference is the dominant mode of evolution when $t_{sub} \ll t_{fix}$ or

$$\frac{t_{\rm fix}}{t_{\rm sub}} = 2NU_b \ln N \gg 1. \tag{19}$$

This condition is often satisfied for populations of bacteria or viruses, which are characterized by large population sizes and large mutation rates. As an example, we consider the famous longterm evolution experiment with *Escherichia coli* that was initiated by Richard Lenski in 1988. The experiment started with 12 genetically identical populations in a minimal growth medium, which allows the bacteria to grow but leaves room for adaptation. Each day the populations are transferred to fresh growth medium following a fixed protocol, and population samples are stored in a freezer every 500 generations. In this experiment the beneficial mutation rate has been estimated to be $U_b \approx 1.7 \times 10^{-6}$, and the population size (averaged over serial transfers) is $N = 3.3 \times 10^7$ [25]. With these numbers the quantity on the right hand side of the criterion (19) is about 2000, and indeed a recent genomic analysis of the first 60000 generations in the Lenski experiment reveals massive clonal interference [26].

3.3 Speed of evolution

Clonal interference implies that beneficial mutations can be outcompeted by others even when they have reached a frequency that is large enough to make genetic drift irrelevant. Thus beneficial mutations are lost that would otherwise have fixed, which reduces the speed at which the fitness of the population increases. If the beneficial mutations cover a range of selection coefficients then those of small effect are more likely to survive competition, and the distribution of fixed mutations is shifted towards larger effect sizes (Fig. 4).



Fig. 4: Signatures of clonal interference in simulations of the Wright-Fisher model with beneficial mutations occurring at rate $U_b = 10^{-6}$. Left: Distribution of selection coefficients of fixed mutations for different population sizes. In these simulations selection coefficients were drawn randomly from an exponential distribution with mean $s_b = 0.02$ [24]. In the periodic selection regime (P.S.) mutations that survive genetic drift fix independently. Since the fixation probability is proportional to s, the distribution of fixed mutations is $\sim se^{-s/s_b}$ For $N \ge 10^5$ clonal interference becomes relevant and the distribution shifts to larger values. Courtesy of Su-Chan Park. Right: Speed of adaptation as a function of population size assuming a single beneficial selection coefficient s = 0.01. The full red line shows Eq. (21), the blue dotted line a related expression derived in [27], and the red squares are simulation results. Adapted from [28].

To quantify the effect of clonal interference on the speed of fitness increase, we assume that all beneficial mutations have the same selection coefficient s, irrespective of when they occur. Then each fixed mutation increases fitness by s, and the rate of fitness increase in the periodic selection regime

$$V = V_{\rm PS} = s\nu_b = s^2 U_b N \tag{20}$$

is proportional to the supply rate U_bN of beneficial mutations. Computing the speed of evolution in the clonal interference regime requires analyzing the complex stochastic dynamics of interacting clones, and no simple closed-form solution exists [11]. An approximate implicit formula for the speed V(N) that covers both regimes is given by [27, 28]

$$\ln N \approx \frac{V}{2s^2} \left[\ln^2 \left(\frac{V}{eU_b s} \right) + 1 \right] + \ln \left(\frac{V}{2s^2 U_b} \right), \tag{21}$$

see Fig. 4 for a comparison to simulation data. For large N the speed increases logarithmically rather than linearly with N, a behavior that persists up to hyperastronomically large population sizes [11].

Exercise: Find the solution V(N) of the implicit equation (21) for small and large N. Note that the result for small N differs slightly from (20), because (21) was derived assuming Wright-Fisher dynamics.

Experiments in which the population size and the beneficial mutation rate are varied systematically are in qualitative agreement with the predictions outlined above [29, 30]. However, the assumption that beneficial mutations of constant effect size emerge at a constant rate, and that therefore the fitness of the population increases linearly in time, is not consistent with long-term experiments. Instead, the fitness increase slows down over time, an effect that is attributed to a decrease of the typical selection coefficient [25]. Since this implies that the mutational effect size depends on the fitness of the population in which it occurs, it constitutes a simple case of *epistasis* (see Sec. 4 for further discussion).

3.4 Muller's ratchet and the evolutionary benefits of sex

The large majority of random mutations is deleterious or lethal [31]. In the preceding subsections we have nevertheless ignored deleterious mutations, because we have seen in Sec. 2.3 that they cannot fix in large populations. However, even if deleterious mutations are constantly removed from the population by selection, their presence reduces the fitness of the population, an effect that is called *mutational load*. To quantify it, we use the deterministic framework developed in Sec. 2.1. Suppose deleterious mutations with selection coefficient $s = -\tilde{s} < 0$ arise at constant rate U_d , and that an individual with k deleterious mutations has fitness $f_k = -\tilde{s}k$. There are no beneficial or neutral mutations. Then the frequency x_k of individuals with k mutations making up the k'th *fitness class* evolves according to [32]

$$\dot{x}_k = -\tilde{s}(k-\bar{k})x_k + U_d x_{k-1} - U_d x_k \text{ with } \bar{k} = \sum_k k x_k.$$
 (22)

This set of equations has a stationary solution \bar{x}_k that takes the form of a Poisson distribution

$$\bar{x}_k = \frac{\Lambda^k \mathrm{e}^{-\Lambda}}{k!} \quad \text{with} \quad \Lambda = \frac{U_d}{\tilde{s}}.$$
 (23)

Remarkably, the mutational load $s\bar{k} = -U_d$ is independent of \tilde{s} .

Exercise: Verify the stationary solution (23) of (22), and show that there is an infinite number of stationary solutions related to (23) by an integer shift of k.

At this point we recall that real populations are finite, and therefore the deterministic description (22) has to be complemented by the effects of genetic drift. If the total population size is N, the number of individuals in the fittest (mutation-free) class is, on average,

$$\bar{n}_0 = \bar{x}_0 N = N \mathrm{e}^{-U_d/\tilde{s}}.$$
 (24)

In fact the number of mutation-free individuals is a random quantity that fluctuations between $n_0 = N$ and $n_0 = 0$. Because beneficial mutations do not occur, the state $n_0 = 0$ is absorbing: Once the mutation-free class has become extinct, it cannot be reconstituted. The time until the (inevitable) extinction of the mutation-free class occurs depends crucially on its mean size (24). If $\bar{n}_0 \gg 1$, we know from the results of Sec. 2.3 that the fixation of the deleterious mutants is exponentially unlikely, and hence the time until extinction is exponentially large in N. Once it occurs, the frequency distribution has ample time to relax back to a shifted Poisson distribution for which the minimal number of mutations is now k = 1. This process repeats itself periodically and is referred to as one *click of the ratchet*. In contrast, when $\bar{n}_0 \sim 1$ extinction events occur rapidly and fitness declines steadily without discernible discrete clicks. The computation of the speed of fitness decline as a function of N, U_d and \tilde{s} is a difficult problem that is the subject of current research [27]. In the regime $\bar{n}_0 \gg 1$ the dynamics is dominated by rare events, and methods from physics based on the WKB approximation of quantum mechanics have proven to be useful [33].

The ratchet-like fitness decline in asexual populations subject to a constant rate of deleterious mutations was first proposed by Hermann Muller as a possible explanation for the evolutionary advantage of sexual reproduction [34]. In sexual populations mutation-free offspring can be created by recombining the genomes of parents which have deleterious mutations at different positions of the genome, and therefore the ratchet can be halted. It is worth pointing out that also the phenomenon of clonal interference described above in Sects. 3.2 and 3.3 was first discussed in the context of comparing sexual and asexual reproduction [4, 35]. In sexual populations clonal interference is alleviated, because beneficial mutations that have originated in different population clones can recombine into a single genome, and therefore sexual reproduction confers an advantage in terms of the speed of adaptation. Theoretical analysis shows that even a small rate of recombination substantially increases the speed [28].

3.5 Spatial structure

All models introduced so far assume that the population is *well-mixed*, in the sense that the competition induced by the constraint of fixed population size acts globally. This assumption is valid when bacteria are grown in shaken liquid culture, but does not apply to the growth of colonies on plates, nor does it generally apply to natural populations. Incorporating an explicit spatial structure where competition between individuals is implemented locally leads to important changes in the scenarios described above. For example, because a clone of beneficial mutants can grow only at the boundary of the region that it inhabits, the expression (14) for the time to fixation is modified. Assuming that the spreading of the clone occurs at a speed proportional to the selection coefficient and that the population density is constant, one obtains [36]

$$t_{\rm fix} \sim \frac{l}{s} \sim \frac{N^{1/d}}{s},\tag{25}$$

where l denotes the linear extension of the system and d the spatial dimension. Thus fixation is considerably slower than in the well-mixed case, and correspondingly clonal interference is enhanced. Moreover, instead of the logarithmic dependence on N discussed in Sec. 3.3 the speed of adaptation attains a finite speed limit for large systems which scales as $V_{\infty} \sim s^2 U_b^{1/(d+1)}$ [36]. Also the dynamics of Muller's ratchet explained in Sec. 3.4 is altered significantly, in that the population fitness declines at a constant rate even in very large populations, provided the rate of deleterious mutations U_d is large enough [37].

4 Epistasis and fitness landscapes

In the discussion of the dynamical regimes in Sect. 3 it was assumed that a mutation can be associated with a selection coefficient that quantifies its effect on fitness independent of the type or *genetic background* of the individual in which it occurs. In most cases this is an oversimplification, and *epistatic interactions* between the effects of different mutations have to be taken into account [38, 39]. In this final section we sketch the mathematical description of epistasis among mutations occurring at multiple genetic loci. This will lead to the important concept of a *fitness landscape*, which defines the arena in which the evolutionary processes described in the preceding sections ultimately take place [7, 40].

4.1 Pairwise interactions

To fix ideas, consider a *wild type* with fitness f_0 and two different mutant types with fitness f_1 and f_2 , respectively. The selection coefficients of the two mutations are $s_1 = f_1 - f_0$ and $s_2 = f_2 - f_0$. If these selection coefficients were independent of the genetic background we would expect that the fitness f_{12} of the double mutant is given by

$$f_{12}^{(0)} = s_1 + f_2 = s_2 + f_1 = f_1 + f_2 - f_0.$$
(26)

Any deviation from this linear relation implies that the two mutations interact epistatically, and the strength and sign of the interactions are quantified by the (pairwise) epistasis coefficient

$$\epsilon_2 = f_{12} - f_{12}^{(0)} = f_{12} + f_0 - f_1 - f_2.$$
(27)

An important qualitative distinction can be made according to whether the sign of the selection coefficient associated with a given mutation depends on the genetic background or not; in the first case one speaks of *sign epistasis*, in the second of *magnitude epistasis* [41]. In the example at hand, the selection coefficient of the second mutation in the background of the first is $s_2^{(1)} = f_{12} - f_1 = s_2 + \epsilon_2$, and hence its sign is affected by the first mutation if $s_2 s_2^{(1)} < 0$ or $\epsilon_2 s_2 < -s_2^2$.



Fig. 5: Different types of epistasis between mutations occurring at two different genetic loci. The panel on the far right shows a case of reciprocal sign epistasis, where the sign of the selection coefficient of the first mutation depends on the presence of the second and vice versa. Note that the corresponding fitness landscape displays two local maxima.

Magnitude and sign epistasis for a pair of mutations is illustrated in Fig. 5. In this figure the four types are encoded by a pair $\sigma = (\sigma_1, \sigma_2)$ of binary variables $\sigma_i \in \{0, 1\}$, where $\sigma_i = 0/1$

implies the absence/presence of the *i*'th mutation. The fitness values of the four types can then be succinctly written in the form

$$f(\sigma_1, \sigma_2) = f_0 + s_1 \sigma_1 + s_2 \sigma_2 + \epsilon_2 \sigma_1 \sigma_2.$$

$$(28)$$

Equation (28) is the simplest example of a fitness landscape which assigns fitness values to a collection of genotypes [40]. In the present case the genotypes consist of two genetic loci i = 1, 2, each of which carries two possible *alleles* $\sigma_i = 0$ or 1. Because the four genotypes are located at the corners of a square, the landscape is easily visualized by arranging the genotypes in a plane and plotting fitness as the third dimension. This property is lost when we extend the description to more than two loci.

4.2 Multiple loci

In the general case of L mutational loci the genotype is described by a sequence of length L, $\sigma = (\sigma_1, \sigma_2, ..., \sigma_L)$. In practice there could be more than two alleles at each site. For example, when describing mutations on the level of DNA sequences each site carries one of 4 nucleotide bases, and in proteins the set of alleles are the 20 amino acids. Here we restrict ourselves to the binary case and let $\sigma_i \in \{0, 1\}$ represent the absence or presence of a specific mutation as before. Then a generic fitness landscape $f(\sigma)$ can be written in the form of a discrete 'Taylor' expansion [39, 42]

$$f(\sigma) = f_0 + \sum_{i=1}^{L} s_i \sigma_i + \sum_{k=2}^{L} \sum_{\{i_1, i_2, \dots, i_k\}} \epsilon_k^{\{i_1, i_2, \dots, i_k\}} \sigma_{i_1} \sigma_{i_2} \cdots \sigma_{i_k},$$
(29)

where the first sum on the right hand side contains the non-epistatic (linear) contributions. The second sum runs over all subsets $\{i_1, i_2, ..., i_k\}$ of $k \ge 2$ out of L loci. The coefficients $\epsilon_k^{\{i_1, i_2, ..., i_k\}}$ generalize the pairwise epistasis coefficient ϵ_2 in (28). As there are $\binom{L}{k}$ coefficients of order L, the total number of coefficients in the expansion (29) is equal to 2^L . Thus the mapping of the 2^L fitness values to the expansion coefficients is one-to-one.

The set of binary sequences of length L can be represented graphically by linking sequences that differ at a single site. The resulting graphs are L-dimensional (hyper)cubes, which are shown for L = 3 and L = 4 in Fig. 6. Unlike the two-dimensional case depicted in Fig. 5, fitness functions defined on hypercubes with $L \ge 3$ cannot be easily plotted. A useful representation that retains partial information about the ordering of fitness values is provided by *fitness graphs*, where the links between genotypes are decorated with arrows pointing in the direction of increasing fitness [43, 44]. For more than 6 loci also this representation becomes unwieldy and any graphical rendering of the fitness landscape is bound to obscure at least some its geometrical structure (Fig. 7).

4.3 Mutational pathways

We are now ready to integrate the elements of evolutionary dynamics developed in Sects. 2 and 3 with the fitness landscape picture. For this, we assign the 2^L genotypes of the fitness landscape to a population of N individuals evolving according to the Moran model. In the periodic selection regime described in Sec. 3.2 a simple dynamics emerges. Most of the time the population is genetically homogeneous and hence occupies a single site in the fitness graph.



Fig. 6: Fitness graphs for L = 3 (left) and L = 4 (right) loci. The graphs are L-dimensional hypercubes, and the arrows on the links point in the direction of increasing fitness. In the left panel the fitness values are additionally indicated by the size of the balls surrounding each node. The right panel shows the experimentally determined fitness values of all combinations of a subset of 4 mutations taken from the 8-dimensional fitness landscape displayed in Fig. 7. Genotypes which constitute local fitness maxima are shown in larger font and underlined. Colored arrows show the steps taken by a greedy dynamics, where the beneficial mutation of largest effect is chosen deterministically [45].

Occasionally beneficial mutations arise and fix, which implies that the population moves to a neighboring genotype of higher fitness. The population thus moves across the fitness graph following the direction of the arrows on the links, and we see that the fitness graph provides a kind of road map of possible evolutionary trajectories.

Within the periodic selection regime the population is constrained to evolve along pathways of monotonically increasing fitness, and such paths have been termed (evolutionarily) accessible [45, 46]. It is of interest to ask how likely accessible pathways are to be found on real fitness landscapes. Consider two binary genotypes that differ at D loci. To evolve from one to the other, mutations have to take place at D sites, and a priori these mutations can occur in any one out of D! orderings. Thus the total number of mutational pathways connecting the two genotypes is D!. In a seminal experimental study, Weinreich and collaborators considered the pathways along which a highly resistant five-fold mutant of the TEM-1 β -lactamase enzyme could evolve [46] (see Fig. 2 for further information on this system). They found that only 18 out of 5! = 120 pathways displayed monotonically increasing resistance, and hence could be considered accessible under a process in which one mutation fixes at a time. Since each accessible pathway consists of a sequence of fixation events, the weight of a pathway can be further quantified by the product of the relative fixation probabilities (15) along the path. Under this measure only a handful of dominating pathways were identified for the β -lactamase system, leading the authors to the conclusion that Darwinian evolution is much more constrained, and hence predictable, than previously recognized.

The periodic selection dynamics terminates at local fitness maxima, which are sinks in the fitness graph where all links carry incoming arrows (Fig. 6). As no beneficial mutations are available at such a genotype, the population cannot escape. In practice this implies that higher order processes such as double mutations or stochastic tunneling events become important, which occur on time scales beyond those considered in the periodic selection scenario [47, 48]. Lo-



Fig. 7: Fitness landscape composed of 8 individually deleterious marker mutations in the filamentous fungus Aspergillus niger [43, 45]. Each mutation resides on a different chromosome, and combinations of mutations were generated by exploiting the parasexual cycle of the fungus. Fitness was measured in terms of mycelial growth rate and normalized to the maximal (wild type) growth rate. Out of the $2^8 = 256$ possible combinations, 70 were not detected in the experiment. The corresponding genotypes were therefore classified as lethal and assigned zero fitness. The fitness values of the $\binom{8}{k}$ genotypes that carry the same number k of mutations are plotted above the same point of the horizontal axis, and the fitness values of genotypes that differ by a single mutation are connected by lines. Local fitness maxima are indicated in red. Courtesy of Ivan G. Szendro.

cal fitness maxima limit the accessibility of high fitness genotypes and present obstacles to the evolutionary process, a concern that was articulated already by Sewall Wright in the pioneering paper which first introduced the fitness landscape concept [7]. The number of local fitness maxima and the number of evolutionarily accessible pathways leading up to them constitute important measures of fitness landscape complexity, which have been applied to a wide range of empirical data sets over the past few years [40, 42].

Acknowledgements: My understanding of evolutionary phenomena has been shaped by many enjoyable collaborations which were funded by DFG within SFB 680 *Molecular bases of evolutionary innovations* and SPP 1590 *Probabilistic structures in evolution*. Special thanks are due to Arjan de Visser and his group at Wageningen University.

References

- [1] D.C. Dennett, *Breaking the spell. Religion as a natural phenomenon* (Viking, New York, 2006)
- [2] T. Dobzhansky, American Zoologist 4, 443 (1964)
- [3] C. Cercignani, *Ludwig Boltzmann. The man who trusted atoms* (Oxford University Press, 1998)
- [4] R.A. Fisher, The genetical theory of natural selection (Clarendon Press, Oxford 1930)
- [5] J.B.S. Haldane, Proc. Camb. Philos. Soc. 23, 838 (1927)
- [6] S. Wright, Genetics 16, 97 (1931)
- [7] S. Wright, Proc. 6th Int. Congr. Genet. 1, 356 (1932)
- [8] P.A.P. Moran, Proc. Camb. Philos. Soc. 54, 60 (1958)
- [9] M.A. Nowak, Evolutionary Dynamics (Harvard University Press, 2006)
- [10] R. Durrett, Probability Models for DNA Sequence Evolution (Springer, New York 2008)
- [11] S.-C. Park, D. Simon and J. Krug, J. Stat. Phys. 138, 381 (2010)
- [12] R.A. Blythe and A.J. McKane, J. Stat. Mech.: Theory Exp. P07018 (2007)
- [13] J. Wakeley, *Coalescent Theory: An Introduction* (Roberts & Company, Greenwood Village 2009)
- [14] N.G. van Kampen, Stochastic Processes in Physics and Chemistry (Elsevier, Amsterdam 2007)
- [15] J. Hofbauer and K. Sigmund, Evolutionary games and replicator dynamics (Cambridge University Press 1998)
- [16] M. Kimura, The neutral theory of molecular evolution (Cambridge University Press 1983)
- [17] M. Kimura, Genetics 47, 713 (1962)
- [18] J.H. Gillespie, Theor. Popul. Biol. 23, 202 (1983)
- [19] H.A. Orr, Evolution 59, 216 (2005)
- [20] L. de Haas and A.F. Ferreira, *Extreme value theory: An introduction* (Springer, New York 2006)
- [21] M.F. Schenk, I.G. Szendro, J. Krug and J.A.G.M. de Visser, PLoS Genet. 8, e1002783 (2012)
- [22] S. Kumar, Nat. Rev. Genet. 6, 654 (2005)
- [23] P.J. Gerrish and R.E. Lenski, Genetica 102-103, 127 (1998)
- [24] S.-C. Park and J. Krug, Proc. Natl. Acad. Sci. USA 104, 18135 (2007)
- [25] M.J. Wiser, N. Ribeck, R.E. Lenski, Science 342, 1364 (2013)
- [26] B.H. Good, M.J. McDonald, J.E. Barrick, R.E. Lenski and M.M. Desai, Nature 551, 45 (2017)
- [27] I.M. Rouzine, É. Brunet and C.O. Wilke, Theor. Pop. Biol. 73, 24 (2008)
- [28] S.-C. Park and J. Krug, Genetics 195, 941 (2013)
- [29] J.A.G.M. de Visser, C.W. Zeyl, P.J. Gerrish, J.L. Blanchard and R.E. Lenski, Science 283, 404 (1999)
- [30] J.A.G.M. de Visser and D.E. Rozen, Genetics 172, 2093 (2006)
- [31] A. Eyre-Walker and P.D. Keightley, Nat. Rev. Genet. 8, 610 (2007)
- [32] J. Haigh, Theor. Pop. Biol. 14, 251 (1978)
- [33] J.J. Metzger and S. Eule, PLoS Comp. Biol. 9, e1003303 (2013)
- [34] H.J. Muller, Mutation Research 1, 2 (1964)
- [35] H.J. Muller, Am. Nat. 66, 118 (1932)
- [36] E.A. Martens and O. Hallatschek, Genetics 189, 1045 (2011)

- [37] J. Otwinowski and J. Krug, Phys. Biol. 11, 056003 (2014)
- [38] P.C. Phillips, Nat. Rev. Genet. 9, 855 (2008)
- [39] F.J. Poelwijk, V. Krishna and R. Ranganathan, PLoS Comp. Biol. 12, e1004771 (2016)
- [40] J.A.G.M. de Visser and J. Krug, Nat. Rev. Genet. 15, 480 (2014)
- [41] D.M. Weinreich, R.A. Watson and L. Chao, Evolution 59, 1165 (2005)
- [42] I.G. Szendro, M.F. Schenk, J. Franke, J. Krug and J.A.G.M. de Visser, J. Stat. Mech.: Exp. Theory P01005 (2013)
- [43] J.A.G.M. de Visser, S.-C. Park and J. Krug, Am. Nat. 174, S15 (2008)
- [44] K. Crona, D. Greene and M. Barlow, J. Theor. Biol. 317, 1 (2013)
- [45] J. Franke, A. Klözer, J.A.G.M. de Visser and J. Krug, PLoS Comp. Biol. 7, e1002134 (2011)
- [46] D.M. Weinreich, N.F. Delaney, M.A. DePristo and D.M. Hartl, Science 312, 111 (2006)
- [47] D.M. Weinreich and L. Chao, Evolution 59, 1175 (2005)
- [48] I.G. Szendro, J. Franke, J.A.G.M. de Visser and J. Krug, Proc. Natl. Acad. Sci. USA 110, 571 (2013)
F 3 Biophysics of Killing

Heiko Rieger Center for Biophysics and Department of Physics Saarland University

Contents

1	Introduction	2
2	Friend or foe: how to recognize an unknown enemy	2
3	Activation and formation of the immunological synapse	3
4	Target search and environmental influence	4
5	MTOC relocation during T cell polarization	7
6	Calcium dynamics during T cell polarization	8
7	Vesicle delivery and intracellular search strategies	11
8	Summary	13

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Immunology is a huge field that covers the study of the immune systems in all organisms [1]. Every one of us has an immune system and it saves our live by fighting hostile intruders, pathogenes, causing infections in our body. Cells of the innate immune system respond recognize and respond to pathogens in a generic way and provide an immediate defense against infection but do not provide long-lasting immunity to the host. The adaptive immune system can recognize novel pathogens (say viruses) which have never before infested our body and develop strategies for their destruction. Moreover it can develop memories of earlier infections and can thus rapidly recognize and fight pathogens intruding our body for the second time.

We will not endeavor to give an overview over the wide area of immunology in a single lecture, we refer to "Janeway's Immunobiology" [1] or the corresponding chapters of Alberts [2]. Instead we will focus here on the crucial aspect in an immune response: the actual elimination or killing of pathogen-infected or tumorigenic cells. The main cytotoxic killer cells of the human body to achieve this are the natural killer (NK) cells from the innate immune system and cytotoxic T-lymphocytes (T cells) from the adaptive immune system. Common to both is that they form a tight contact, the immunological synapse (IS), with targets and release their lytic granules containing perforin/granzyme and cytokine containing vesicles via exocytosis reminiscent of synaptic contacts and neurotransmitter release from one neuron to another. Although T cells can also kill via non-directional release of cytotoxic material the "kill-shot" via IS formation is most frequent and effective.

Successful killing involves the following stages:

- activation of the T cell (not necessary for NK cells)
- target search in an inhomogeneous environment
- recognition of the target upon contact
- polarization and establisment of the IS
- delivery of the cytotoxic material

The sketch on the right (from [17]b) illustrates the events from target recognition to killing, in particular the establishment of an immunological synapse (IS).

In each stage physical processes can be identified that we aim to elucidate in this lecture. In the following we treat these stages in separate sections, each providing a brief overview over its biology and biophysics.

2 Friend or foe: how to recognize an unknown enemy

A central problem to be solved by the immune system is to discriminate between the organism's own healthy cells and tissue – and so to avoid harming its own organism by "friendly fire" – and virus-infected, tumor or intruders. For completeness we give a brief and simplified summary of the fascinating story of friend or foe discrimination in the immune system – an introduction with a physics perspective is [3], the full details of the current biomolecular understanding can be found in [1]. A key player on the target side is the major histocompatibility complex (MHC), which is a set of cell surface proteins, which is essential for the adaptive immune system to



recognize foreign molecules, which in turn determines histocompatibility. Natural killer (NK) cells as part of the innate immune system recognize the absence of MHC complex on the surface of a cell, for example bacteria, which they will then eliminate. T cells on the other hand check the MHC complex for signs of foreign cell material as follows:

The MHC complex binds to antigens (here peptides, i.e. amino acid chains) derived from the cell owning the MHC and display them on the cell surface for recognition by the appropriate T cells. Each MHC molecule on the cell surface displays a molecular fraction of one of the proteins continually synthesized and degraded within the cell, called an epitope. In its entirety the MHC population on the cell surface indicates the balance of proteins within the cell. Virus-infected or tumor cells display antigens that do not belong to the epitope repertoire of the healthy organism and can thus be recognized by T cells expressing the corresponding T cell receptor (TCR) able to bind to this antigene. The population of naive T cell express a broad palette of antigen specific T cell receptors which can recognize a substantial number of distinct antigene-MHC-complexes. Due to this diversity only a few from the $O(10^{12})$ T cells of the human body are expected to bind to a particular non-self antigen - estimates are that about one T cell in 10^{5} - 10^{6} T cells is specific for a given antigen [4].

The immune repertoire, i.e. the diversity of immune cell receptors that exist in T cells (note that each T cell expresses only a single receptor kind), is generated during the formation of the T cell in the thymus: it is the outcome of a remarkable process in which germline DNA is edited to produce a repertoire of T cells with varied antigen receptor genes [1]. The process is called VDJ recombination because the germline contains multiple versions of so-called V- (for "variable"), D- (for "diversity"), and J- (for "joining") genes, particular instances of which are quasi-randomly selected, stochastically edited, and joined together to produce a new surface receptor gene each time a new immune system cell is generated. Following that, the repertoire is edited in the thymus. Cells whose protein would cause an immune reaction with the body, are removed. Cells which actually detect an invading organism, become more numerous. And cells with new types, may be added. The statistical distribution of these biochemical events (and the resulting receptor coding sequences) has recently be inferred from the large T cell sequence repertoires that are becoming available via high-throughput sequencing technology [5] providing insight into the molecular mechanisms involved [6].

3 Activation and formation of the immunological synapse

In contrast to NK cells a naive population of T-lymphocytes needs to be stimulated by specific antigen presenting cells (APCs), such as dendritic cells (DC) or macrophages. DCs posses different mechanisms to internalize extracellular material (phagocytosis, makropinocytosis, virus infection, transfer from other DCs) hence they can present antigens from practically all pathogens. When they do this, for instance in the micro-environment of an inflammatory site, they wander into the lymphatics. In the lymph nodes they present their surface to naive T cells and act as so-called antigen presenting cells (APCs). Once a naive T cell with the appropriate T cell receptor (TCR) encounters a DC with the corresponding MHC-antigene complex a tight contact between DC and T cell, a immunological synapse (IS), is established. At the IS signaling events take place (also involving calcium as a second messanger as discussed in section 6 below) that will activate the T cell (see [7] for a review), which means that it will proliferate and search for target cells presenting the antigen that fits its TCR.

The specific binding of the MHC-antigen complex to the TCR leads to the recruiting of several



Fig. 1: Left: Model of activation of T cells by micro-clusters formed during the initial phase of T cellAPC encounters (i.e. before the formation of SMAC). The 70kDa kinase ZAP-70 is activated by binding to the phosphorylated tyrosine groups of the CD3 ζ -chains. **Center**: abolishment of ZAP-70 activation by CD45-mediated continuous dephosphorylation of CD3. Right side: formation of active immune synapse (IS) by lateral segregation of bound TCRMHCpeptide pairs, resulting in the expulsion of the inhibitor CD45. **Right**: Model of the formation of the peripheral adhesion zone (such as p-SMAC) serving in the long-time stabilization of cellcell contacts enclosing a central zone (such as a c-SMAC) in which signalling micro-domains of tight adhesion formed by MHCAG binding (left side) could coexist with areas where mature TCRCD3 complexes are taken up by endocytosis (right side). The inset on the left shows the coupling of the MT plus end to the actin cortex by MTactin binding proteins. From [8].

other membrane-bound helper molecules (such as CD4 or CD8, where CD stands for "cluster of differentiation"). This results in the generation of micro-domains of tight adhesion that are composed of a mixture of the receptor-ligand pairs (formed between MHC-antigen and the TCR) and the helper molecules (see [8] for more details). The adhesion of the T cell on the APC leads to the assembly of proteins involved in T cell activation into micro-domains (also called adhesion domains) from which proteins counteracting the activation (such as CD45) are expelled by steric forces, see Fig. 1. This activation of the T cell by numerous adhesion domains formed within the contact zone is reminiscent of the activation of nerve or muscle cells by synapses contacting the target cell, which motivated their name immunological synapses. In a second phase of the T cell - APC contact establishment, long-lived global reaction spaces (called supramolecular activation complexes (SMAC)) form by talin-mediated binding of the T cell integrin (LFA-1) to the counter-receptor ICAM-1, resulting in the formation of ring-like tight adhesion zones (peripheral SMAC), c.f. Fig. 1 right). The adhesion domains move to the center of the intercellular adhesion zone forming the central SMAC, which serve in the recycling of the adhesion domain by down-regulation of the signaling by internalization (endocytosis) of the TCR-CD3 complexes. In [8] a model is proposed for the formation of p- and c-SMAC and the Mexican hat-like adhesion caps formed at the IS. Similar global reaction platforms are formed by killer cells to destruct target cells (see section 5 and 6 below), which implies that in this way closed reaction spaces for the localized secretion of lytic proteins, such as perforin, are formed, avoiding the loss of toxic material into the surrounding extracellular space.

4 Target search and environmental influence

During an immune response a T cell has to encounter its cognate antigen many times in different contexts and tissues which necessitates their high motility. The movement of T cells in the lymph nodes has been well studied with two photon microscopy, which showed that they indeed migrate very fast with average speeds of 10 $\mu m/min$ and peak speeds of up to $25 \ \mu m/min$ [9]. They are highly versatile migrators that can effectively navigate almost any tissue of the body. As other leukocytes, T cells migrate using an amoeboid mode of locomotion in three-dimensional in vivo environments. Characteristic for this T cell migration mode is a rounded but polarized morphology with frequent and rapid extensions of pseudopodia at the leading edge. Those are driven by the polymerization of filamentous actin (F-actin) that extends the plasma membrane in the direction of migration. The lagging tail contains stable F-actin networks rich in myosin motors and undergoes little overall deformation. In order to adapt to the varying tissue architecture and molecular composition of their surrounding T cells can rapidly alternate between adhesion-dependent and adhesion-independent motility. During adhesion-driven migration the cell is anchored to the cell to the substrate by integrin-mediated attachments in the extended pseudopodia and the subsequent forward displacement of the cell is driven by the retrograde flow of treadmilling actin filaments coupled to the extracellular matrix (for the underlying biochemical and mechanical processes see chapter C6 on the active cytoskeleton). The contractile lagging tail is only loosely adherent and actomyosin-mediated forces help detach mature adhesion sites and propel the cell body forward. In non-adhesive environments, leading edge pseudopodia are wedged into confined spaces in the extracellular matrix, which, coupled with actomyosin contraction of the rear, allows for a concerted translocation of the cell.

The principal purpose of T cell motility is to search for cognate antigen on APCs and for target cells. The process of search is a fundamental requirement in almost any biological system in which many agents (that is, cells and organisms) reside within an ecosystem much larger than their perceptual capabilities [10]. As with many biological movements, T cell motility was initially characterized as resembling a random walk or diffusion [11]. Two different types of random migration have been used to model T cell motility: diffusive (Brownian-type) random walks [11] and superdiffusive (Lévy-type) random walks [12]. When pauses drive the overall movement, subdiffusive patterns may also be observed. In addition, under some circumstances and for short times, T cells can undergo fully ballistic migration (that is, in a nearly straight line). This range of T cell motility is generated by a combination of cell-intrinsic locomotion events (for example, those controlled by rates of actin polymerization and location of cortex contraction), physical guidance cues from the microenvironment (for example, cues provided by collagen fibres and the orientation of stromal cells) and chemical information provided by the microenvironment (for example, chemokines, antigen dose and through co-stimulatory molecules) [10].

In a pathological scenario, not all cells in a given patrolling area are necessarily target cells. For example, NK cells encounter stromal cells, infiltrated immune cells, as well as malignant cells with expression of MHC class I molecules. These bystander cells pose a challenge to NK cells to efficiently identify their targets in a complex microenvironment. Whether and how the presence of bystander cells can affect the efficiency for NK cells to find and kill their targets has been investigated in [13]. There it has been shown that the presence of non-target bystander cells unexpectedly enhance the killing efficiency as well as NK cell migration. In Fig. 2a-c the analysis of a simple reaction-diffusion model for target search in an obstacle park shows that obstacles generally decrease the search efficiency. Thus the observation of an increase in search efficiency in the presence of bystander cells implies that they cannot be just obstacles but must somehow accelerate the search of the NK cells. In Fig. 2d a corresponding simplified reaction diffusion model is sketched in which bystanders are surrounded by an circular area in which the diffusion constant of the searchers is increased. 2e-g shows the analysis of this model and

supports quantitativley the prediction that the assumed local acceleration of the NKs close to bystander indeed leads to a search efficiency increase.

When the NK migration speed was experimentally measured and a comparison between speeds close to and far away from bystanders were made it was observed that NK motility close bystanders was significantly increased. It turned out that bystander cells increase the H_2O_2 concentration in the surrounding environment, which is known to have an accelerating effect on the migration of killer cells [13]. Thus bystander cells could be used to manipulate the microenvironment of killer cells searching for targets increasing their search efficiency.



Fig. 2: (a) Sketch of the mathematical model of a Brownian (diffusive) disk-like particle performing a random search for disk-like targets among disk-like obstacles: the black disk represents one randomly moving NK cell, the black wiggly lines its diffusive motion with diffusion constant D, grey disks the immobile obstacles of radius robs, and green disks immobile targets of radius R. Targets that are hit, here at time t_{kill} , are removed upon contact with the NK cells. (b) The ratio of killed targets is depicted as a function of time for $N_k = 20$ killers, $N_t(0) = 20$ targets, different numbers of obstacles N_o , and $r_{obs} = 0.5$. (c) The average half time $t_{1/2}$ is given as a function of the number of obstacles N_0 for different values of robs for the same number of killers and targets (left) and for the case of only one killer (right). (d) Same as (a) but now with bystanders (blue). Diffusion of the N_K cells is accelerated by an increased diffusion constant, $D_{\rm acc} > D$, within the circles with radius Δ around bystanders, as indicated by red portions of the wiggly lines representing the random killer motion. Otherwise the symbols are similar to (a). (e) The ratio of killed targets for the model shown in (d) as a function of time for $N_k = 20$ killers, $N_t(0) = 20$ targets, $N_o = 200$ obstacles, $D_{acc} = 4$ and $\Delta = 3$ and different numbers of obstacles. (f) Average half time $t_{1/2}$ as a function of N_b for $N_k = 20$ killers, $N_t(0) = 20$ targets, $N_o = 200$ obstacles and $D_{\rm acc} = 4$ and different values of Δ . (g) Average half time $t_{1/2}$ as a function of $D_{\rm acc}$ for $N_k = 20$ killers, $N_t(0) = 20$ targets, $N_o = 200$ obstacles and $\Delta = 3$ and different numbers of obstacles. From [13]

5 MTOC relocation during T cell polarization

Once a NK or T cell has identified a target cell (see section 2) and established a contact zone (see section 3) one observes a re-polarization of the cell involving the rotation of the microtubule (MT) spindle and a movement of the centrosome or microtubule organizing center (MTOC), the organelle in which the MT fibers of the cytoskeleton are anchored, to a position that is just underneath the plasma membrane at the center of the IS [14, 15, 19] (c.f. Fig. 3). Concomitantly a massive relocation of organelles attached to MTs is observed, including the Golgi apparatus [16] lytic granules (the vesicles containing the cytotoxic material) [17] and mitochondria [18]. As a consequence granules come closer to the point where their cytotoxic content can be release towards the target (see section 7) and mitochondria modulate the redistribution of calcium in the cell leading to increased cytosolic calcium concentrations necessary for signaling during activation and killing (see section 6).



Fig. 3: (A) *Time-lapse images of a JurkatRaji conjugate (transmitted image) in which the Jurkat has been labeled with 3GFP-EMTB to mark MTs and RFP-Pericentrin to mark the MTOC (which appears yellow) (fluorescent images). Zoomed images of the boxed regions in the fluorescent images are shown underneath. The white arrows point to the straightening MTs, and the yellow arrowheads point to the center of the IS where these MTs abut.* (B) As in A, but including a kymograph demonstrating the shortening of the straight MT stalk. (C) Cartoon depicting the two apparent kinetic phases.From [19]

Experiments suggest the massive relocation of the MTOC is driven by dynein, the molecular motors that move along MTs towards the minus end (towards the MTOC). According to this hypothesis, repositioning is driven by a cortical sliding mechanism involving dynein motors anchored at the IS that reel in the centrosome by pulling on MTs that pass over the interface [15]. A theoretical model analyzing the cortical sliding mechanism has been developed and analyzed in [20], see Fig. 4. The cell outline consists of an unattached round part and of a flat part which is attached to the target cell (called synapse, or synaptic plane). The large nucleus is coupled to the aster of MTs converging near its surface, and the mobility of both is constrained by the cell outline. MTs slide along the cell outline in the areas of contact with the targets. This active sliding drives all movements that are observed. The movements are opposed by MT bending elasticity and by viscous drag in the cytoplasm [20].

qualitatively the motion of the MT-network together with large fluctuations of the centrosome next to the cell-cell interface and fluctuations between interfaces with simultaneously engaged targets as observed experimentally in [15].

The experiments of [19], displayed in Fig. 3 examined in closer detail the dynamic organization of the MT cytoskeleton during MTOC re-positioning. Immediately after conjugation of the T cell with the APC one observes several MTs that project from the centrosome toward the APC (Fig. 5A, inset, white arrows, 0s frame) and that terminate in end-on fashion at the approximate center of the IS (Fig. 5A, compare yellow arrowheads in the phase and fluorescent images). Over the next 30s, these MTs straighten and come together to form an apparent bundle of MTs (referred to as the MT "stalk") that terminates in end-on fashion near the center of the IS (follow the white arrows in the insets). Most importantly, MT stalks were then observed to undergo shortening (Fig. 5B; follow the yellow arrowhead marking the point of contact of the stalk with the IS relative to the position of the yellow centrosome), causing the IS and centrosome to approach each other (in this case the IS interface came toward the centrosome). It could be shown that shortening MT stalks pulled the centrosome all the way to the IS and that such perpendicular, centrally anchored, shortening MT stalks coincident with MTOC repositioning [19]. These results are consistent with MTOC repositioning being driven by a capture-shrinkage mechanism focused at the center of the IS and not with the cortical sliding mechanism operating at the IS periphery, c.f. Fig. 4.



Fig. 4: Left: Sketch of the cortical sliding model. MTs (modelled as semi-flexible filaments, are blue, dynein motors anchored at the IS are red. From [20]. **Right:** Corresponding sketch of the capture end-on shrinkage mechanism. Red indicates the point in the IS where the end of a MT bundle is captured and pulled by dynein.

It is a force dependent depolymerization rate of MTs that realizes the chemo-mechanical shrinkage of the MTs when their plus end is pulled against the cell membrane by dynein, as has been demonstrated recently in vitro [21]. It was also observed that MTOC repositioning involves two distinct kinetic phases comprising a fast $(3.3\mu m/min)$ polarization phase and a slower $(0.9\mu m/min)$ docking phase, with a transition between them occurring at $2.2\mu m$ from the IS. The biomechanical origin of the observed biphasic repositioning is still elusive.

6 Calcium dynamics during T cell polarization

Calcium (Ca²⁺) is one of the most important intracellular and intercellular messengers. It transmits signals that arrive at the cell surface to intracellular targets via the transient increase of the intracellular concentration. Information is often transmitted as a Ca²⁺wave travelling through

the cytoplasm of a cell or a group of cells. Ca^{2+} is a simple ion which cannot transmit information by its binding specificity or simply by its presence. Consequently, the signal is encoded in temporal and spatial patterns very similar to the use of electric voltage or current signals in information technology [22]. A comprehensive review on the physics of Ca^{2+} signaling can be found in [22].

Several of the signaling steps governing target killing through T cell and/or NK cells are Ca^{2+} dependent. 1) MTOC relocalization to the IS [19]; 2) mitochondria relocalization to the IS [18]; 3) secretion of LG at the IS [23]; 4) perforin-dependent lysis of target cells [24]. Ca^{2+} signaling in NK and T cells is initiated by binding of the TCR to to its cognate antigen (see section 2 and 3). The following signaling cascade leading to Ca^{2+} release from the endoplasmatic reticulum (ER) is quite common in eukaryotic cells (see [22]): The TCR stimulates PLC γ 1 activity that generates IP₃, which diffuses to the ER where it binds to Ca^{2+} permeable ion channel receptors (IP₃R) on the ER membrane leading to the release of ER Ca^{2+} into the cytoplasm (see Fig. 5A for a sketch). What comes next is characteristic for NK and T cell Ca^{2+} dynamics: Depletion of ER Ca^{2+} triggers a sustained influx of extracellular Ca^{2+} through the activation of plasma membrane Ca^{2+} release-activated Ca^{2+} (CRAC) channels in a process known as store-operated Ca^{2+} entry (SOCE). This SOCE has to be sustained during T cell activation.



Fig. 5: (A) Schematic representation of ion channels, the ER with the STIM proteins and the mitochondria. (B) Schematic representation of CRAC / Orai channel formation, After ER store depletion STIM1 proteins multimerize and move and attach to the PM junction where they trap Orai1 dimers to form CRAC channels and to mediate the Ca^{2+} influx. (C) Schematic representation of the reaction scheme underlying CRAC / Orai channel formation . Reaction 1 describes the STIM multimerization, reaction 2 the attachment of these multimers to the PM. Reactions 3 to 5 represent the stepwise CRAC channel formation. Reaction 6 is a additional reaction between CRAC channel subunits and fully open CRAC channels.

For decades the CRAC channels had only been identified by their biophysical properties, and it has just been in the past few years that the pore-forming subunit of the channels was identified as the four-transmembrane domain molecule Orai1. Additionally, studies in the past few years have revealed the sensor for depleted ER Ca^{2+} stores and the activator of CRAC channels as stromal interaction molecule (STIM). STIM is an ER-resident transmembrane protein which has a Ca^{2+} binding site within the ER lumen, serving as Ca^{2+} sensor, and a domain that can bind to Orai when the Ca^{2+} binding site in the ER lumen is empty. Thus the drop in ER luminal Ca^{2+} is sensed by STIM, which undergo a conformational change, multimerize and move to regions near the plasma membrane, so called plasma membrane junctions (PMJ). Here STIM1 proteins trap Orai1 ion channel proteins diffusing within the plasma membrane and depending on the STIM1-Orai1 stochiometry different Orai1 conductance states are reached open and selectively conduct Ca^{2+} ions into the cell (see Fig. 5B). A quantitative analysis of these membrane bound events leading to SOCE has been perfromed in [27] based on a reaction-diffusion model with reaction schemes as indicated in Fig. 5C.

Ca²⁺dynamics during activation of T cells is particularly interesting, because it involves a sustained increase of intracellular Ca²⁺levels that results in the activation of Ca²⁺ and calmodulindependent transcription factors hat in turn activate a variety of transcription programs. Activated calcineurin dephosphorylates members of the nuclear factor of activated T cells (NFAT) family, leading to their translocation to the nucleus resulting in differential gene expression patterns, which eventually transform the naive T cell into a killer. The duration of the Ca²⁺increase necessary for a succesfull completion of these steps (up to 1 hour) is remarkable, because Ca²⁺pumps and Ca²⁺buffers in living cells keep the cytosolic Ca²⁺level very low (ca. 80 nanomolar, as compared to, for instance, Ca²⁺concentration in mineral water which is a few millimolar, i.e. 10^6 times higher). The interesting question is, how T cells succeed to outmaneuver its Ca²⁺removing machinery to sustain 10 fold increased Ca²⁺concentrations for up to 1 hour.



Fig. 6: Top: Sketch of the reaction-diffusion model for the intracellular Ca^{2+} dynamics in T cells represented by the Ca^{2+} fluxes between the different compartments: cytosol, mitochondria, ER, and extracellular space. Bottom: Ca^{2+} concentrations depend on mitochondria ratation angle - model prediction with working SERCAs and spindle ER geometry. **Bottom:** (A) Sketch of the mitochondria (green) and a spindle ER geometry (red). (B) Time course of Ca^{2+} concentrations in the cytosol (top), at the IS (middle), and in the ER (bottom) for the geometry in A. CRAC channels are activated at t = 10s and deactivated at t = 60s. (C) Plateau value (at t = 40s) of cytosol (global) and IS (local) Ca^{2+} concentrations as a function of the rotation angle of the mitochondrial network for different ratio of PMCA accumulation at the IS. From [28].

It turned out that a key factor for the sustained increased cytosolic Ca²⁺concentration is the relocation of mitochondria towards the IS [26], which are powerful Ca^{2+} buffers and can distribute Ca^{2+} within the cytool via a filamentous network. The reason for the mitochondria relocation was identified to be the cytoskeleton rotation discussed in the previous section 5 [28]: mitochondria attached to MTs (via molecular motors, by which they can also move along MTs) are simply dragged with the moving cytoskeleton. Activated Orai channels also accumulate at the IS such that it appears plausible that mitochondira modulate the strength of the Ca^{2+} influx. However, since activated Orai channels are directly linked (via STIM) to the ER, it is questionable that mitochondria can come closer than 100-200 nanometer to the Orai channel, which impedes a direct modulation of channel inhibition [26]. Since plasma membrane calcium ATPase (PMCA) pumps (which expel Ca^{2+} from the cytosol against a concentration gradient into the extracellular space) also accumulate at the IS [26] an alternative explanation has been proposed in [27, 28]: the mitochondria absorb the Ca²⁺entering the IS microdomain through Orai/CRAC channels before it can be extruded again by the PMCA clusters in the IS. With the help of a mathematical model for the intracellular Ca^{2+} dynamics it was indeed shown in [28] that the repositioning of mitochondria alone can modulate the global cytosolic Ca^{2+} concentration, independent of any influence of mitochondrial position on CRAC/Orai channel activity.

A sketch of the model is shown in Fig. 6: it consists of a reaction-diffusion model for the intracellular Ca^{2+} dynamics summarized in the top panel of Fig. 6 and a specific geomtry for various cell compartmens (cell body, mitochondria, ER) and location of the CRAC channels shown in Fig. 6A. The assumed geometry of mitochondria is inspired by the fact that they form filaments along MTs, which themselves are arranged in a spindle-like geometry (c.f. 3). Results for the local (at the IS) and global (cytosolic) Ca^{2+} concentration are shown in Fig. 6B and C and demonstrate that global Ca^{2+} concentration increases with decreasing rotation angle of the cytoskeleton network, i.e. with decreasing distance of mitochondria from the IS.

7 Vesicle delivery and intracellular search strategies

In order to kill the target cell by NKs and T cells the lytic granules have to be transported to the IS where they can dock and release their cytotoxic content via exocytosis. The cytoskeleton of living cells is a self-organizing filamentous network that shapes the mechanical and rheological characteristics of the cell and coordinates cargo transport between different cellular regions (see also chapters C3, C4, C5 and C6). Intracellular transport of particles equipped with one or several motors [30] switches between two modes of motion: free diffusion within the cytosol and ballistic motion during intermittent bindings of the motor(s) to a cytoskeleton filament. Typically cargo, like proteins, vesicles, and other organelles, is produced or emerges in one region of the cell and is needed in some other region or has to fuse with a reaction partner being produced somewhere else. In the absence of a direct connection between origin and destination the transport is a stochastic process [31] with random alternations between ballistic and diffusive motion, which is denoted as intermittent search [32]. A particular set of parameters defining the stochastic process, like the switching rate between ballistic and diffusive transport, represents an intermittent search strategy – and it has been shown that an optimal choice of this parameter enhances the search efficiency in a homogeneous and isotropic environment, i.e. under the assumption of a constant density of filaments with no preferred direction [33].

Real cell cytoskeletons display a complex spatial organization, which is neither homogeneous nor isotropic. For instance in cells with a centrosome, like NK and T cells, the microtubules

emanate radially from the microtubule organizing center (MTOC) and actin filaments form a thin cortex underneath the plasma membrane with broad distribution of directions again centered around the radial direction (c.f. Fig. 7a). Since the MTOC is frequently located close to the nucleus, transport of cargo between the plasma membrane and the nucleus, as for instance necessary for the establishment of synaptic junctions, appears to be facilitated by this cytoskeleton organization. In essence the specific spatial organization of the cytoskeleton filaments represents, in connection with intracellular transport, an intermittent search strategy, probably highly efficient for specific frequently occurring tasks, but less well suited for others. In [34, 35] the efficiency of spatially inhomogeneous intermittent search strategy for particular task has been analyzed and compared with homogeneous search strategies with the help of a model for intermittent dynamics in spatially inhomogeneous environments. Three basic search tasks frequently encountered in intracellular transport have been investigated in [34, 35]:

1) Transport of cargo from an arbitrary position within the interior of the cell, typically from a location close to the nucleus, to a specific area on the cell boundary, the plasma membrane. The transport of lytic granules to the IS involves the cytoskeleton [17] and represents such a task. The stochastic search for a specific small area on the boundary of a search domain is reminiscent of the so-called narrow escape problem [36].

2) The enhancement of the reaction kinetics between two reaction partners by motor assisted ballistic transport. It has already been demonstrated that spatially homogeneous and isotropic intermittent search strategies can be efficient, but can only realized in those parts of a biological cell, where cytoskeleton filaments are homogeneously and isotropically distributed, which is certainly not true for the whole cell body.

3) Finally the combination of the reaction and escape problem, where cargo has first to bind to a reaction partner before it can be delivered or dock at a specific area in the cell boundary as, for instance, a synapse. A prominent realization is the docking of lytic granules at the IS that requires the pairing with CD3 endosome beforehand [29]. Single lytic granules have a low docking probability at the IS, whereas endosomes containing CD3 molecules have a high docking probability. Apparently it represents an advantageous strategy to guarantee the delivery of cytotoxic cargo exclusively to the IS via binding to these endosomes in advance.

A search strategy that idealizes the cytoskeleton structure in a spherical cell of radius R consists of microtubule filaments emanating radially from the MTOC in the cell center and randomly oriented actin filament in a cortex of width Δ underneath the plasma membrane, as sketched in Fig. 7. Mathematically such a filament distribution is defined by the probability density $\rho(\Omega, r)$ to choose direction Ω conditional on the switch from the diffusive to a ballistic mode at position r and can, for simplicity, be parameterized by a width $\Delta \in [0, R]$ and a probability $p \in [0, 1]$: In the central region, $0 < |r| < R - \Delta$, $\rho(\Omega, r) = p$ for radial outward transport, i.e. Ω is the direction of the position vector, and $\rho(\Omega, r) = 1 - p$ for radial inward transport; and in the periphery, $R - \Delta < |r| < R$, $\rho(\Omega, r) = 1/4\pi$ uniform for all Ω . The intracellular motion of cargo is modeled by an intermittent search process [32], in which a particle per-forms random motion in two alternating modes: Brownian motion with a diffusion coefficient D, and ballistic motion with velocity v. Transitions between the two modes occur stochastically with rates k and k', respectively, see Fig. 7 for a sketch in case of the narrow escape problem, and deterministically from ballistic to diffusive at |r| = 0, $R - \Delta$, and R.

The efficiency of a search strategy, or a specific directional distribution, $\rho(\Omega, r)$, is measured in terms of the time that a searcher needs on average until it hits the target the first time, the so-called mean first passage time (MFPT). A major progress in the analytical determination of the MFPT in 2d and 3d bounded domains was made in [37] with a computational method that is



Fig. 7: (a) Sketch of the cellular cytoskeleton with microtubules in green and actin filaments in red. (b) Sketch of the idealized direction distribution, red only radial transport, blue transport in all directions, green the narrow escape region. (c) Sketch of the random intermittent search process with the idealized direction distribution, Δ =thickness of the idealized actin cortex. Green radial ballistic transport, grey wiggly lines: diffusive motion, red ballistic transport in arbitrary directions, =thickness of the idealized actin cortex. (d) Mean first passage time (MFPT) for the narrow escape problem as a function of the actin cortex thickness in the random intermittent search process sketched left. For small diffusion constants the MFPT has a pronounced minimum between $\Delta/R = 0$ and 0.2 corresponding to a thin actin cortex. From [34].

based on an expression for the MFPT between two nodes of a general graph developed in [38]. Still a spatially inhomogeneous problem like the one defined by above appears analytically intractable, although some regions in the parameter space may be approximately treatable analytically by combining the known results for homogeneous searches. Numerically it was found in [34, 35] that the confinement of randomly oriented cytoskeleton filaments to a thin actin cortex is not a handicap for the cell, but can actually, in conjunction with radial transport along microtubules, substantially increase the efficiency of transport tasks. For the narrow escape problem, which is exemplarily shown in Fig. 7D, it turns out to be a superior strategy to allow only radial outward ballistic transport from the center towards a thin sheet of thickness Δ underneath the boundary, where ballistic transport in all directions is possible. This thin boundary layer allows an accelerated random motion along the boundary to find the escape region, somewhat reminiscent of purely diffusive search with an accelerated surface mediated diffusion [39]. A similar result holds for the reaction kinetics problem, in which the target is not located on the boundary: here again optimal strategies with small thickness Δ exist, in particular better than the homogeneous strategy [36], but the optimal probability for forward/backward radial transport is now around p = 1/2. This result is reminiscent of an acceleration of purely diffusive search kinetics by following boundaries with an increased diffusion constant [40]. The reaction-escape problem combines both scenarios and the optimal forward/backward radial transport probability depends on the size ratio of target and escape region.

8 Summary

We have discussed various biophysical processes involved in the successful killing of pathogeninfected or tumorigenic cells by NK or T cells: the self-organized arrangement of membrane proteins into characteristic patterns at the IS, the mechanical relocation of the MT cytoskeleton towards the IS by molecular motors, the modulation of the cytosolic calcium concentration by mechanically moving compartments, and the delivery of lytic granules to the IS. Important open questions and issues of ongoing research concern for instance the role of actin in particular and the forces in general in the formation of the IS and the bio-mechanics of the different migration modes of NK and T cells during their search for targets in different environments. Of particular interest here is the question how the cell processes signals from the environment and transforms them into mechano-chemical events connected with the cell cytoskeleton that lead to directional changes during migration and shape and polarization changes upon contact with a target cells. It would be highly desirable to understand the ways in which environmental factor influence the efficiency NK and T cells in searching and killing target cells since it bears the potential to improve the immune response of an organism in fighting a disease or cancer.

Acknowledgements: I thank my collaborators Markus Hoth, Bin Qu and Barbara Niemeyer from the Biophysics Department of the Medical School of the Saarland University as well as my students Anne Hafner, Barbara Schmidt, Karsten Schwarz, Ivan Hornak and Thierry Fredrich for their valuable inputs to this lecture and to my understanding of the biophysics of killing.

Abbreviations:

APC	antigen presenting cell
CRAC	calcium release activated channel
CTL	cytotoxic T lymphocyte
DC	dendritic cell
ER	endoplasmatic reticulum
MFPT	mean first passage time
IS	immunological synapse
LG	lytic granule
MHC	major histocompatibility complex
MT	microtubule
MTOC	microtubule organizing center
NFAT	nuclear factor of activated T cells
NK	natural killer (cell)
PM	plasma membrane
PMJ	plasma membrane junction
PMCA	plasma membrane calcium ATPase
SMAC	supramolecular activation complex
SOCE	store operated calcium entry
STIM	stromal interaction molecule
TCR	T cell receptor

References

- [1] K. Murphy and C. Weaver: Janeway's Immunobiology (Taylor & Francis Ltd.; 9th ed.).
- [2] Bruce Alberts, et al.: *Molecular Biology of the Cell*. (Taylor & Francis Ltd.; 6th ed.).
- [3] A. S. Perelson and G. Weisbuch: *Immunology for physicists*. Rev. Mod. Phys. 69, 1219-1267 (1997).
- [4] J. N. Blattman, et al.: *Estimating the precursor frequency of naive antigen-specific CD8 T cells.* J. Exp. Med. 195, 657664 (2002).
- [5] G. Georgiou, et al.: *The promise and challenge of high-throughput sequencing of the antibody repertoire.* Nature Biotech. 32, 158168 (2014).
- [6] A. Murugana, et al.: *Statistical inference of the generation probability of T cell receptors from sequence repertoires.* Proc. Natl Acad. Sci. USA 109, 16161-16166 (2012).
- [7] J. E. Smith-Garvin, G. A. Koretzky, and M. S. Jordan: *T Cell Activation*. Annu. Rev. Immunol. 27, 591-619 (2009).
- [8] E. Sackmann: Quantal concept of T cell activation: adhesion domains as immunological synapses. New J. Phys. 13, 065013 (2011).
- M. J. Miller, S. H. Wei,1 I. Parker, M. D. Cahalan: *Two-Photon Imaging of Lymphocyte* Motility and Antigen Response in Intact Lymph Node. Science 296, 1869-1873 (2002); M. D. Cahalan and I. Parker: Choreography of cell motility and interaction dynamics imaged by two-photon microscopy in lymphoid organs. Annu. Rev. Immunol. 26, 585626 (2008).
- [10] M. F. Krummel, F. Bartumeus and A. Gerard: T cell migration, search strategies and mechanisms. Nature Rev. Immunol. 16, 193 (2016).
- [11] M. J. Miller, S. H. Wei, M. D. Cahalan, I. Parker: Autonomous T cell trafficking examined in vivo with intravital two-photon microscopy. Proc. Natl Acad. Sci. USA 100, 2604-2609 (2003).
- [12] T. H. Harris, et al.: Generalized Lévy walks and the role of chemokines in migration of effector CD8+ T cells. Nature 486, 545-548 (2012).
- [13] X. Zhou, R. Zhao, K. Schwarz, M. Mangeat, E. C. Schwarz, M. Hamed, I. Bogeski, V. Helms, H. Rieger and B. Qu: Bystander cells enhance NK cytotoxic efficiency by reducing search time. Scientific Reports 7, 44357 (2017)
- [14] B. Geiger, D. Rosen, and G. Berke: Spatial relationships of microtubule- organizing centers and the contact area of cytotoxic T lymphocytes and target cells. J. Cell Biol. 95, 137-143 (1982).
- [15] J. R. Kuhn and M. Poenie: *Dynamic polarization of the microtubule cytoskeleton during CTL-mediated killing*. Immunity 16, 111-121 (2002).
- [16] A. Kupfer, G. Dennert, and S.J. Singer. Polarization of the Golgi apparatus and the microtubule-organizing center within cloned natural killer cells bound to their targets. Proc. Natl. Acad. Sci. USA 80, 7224-7228 (1983).
- [17] a) J. C. Stinchcomb et al.: Centrosome polarization delivers secretory granules to the immunological synapse. Nature 443, 462-465 (2006); b) A. T. Ritter et al.: Actin depletion initiates events leading to granule secretion at the immunological synapse. Immunity 42, 864-876 (2015).
- [18] A. Quintana, C. Schwindling, A. S. Wenning, U. Becherer, J. Rettig, E. C. Schwarz, and M. Hoth: *T cell activation requires mitochondrial translocation to the immunological synapse.* Proc. Natl. Acad. Sci. USA 104, 14418-14423 (2007).
- [19] J. Yi, X. Wu, A. H. Chung, J. K. Chen, T. M. Kapoor, and J. A. Hammer: *Repositioning in T cells is biphasic and driven by microtubule end-on capture-shrinkage*. J. Cell Biol 202,

779-792 (2013).

- [20] M. J. Kim and I. V. Maly: Deterministic Mechanical Model of T-Killer Cell Polarization Reproduces the Wandering of Aim between Simultaneously Engaged Targets, PLOS Comp. Biol. 5, e1000260 (2009).
- [21] L. Laan et al.: Cortical dynein controls microtubule dynamics to generate pulling forces that position microtubule asters. Cell. 148, 502-514 (2012).
- [22] M. Falcke: *Reading the patterns in living cells the physics of Ca*²⁺*signaling.* Advances in Physics 53, 255-440 (2004).
- [23] A.T. Pores-Fernando and A. Zweifach: Calcium influx and signaling in cytotoxic Tlymphocyte lytic granule exocytosis. Immunol. Rev. 231, 160-173 (2009).
- [24] I. Voskoboinik, M. J. Smyth, and J. A. Trapani: *Perforin-mediated target-cell death and immune homeostasis*. Nature Rev. Immunol 6, 940-952 (2006).
- [25] B. Qu, D. Al-Ansary, C. Kummerow, M. Hoth, and E. C. Schwarz: ORAI-mediated calcium influx in T cell proliferation, apoptosis and tolerance. Cell Calcium 50, 261-269 (2011).
- [26] A. Quintana, M. Pasche, C. Junker, D. Al-Ansary, H. Rieger, C. Kummerow, L. Nunez, C. Villalobos, P. Meraner, U. Becherer, J. Rettig, B.A. Niemeyer, M. and Hoth: *Calcium microdomains at the immunological synapse: how ORAI channels, mitochondria and calcium pumps generate local calcium signals for efficient T cell activation*. EMBO J. 30, 3895 (2011).
- [27] M. Peglow, B.A. Niemeyer, M. Hoth, and H. Rieger: *Interplay of channels, pumps and organelle location in calcium microdomain formation.* New J. Phys. 15, 27 (2013).
- [28] I. Maccari, R. Zhao, M. Peglow, K. Schwarz, I. Hornak, M. Pasche, A. Quintana, M. Hoth, B. Qu, and H. Rieger: *Cytoskeleton rotation relocates mitochondria to the immuno-logical synapse and increases calcium signals*. Cell Calcium, 60, 309 (2016).
- [29] V. Pattu, C. Junker C, E. C. Schwarz, S. S. Bhat, C. Kummerow, M. Marshall, U. Matti, F. Neumann, M. Pfreundschuh, U. Becherer, H. Rieger, J. Rettig, and M. Hoth: *Docking* of Lytic Granules at the Immunological Synaps in Human CTL Requires Vti1b-Dependent Paiiring with CD3 Endosomes. J. Immunol. 186, 6894 (2011)
- [30] M. Schliwa and G. Woehlke: Molecular motors. Nature 422, 759 (2003).
- [31] P. C. Bressloff, J. M. Newby: Stochastic models of intracellular transport. Rev. Mod. Phys. 85, 135 (2013).
- [32] O. Bénichou, C. Loverdo, M. Moreau, R. Voituriez: *Intermittent search strategies*. Rev. Mod. Phys. 83, 81 (2011).
- [33] C. Loverdo, O. Bénichou, M. Moreau, R. Voituriez: Enhanced reaction kinetics in biological cells. Nature Physics 4, 134 (2008).
- [34] K. Schwarz, Y. Schröder, B. Qu, M. Hoth, H. Rieger: Optimality of spatially inhomogeneous search strategies. Phys. Rev. Lett. 117, 068101 (2016); K. Schwarz, Y. Schröder, H. Rieger: Numerical analysis of homogeneous and inhomogeneous intermittent search strategies. Phys. Rev. E 94, 042133 (2016).
- [35] A. E. Hafner, H. Rieger: Spatial organization of the cytoskeleton enhances cargo delivery to specific target areas on the plasma membrane of spherical cells. Phys. Biol. 13, 066003 (2016); A. E. Hafner, H. Rieger: Spatial cytoskeleton organization supports targeted intracellular transport. arXiv:1709.05133, Biophys. J., in press (2018).
- [36] Z. Schuss, A. Singer, and D. Holcman: *The narrow escape problem for diffusion in cellular microdo-mains*. Proc. Nat. Acad. Sci. USA, 104, 16098 (2007).
- [37] S. Condamin, O. Bénichou, and M. Moreau: First-Passage Times for Random Walks in

Bounded Domains. Phys. Rev. Lett. 95, 260601 (2005).

- [38] J.-D. Noh and H. Rieger: *Random Walks on Complex Networks*. Phys. Rev. Lett. 92, 118701, (2004).
- [39] O. Bénichou, D. Grebenkov, P. Levitz, C. Loverdo, and R. Voituriez. *Optimal Reaction Time for Surface-Mediated Diffusion*. Phys. Rev. Lett., 105, 150606 (2010).
- [40] T. Calandre, O. Bénichou, and R. Voituriez: *Accelerating search kinetics by following boundaries*. Phys. Rev. Lett. 112, 230601 (2014).

F4 Monogenetic diseases

Christoph Fahlke Cellular Biophysics Institute of Complex Systems Forschungszentrum Jülich

Contents

1	Introduction	2
2 fui	Myotonia congenita demonstrates cellular roles and molecular nctions of muscle chloride channels	2
3 tra	Disturbed anion channel function of secondary-active glutamate unsporters results in ataxia and epilepsy	5
4	Abbreviations	6
Re	ferences	6

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Genetic factors contribute to the disease risk and to the severity of many human diseases. The majority of such conditions do not have a single genetic cause, but are associated with multiple genes combined with environmental and sometimes with lifestyle factors. Examples for such complex or multifactorial disorders are common medical problems such as cardiovascular heart disease, diabetes, or cancer. There are other diseases, which are caused by mutations in indivudal genes and which are inherited according to Mendel's Laws with autosomal dominant, autosomal recessive or X-linked recessive inheritance patterns. Monogenetic diseases usually occur only rarely. However, they link a defined dysfunction of individual proteins to specific changes in organ function, and this particular property makes them a valuable tool to study protein function and cellular roles of affected proteins.

In recent years, improved clincial descriptions together with the introduction of whole exome and whole genome sequencing to human genetics has permitted identification of a large number of gene alterations in monogenetic as well as in multifactorial diseases. Disease-causing mutations affect relevant protein functions and can therefore be used to understand sequence determinants of the affected protein function. We are still far away from a comprehensive understanding of all all proteins involved in cell and organ functions. In the past, monogenetic diseases have helped identifying novel proteins that play important roles in certain cell processes. Lastly, many proteins fulfill more than one function, and changes in protein function by disease-causing mutations provide hints in the particular importance of defined protein functions. This chapter provides examples for using naturally occuring disease-causing mutations for a better understanding of the normal function of proteins and cells.

2 Myotonia congenita demonstrates cellular roles and molecular functions of muscle chloride channels

Myotonia congenita is an inherited human disease characterized by muscle stiffness upon sudden forceful movement. It is transmitted either dominantly (Thomson's myotonia) or recessively (Becker-type myotonia) [1]. Myotonia is a rare disease and often without serious symptoms. However, it represents the first disease, in which changes in ion channel function have been identified as disease cause, and it therefore represents the founding member of the growing family of "channelopathies". Its pathophysiology permitted appreciation of the role of chloride channels in muscle excitability and exemplified for the first time how so-called background conductances can regulate cellular excitability.

The existence of an animal model of myotonia congenita, the "fainting goats", a strain of goats with bizarre attacks of stiffness and rigidity, permitted detailed analysis of the molecular and cellular pathophysiology of this disease, even before the genetic basis of myotonia was understood, Muscle preparations from the myotonic goat permitted the analysis of myotonic muscle fibers with electrophysiological tools. This studies, mainly performed by Shirley Bryant and colleagues, established that muscle stiffness originates from an alteration of the muscle membrane excitability [2, 3]. They demonstrated a drastically reduced chloride conductance (gCl) in muscle fibers from myotonic goats and showed that normal muscle can be made myotonic by blocking muscle chloride channels pharmacologically [4]. These results firmly established chloride channel dysfunction as the cause of myotonia. Later, changes in muscle chloride channels were also shown to be the basis of human myotonias [5].

The muscle chloride conductance is the dominant sarcolemmal conductance at resting potentials [6]. It is small compared with voltage-activated sodium and potassium currents at more depolarized potentials, and muscle chloride channels therefore do not play a role in the repolarization of the action potential. Their main role is to modify the length constant of the t-

tubuluar membrane and the sarcolemma. Any electrical signal is propagated passively along the muscle surface and t-tubular membranes, and during this electronic propagation the amplitude of the electrical signal decays exponentially with increasing distance from the site of its initiation. The length constant describes the steepness of this electrical signal decay, and length constants of membranes are proportional to the square root of their resistance (see chapter D3).



Fig. 1: *A*, In the t-tubules of skeletal muscle, trains of action potentials results in K^+ accumulation. In normal muscle, the subsequent changes in t-tubular membrane potential are offset by the high sarcolemmal chloride conductance. In myotonic muscle fibers, t-tubular depolarizations propagate to the surface membrane and can trigger the spontaneous generation of new action potentials even after the end of the voluntary movement. This "after-depolarization" is marked in red on the right side. B,C, Representative recordings from WT and G230E hClC-1 channels.

Myotonia is due to potassium accumulation and membrane depolarization in the t-tubule that is propagated to the sarcolemma (Fig. 1A). In normal muscle, this depolarization of the t-tubular membrane does not affect the potential of the surface membrane, because the length constant of healthy muscle fibers is small due to a large resting chloride conductance. In myotonic muscle, dysfunctional chloride channels increase the length constant of the t-tubular membrane and the sarcolemma, so that t-tubular depolarizations also depolarize the surface membrane. Muscle chloride channels simply reduce the length constant of muscle and thus allow muscle to electrically tolerate the t-tubule depolarization [3].

The muscle chloride channel, ClC-1, was cloned in 1991 [7], and soon after, linkage of both forms of myotonia congenita to *CLCN1* (the human gene encoding ClC-1) was shown and first disease-causing mutations identified [8, 9]. In the following years, a great number of disease-causing mutations were reported for Thomsen's as well as for Becker's myotonia, spreading over the whole coding region of ClC-1. ClC-1 belongs to a large family of anion channels and transporters, encompassing various proteins with important cellular roles. For many years, the molecular basis of anion selectivity was incompletely understood, and the large size and the unclear membrane topology of CLC channels prevented the use of site-directed mutagenesis and functional testing for finding pore-forming residues. The functional analysis of a myotonia-causing mutations identified the first amino acid residue that contributes to the formation of the anion coinduction pathway of ClC-1 channels.

This mutation was found in a patient with dominant myotonia and resulted in the substitution of a conserved glycine by glutamate, G230E [9]. Functional analysis of mutant channels in heterologous expression systems revealed changes in the ion conduction and selectivity process [10]. WT ClC-1 exhibit inwardly rectifying unitary current conductance with a value of approximately 1.5 pS below -85 mV at symmetrical chloride concentration (Fig. 1B) and tenfold lower values at positive potentials [11]. G230E hClC-1 exhibits a bidrectional rectification with increased conductances at positive and at negative voltages (Fig. 1C). Moreover, whereas WT hClC-1 chloride currents are blocked by larger and poyatomic anion such as Γ or NO₃⁻, these ions show higher permeability than Cl⁻ in mutant channels [10]. Most importantly, G230E permitted cation permeation through hClC-1 [10].



Fig. 2: Structure of the bovine ClC-K channel in a side view (A) and in a view from the extracellular side (B) [13]. The locations of corresponding G230E hClC-1 mutations is shown in yellow and permeating Cl ion as red sphere.

The functional alterations of G230E hClC-1 channels suggested that G230 is close neighborhood to the conduction pore of ClC-1. G230 is part of a sequence motif that is conserved in all CLC channels and transporters. Mutagenesis and site-directed chemical modification demonstrated that residues within this motif modify anion permeation and gating of CLC channels [12]. The first high-resolution structure of a CLC anion channels, a bovine renal ClC-K channel, fully verified this concept [13]. Fig. 2 shows the localization of the homologous residue in bClC-K, together with an anion and the anion permeation pathway, proving direct interaction of this residuee with the anion conduction pathway. Thus, myotonia

congenita was not only instrumental in identifying the cellular roles of ClC-1, but also to clarify the mechanisms of anion selectivity and underlying sequence determinants.

3 Disturbed anion channel function of secondary-active glutamate transporters results in ataxia and epilepsy

Glutamate is the major excitatory neurotransmitter in the mammalian central nervous system. Glutamatergic synaptic transmission is triggered by action potentials eliciting Ca^{2+} influx and subsequent fusion of glutamate-containing synaptic vesicles. Consequently, glutamate diffuses across the synaptic cleft and binds to ionotropic or metabotropic glutamate receptors at postsynaptic membranes. Glutamatergic synaptic transmission is terminated by the removal of glutamate from the synaptic cleft, and this task is performed by a family of glutamate transporters, the "excitatory amino acid transporters (EAATs)". Effective EAAT-mediated glutamate uptake is crucial for high temporal resolution of glutamate receptors, whose activation results in elevated $[Ca^{2+}]$ and neuronal death, they do not survive even slight elevation of resting glutamate concentrations. EAATs are not only secondary-active glutamate transporters, but also anion-selective ion channels. The structural basis of this curious dual function behaviour was resolved recently [14], but the cellular roles of EAAT anion channels remain insufficiently understood.

EAATs are physiologically very important, but only few genetic disease have been linked to this class of transporters. Thus far, four different diseases were linked to *SLC1A3*, the gene encoding the glial glutamate transporter EAAT1: episodic ataxia [15-18], Tourette syndrome [19], Attention Deficit Hyperactivity Disorder (ADHD) [20] and familial hemiplegic migraine [21]. Episodic ataxia (EA) is a rare genetic condition characterized by paroxysmal cerebellar incoordination associated with additional other neurological symptoms. Based on differences in clinical symptoms and in the underlying genes, several different types of episodic ataxia have been defined. Episodic ataxia type 6 clinically differs from other forms in lasting attacks of ataxia and epilepsy and absent myokymia, nystagmus and tinnitu, and mutations in *SLC1A3* were found in patients with EA6. The first naturally occurring *SLC1A3* mutation was found in a 10 year old patient suffering from paroxysmal ataxia, epilepsy and hemiplegia and predicted the substitution of proline by arginine in the fifth transmembrane domain (TM5), P290R [15].

In heterologous expression systems, P290R reduces the number of EAAT1 in the surface membrane and impairs EAAT1-mediated glutamate uptake (Fig. 3A) [22, 23]. Despite the lower number of mutant transporters in the surface membrane, cells expressing P290R EAAT1 exhibit larger anion currents than WT cells in the absence as well as in the presence of external L-glutamate (Fig. 3B,C). Noise analysis revealed unaltered unitary current amplitudes, indicating that P290R modifies opening and closing, and not anion permeation through mutant EAAT1 anion channels. These findings identified gain-of-function of EAAT anion conduction as causal molecular process in the pathogenesis of episodic ataxia.

EAAT1 is mostly expressed in glial cells, and since EAAT1-expressing glia cells exhibit rather small anion conductances and large potassium conductances, changes in EAAT1 anion current amplitude are not expected to modify electrical parameters of these types of astrocytes. Since it appeared possible that gain-of-function of EAAT1 anion channels reduces glial [Cl⁻]_{int}, we measured glial [Cl⁻]_{int} in acute cerebellar slices using the chloride-sensitive dye MQAE and fluorescence lifetime imaging microscopy [24]. We found Bergmann glial [Cl⁻]_{int} to be in a dynamic equilibrium between accumulation by Na⁺–K⁺–2Cl⁻ cotransporter NKCC1 and efflux through EAAT1 and EAAT2 anion channels. Changes in expression level or in the activity of EAAT anion channels modify [Cl⁻]_{int}. Bergmann glial cells express secondary-active GABA transporters (GATs) that couple the inward transport of one GABA to two Na⁺ and one Cl⁻. Thus, reduced [Cl⁻]_{int} will increase the driving force for GABA uptake and reduce synaptic levels of the inhibitory neurotransmitter. Since external glutamate activates EAAT anion channels, these channels permit a crosstalk between excitatory and inhibitory neurotransmission. This hypothetical mechanism still awaits experimental verification. Knock-in animal models carrying gain-of-function mutations of various EAAT anion channels will help understanding how the EAAT anion conduction mode modifies excitability of the human brain.



Fig. 3: *A, WT and mutant EAAT1 uptake currents were measured upon heterologous expression in Xenopus oocytes as difference of current amplitudes with and without glutamate in a chloride-free solution. B, Representative current recordings from HEK293T cells expressing WT or P290R hEAAT1 with KNO3-based internal solution and a NaNO3-based external solution before (top) and after application of 1 mm glutamate (bottom). C, Voltage dependence of mean current amplitudes from cells expressing WT or P290R hEAAT1 in the presence as well as in the absence of glutamate. Modified from [22].*

4 Abbreviations

[Cl ⁻] _{int}	intracellular chloride concentration
ADHD	attention deficit hyperactivity disorder
bClC-K	bovine kidney-specific ClC chloride channel
CLCN1	human gene encoding the muscle chloride channel ClC-1
EAAT	excitatory amino acid transporters
GAT	GABA transporter
gCl	resting sarcolemmal chloride conductance
HEK293T	HEK293 varainat that expresss the SV40 large T antigen
SLC1A3	human gene encoding the glial glutamate transporter EAAT1

References

- [1] R. Rüdel and F. Lehmann-Horn, Pysiol. Rev. 65, 310 (1985)
- [2] S.H. Byrant, Disorders of the Motor Unit (Wiley, New York, 1982)
- [3] R.H. Adrian and S.H. Bryant, J. Physiol. 240, 505 (1974)

- [4] S.H. Bryant and A. Morales-Aguelera, J. Physiol. 219, 361 (1971)
- [5] R.J. Lipicky and S.H. Bryant, Electromyography and Clinical Neurophysiology (Karger, Basel, 1973)
- [6] A.L. Hodgkin and P. Horowicz, J. Physiol. 148, 127 (1959)
- [7] K. Steinmeyer et al. Nature 354, 301 (1991)
- [8] M.C. Koch et al, Science 257, 797 (1992)
- [9] A.L. George et al, Nature Genet. 3, 305 (1993)
- [10] C. Fahlke et al. Proc. Natl. Acad. Sci. USA 94, 2729 (1997)
- [11] G. Stölting et al, Pflügers Arch. Eur. J. Physiol. 466, 2191 (2014)
- [12] C: Fahlke et al, Nature 390, 529 (1997)
- [13] E. Park et al, Nature 541, 500 (2016)
- [14] J.P. Machtens et al, Cell 160, 542 (2015)
- [15] J.C. Jen et al, Neurol. 65, 529 (2005)
- [16] B. de Vries et al, Arch. Neurol. 66, 97 (2009)
- [17] K.D. Choi et al, J. Hum. Genet. 62, 443 (2017)
- [18] K. Iwama et al, J. Hum Genet. doi: 10.1038/s10038-017-0365-z (2017)
- [19] A. Adamczyk et al, Psychiatr. Genet. 21, 90 (2011)
- [20] C.J. van Amen-Hellebrekers et al, Eur. J. Med. Genet. 59, 373 (2016)
- [21] P. Kovermann et al, Sci. Rep. 7, 13913 (2017)
- [22] J. Hotzy et al, J. Biol. Chem. 288, 36492 (2013)
- [23] N. Winter et al, Brain 135, 3416 (2012)
- [24] V. Untiet et al, Glia 65, 388 (2017)

F 5 Physics of the malaria parasite

Ulrich S. Schwarz Institute for Theoretical Physics and BioQuant-Center for Quantitative Biology Heidelberg University

Contents

1	Introduction	2			
2	The malaria lifecycle2.1Skin stage2.2Liver stage2.3Blood stage2.4Other stages	3 3 4 5 6			
3	Gliding motility of sporozoites	6			
4 Mechanics and remodelling of infected red blood cells					
5	Cytoadhesion of infected red blood cells	11			
6	Conclusions and outlook	13			

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Malaria is one of the most devastating diseases that plague mankind [1]. It is caused by a unicellular eukaryotic parasite from the genus *Plasmodium* that is transmitted to humans through the bite of a female *Anopheles* mosquito. Although the malaria incidence rates have gone down significantly over the last years due to improved prevention measures (like use of mosquito nets), according to the latest estimate, in 2015 there were still 212 million cases and 429.000 deaths [2]. Since 2001, a total of 6.8 million malaria deaths have been estimated, with the main victims being under 5 years old children in Africa. Despite many efforts in this direction, there is still no vaccine available against malaria. The malaria parasite is extremely well adapted to its host organisms and its permanent struggle with the human host has strongly shaped our genome. In particular, there are several genetic diseases that in fact are favored by the presence of the malaria parasite, including mutations in the hemoglobin genes (such as the one leading to sickle cell amenia) and hereditary ovalo-, ellipto- and spherocytoses [3, 4].

Although research on malaria is still largely motivated by the search for new therapies and strongly focused on epidemiology, immunology and genetics, during the last decade there has been a growing effort to also address biophysical questions arising in the context of this disease [5, 6, 7, 8]. The relevance of biophysics becomes obvious if one considers the lifecycle of the malaria parasite in the human host, as shown in Fig. 1. It starts with an Anopheles mosquitos injecting several malaria sporozoites into the skin of the host during a blood meal. These then search for blood vessels and use the blood flow to travel to the liver, where one sporozoite can multiply into thousands of merozoites, that then are released into the blood, where they invade red blood cells (RBCs) (top right inset). The infected red blood cells (iRBCs) gets remodeled by the parasite and starts to become adhesive, e.g. to placenta and vascular endothelium (middle and bottom right insets, respectively). This increases the residency time in the vasculature and avoids clearance by the spleen, where RBCs that do not manage to squeeze through the 2 μ m narrow interendothelial slits are sorted out by macrophages from the immune system [9, 10]. After 48 h, the iRBC ruptures and around 20 new merozoites are released, thus closing the asexual cycle. A small portion of the parasites become gametocytes. If taken up by a female mosquito, the parasites go through several mosquito stages, until they are injected again into a human host, so that the full infectious cycle is closed. Including all human and mosquito stages, the complete malaria cycle takes several weeks.

The first obvious questions to address with concepts and method from physics are how the malaria parasite manages to physically move through so many different parts of the human body (mainly skin, liver and blood) and how it invades and remodels compartments in the host (in particular RBCs). Because the medical symptoms of the disease (like fever and anemia) are mainly related to the blood stage, the second set of interesting biophysical questions centers around the way the malaria parasite changes the hydrodynamic movement of RBCs and their interactions with other cells in the vasculature (other RBCs, white blood cells, platelets and vascular endothelial cells). Most of these biophysics questions concern cell mechanics, cell adhesion and motion in hydrodynamic flow, which are well-developed subfields of cellular biophysics. Interestingly, similar biophysics questions as addressed here for the malaria parasite are increasingly asked also for other parasites, including *Toxoplasma*, *Leishmania* or *Trypanosoma* (the causative agents of sleeping sickness) [11, 12].

Here we will review recent progress regarding the biophysics of the malaria parasite. We start with an introduction to the malaria lifecycle and then discuss the most important feature of the skin stage, namely the surprisingly rapid movement through the skin of the host based



Fig. 1: Lifecycle of the malaria parasite in the human body. The malaria parasite is injected into the skin in the form of sporozoites and first replicates in the liver. It then invades red blood cells (RBCs) in the form of merozoites. By replicating within RBCs and then rupturing them, it forms an asexual 48 h cycle in the blood. Some of the parasites become gametocytes and are taken up by another mosquito. In this review we discuss the skin and blood stages in more detail, which are here marked by red circles. Adapted from [1].

with a special mode of locomotion called *gliding motility*. Interestingly, the same machinery underlying gliding motility is also used to invade RBCs. We will then discuss how the parasite remodels the iRBC and its interactions with the environment. In particular, we will discuss why and how it makes the iRBC adhesive (*cytoadherence*) and which consequences this will have for the movement of iRBCs in the vasculature. We finally conclude with a summary of some open questions.

2 The malaria lifecycle

2.1 Skin stage

Several species from the genus *Plasmodium* can transmit malaria, but the most fatal and therefore medically most important one is *Plasmodium falciparum*. A large range of imaging modalities, including confocal, intravital, two-photon, super-resolution, light sheet, electron and atomic force microscopies, have been employed to reveal the details of how the parasite moves and develops over the different stages of the lifecycle [13]. As shown in Fig. 1, the lifecycle starts when during the blood meal of a female mosquito tens of malaria parasites are injected into the skin of the host in the form of crescent-shaped *sporozoites*. A convenient model system to study sporozoite migration is the rodent parasite *Plasmodium berghei*, which does not infect humans. Fig. 2A shows the architecture of a mature sporozoite. Typical values for length, width and radius of curvature are 10 μ m, 1 μ m and 5 μ m, respectively. At the right, one sees the apical polar ring (APR), that defines the front of the cell and through which it secretes various components required for motility and invasion. The shape of the sporozoite is fixed by the inner membrane complex (IMC), a system of flattened vesicles underlying the membrane, and the



Fig. 2: (A) Organization of a sporozoite, the crescent-shaped and highly motile form of the parasite during the skin stage. The cartoon clearly shows the polar structure of the cell, with an apical polar ring (APR) at the front and a posterior polar ring (PPR) at the back. (B) Organization of the iRBC over the 48 h asexual cycle of the blood stage. As the parasite mass grows, it moves into the center. The digested hemoglobin is collected in a food vacuole. Cartoons taken from [13].

basket of microtubules anchored to the APR. The nucleus is located two thirds towards the back of the cell, which is defined by the posterior polar ring (PPR). Invasion and motility is closely related to myosin molecular motors, which together with a system of short actin filaments are located between the IMC and the plasma membrane. Together they effect a continuous flow of adhesion molecules from the APR to the PPR. Once these adhesion molecules engage with ligands in their environment, the cell itself is pushed forward.

Although earlier it had been believed that the sporozoites are directly injected into blood capillaries, today we know that usually they are deposited into the skin [14]. They then move rapidly through the connective tissue, with a typical speed of 1-2 μ m/s and in locally helical trajectories, that arise from their crescent shape [15]. The crescent shape seems to be beneficial for circling around capillaries and finally invading them [16, 17]. In order to appreciate the high speed of these cells, one has to note that the typical speed for keratocyte and fibroblasts, which are the standard model systems for fast and normally migrating animal cells, are 0.2 μ m/s [18] and less than 1 μ m/min [19], respectively. Thus the malaria sporozoite holds the world record for a migrating cell (which however is still much lower than the typical speed range of 10-100 μ m/s for microswimmers such as sperm or flagellated bacteria).

2.2 Liver stage

Once inside a blood capillary, blood flow passively carries the sporozoites towards the liver. There seem to be multiple entry pathways into the liver, including the Kupffer cells from the immune system that connect blood flow and liver. In the liver, malaria parasites multiply inside hepatocytes. One sporozoite is sufficient to have thousands of merozoites being released into the blood stream after 7-10 days, where they start to invade RBCs. At this stage, there are no

clinical symptoms yet of the infection.

2.3 Blood stage

The merozoite in the blood stage is the smallest cell of the lifecyle, with a typical size of 1-2 μ m and an egg shape. The time before RBC-invasion is the only part of the lifecycle in which it is directly exposed to the host immune system, and it lasts only for a few minutes, because then merozoites quickly loose the ability to invade RBCs. Merozoite-invasion of RBCs is an intriguing process and has been studied in great detail as it might provide a way to stop propagation of the disease [20, 21, 22].

The first cartoon in Fig. 2B shows this very first part of the blood stage. After attachement, the merozoite quickly reorients with its apex towards the host membrane, with deformations waves emanating from the contact site [23]. A tight junction forms that during invasion moves over the merozoite within tens of seconds, driven by the same myosin molecular motor that also underlies sporozoite motility. Recently it has been suggested on theoretical grounds that the motor contribution can be relatively small as adhesive interactions with the membrane can account for large parts of the parasite re-orientation and wrapping [24]. Using optical tweezers to control contact between parasite and RBC, it has been shown that the adhesive forces are sufficiently strong to balance 40 pN forces [25]. Once the tight junction has reached the base of the cell, the membrane seals behind the parasite and forms the parasitophorous vacuole that the parasite now uses for its further development. Resealing is followed by another period of dramatic shape changes, during which the RBC forms multiple and evenly spaced projections on its surface (*echinocytosis*). It then returns to its normal biconcave shape within 10 minutes and the parasite starts to develop inside the iRBC.

During the 48 h until the iRBC is ruptured, the parasites produces and exports many proteins through the parasitophorous vacuole membrane (PVM), which together completely remodel the RBC. To feed its own metabolism, but also to convert the interior of the RBC into a normal cytoplasm, the parasite starts to digest hemoglobin . Because this increases osmotic pressure, at the same time the parasite establishes new permeation pathways in the host membrane, to control the osmotic pressure of the iRBC as described by the colloid-osmotic model [26]. It also starts to establish a systems of adhesive knobs on the surface of the iRBCs, by exporting structural proteins like the *knob-associated histidine-rich protein* (KAHRP) [27, 28, 29, 30], that self-assemble into spiral-shaped platforms below the membrane, and the adhesion protein *P. falciparum erythrocyte membrane protein 1* (PfEMP1), that inserts into these and can bind to a large range of extracellular adhesion molecules, including CSA, CD36 and ICAM1 [31, 32]. The spectrin network of the RBC is also completely remodeled: it becomes sparser away from

the knobs and denser around the knobs [33, 34]. The junctional complexes in the spectrin network are dissolved and the actin of the protofilaments is used by the parasite to build actin filaments between the cell surface and newly induced membrane structures (*Maurer's clefts*) in the cytoplasm. Recently it has been argued that these actin filaments are essential transport pathways for the parasite, and that the sickle cell disease protects its carriers from malaria by impairing the build-up of these filaments [35, 4]. Effectively this then leads to reduced cytoadherence and increased clearance by the spleen, as observed earlier [32, 36].

The 48 h development inside the iRBC can be divided into ring (0-24 h), trophozoite (24-36 h) and schizont (40-48 h) stages, as shown in Fig. 2B [37, 13]. During the ring stage, the parasite stays close to the site of invasion, at the rim of the iRBC. At late ring stage, the first knobs start to appear on the iRBC-surface and the iRBC starts to adhere to the blood vessel walls. During

the trophozoite stage, the parasite mass becomes more rounded, moves to the center of the RBC and the end products of the hemoglobin digestion are collected in a growing food vacuole inside the parasite (*hemozoin*). Knob density increases to a value of around $10-30/\mu m^2$ (depending on strain) and their typical diameter is 160 nm. Average spectrin length grows from 42 nm in wildtype to 64 nm in trophozoite [33].

During the schizont stage, the nucleus starts to divide asynchronously in a common syncytium, forming a flower structure with around 20 budding merozoites. Knob density further increases to a value of around $40-60/\mu m^2$, while the diameter goes down to around 100 nm. The average spectrin length increases to 75 nm [33]. The Maurer's clefts initially move freely in the cytoplasm, but later are anchored to the cell surface.

Eventually the schizont rounds up and ruptures to release the new merozoites. Rupture is synchronized across iRBCs and leads to periodic fever in the patient. It has been shown that first the PVM and then the RBC-membrane opens, and that egress has a dispersive character, with merozoites being ejected in less than a second up to 10 μ m into the environment [38]. Surprisingly, during merozoite ejection the host membrane curls away from the opening, indicating that the iRBC has built up some spontaneous curvature towards the outside [39]. Membrane curling is an obvious solution to quickly remove the membrane such that the merozoites have a high chance to encounter and invade nearby RBCs, thus closing the asexual cycle.

2.4 Other stages

The sexual part of the blood stage is not completely understood. Some parasites develop into *gametocytes* that seem to mature in RBCs in the bone marrow, in order to avoid clearance by the spleen. Mature gametocytes have to come back into the vasculature and then seem to be softer [40, 41]. Therefore they might be able to stay longer in the vasculature, until they are taken up during the blood meal of a female mosquito.

Once in the mosquito, female and male gametocytes fuse into *ookinetes*. These transverse the lumen of the mosquito midgut and develop into large *oocysts* on the outside of the mosquito gut. Until the oocysts rupture after approximately one week, hundreds of sporozoites are produced by asexual replication. In oocytes, the sporozoites are not able yet to move individually, but they do so collectively [42]. The sporozoites then move with the hemolymph to the salivary glands, where they acquire the individual ability for gliding motility. Although also very interesting from the biophysics point of view, the mosquito stages are not very well investigated, presumably because their medical relevance is not as large as the one of the stages in the vertebrate hosts.

3 Gliding motility of sporozoites

Like *Toxoplasma*, *Plasmodium* belongs to the genus of *Apicomplexa*, which move over external surfaces in a peculiar mode of locomotion called *gliding motility* [43, 44, 45, 22]. Gliding motility of *Plasmodium* and *Toxoplasma* also share some similarities with the adventurous motility of the bacterium *Myxococcus xanthus*, but there are also essential differences (in particular the chiral nature of the myxococcus gliding machinery) [46, 47, 48].

As shown in Fig. 2A, the sporozoite is a strongly polarized cells. Its apical ring at the front is used to secret various factors, which then are driven backwards over the surface of the cell by retrograde flow. Recently this retrograde flow has been measured by optical tweezer exper-



Fig. 3: (A) A sporozoite that got stuck with the rear end stretches itself until the connection ruptures and motion ensues again. (B) Both a fast (blue) and a slow (red) sporozoite exhibit speed peaks that are characteristic for stick-slip motion. (C) The number of speed peaks correlates with the number of adhesion cycles, demonstrating the close relation between motility and adhesion. Taken from [44].

iments and it was found that it can be much faster than sporozoite motility, namely 15 μ m/s versus 1-2 μ m/s [49]. It is driven by an ancient myosin motor, myoA, that interacts with short actin filaments in the narrow space between the inner membrane complex and the plasma membrane, together forming some kind of active fluid , similar to the actin cytoskeleton of animal cells. Because malaria actin does not polymerise into long filaments and also is not known to branch or crosslink, however, this system seems to be quite different from the retrograde flow usually driving migration of animal cells.

While the crescent shape of sporozoites in tissue leads to locally helical trajectories [15], for single sporozoites on planar substrates the combination of crescent shape and retrograde flow along the cell body leads to circular movement [44]. Motion is usually counterclockwise, presumably because chiral symmetry is broken by the microtubule basket anchored at the apical ring. Closer investigation of sporozoite trajectories revealed that circular motion is not homogeneous, but often interrupted by adhesive events. A typical example is shown in Fig. 3A, where the sporozoite gets stuck at the back (yellow arrow). The front continues to move forward (white arrow), thus the cell body stretches and finally the adhesion at the rear is broken and motion ensues again. Fig. 3B shows that many such speed peaks appear during sporozoite motion, irrespective of the average speed of the parasite (here a fast and a slow parasite are shown in blue and red, respectively). Using reflection interference contrast microscopy (RICM) and traction force microscopy (TFM), it has been shown that small regions of strong adhesion exist between sporozoite and substrate, and that these adhesion sites are highly dynamic. As demonstrated in Fig. 3C, speed peaks correlate with adhesion cycles, demonstrating the close relation between movement and adhesion. Although sporozoites adhere through specific adhesion molecules to their environment, their speed and the stick-slip type motion pattern seem not to depend strongly on the exact nature of the extracellular ligand. In fact such a motion pattern is generic for sliding friction and it has been modelled before also in the context of retrograde



Fig. 4: (A) Different motility modes are observed for sporozoites in pillar assays: (i) circling, (ii) wavering, (iii) linear and (iv) meandering. (B) An agent-based model for sporozoite motility can be used to predict these different motility modes as a function of lattice constant: circling around pillars (red), circling between pillars (blue), linear (green) and meandering (black). Taken from [15].

flow of animal cells [50].

Given the circular movement on planar substrates, it is intriguing to ask how the complex motion patterns arise that one can observe for sporozoites *in vivo*. A surprising answer has been provided by the help of pillar assays [16, 17]. PDMS-pillars with similar radii as sporozoites have been microfabricated and used as obstacle arrays for sporozoite migration. As shown in Fig. 4A, many different motion patterns were observed even for the same geometry. This suggests that sporozoite trajectories are mainly determined by the geometry of their extracellular environment. A simple agent-based model for sporozoite motility was therefore used to simulate sporozoite motility in pillar arrays and indeed gave very similar results, compare Fig. 4B [15]. Here the sporozoite was modelled as a self-propelled particle with curvature, bending energy and excluded volume interaction with the pillars. In addition, it was required to allow for complete re-orientation if collisions could not be resolved by bending, in agreement with experimental observations that sporozoites can buckle and loose substrate contact during collisions. Remarkably, it was also found that sporozoites tend to accumulate around pillars with matching radii, suggesting that their curvature has evolved through the interaction with blood capillaries, which have a similar radius.

4 Mechanics and remodelling of infected red blood cells

RBCs are the most abundant cell type in our body. From the estimated $3.1 \cdot 10^{13}$ cells in our body, $2.6 \cdot 10^{13}$ are RBCs [51]. With an average lifetime of 120 days, this implies that we produce 2.6 million new RBCs every second. For a parasitemia of 10%, there will be $2.6 \cdot 10^{12}$ infected RBCs (iRBCs) in the circulation (ca. 200 g of parasitic mass). In principle, this means that in extreme cases, the malaria parasite can outnumber all other cells in our body, including *E. Coli*, of which humans carries approximately equal numbers as own cells (amounting also to 200 g) [52]. Not only are there so many RBCs, each of them is also an ideal host for the parasite, offering nutrients (mainly hemoglobin, which the malaria parasites digest during the 48 h in the iRBC) and protection from the immune system (especially because the parasite is shield by two membranes, the PVM and the RBC-membrane).

Because RBCs do not have a nucleus and are filled with hemoglobin, their mechanics is mainly determined by the cell envelope, a composite of plasma membrane and spectrin cytoskeleton [53]. Moreover their biological function is strongly shaped by physical factors. For these

reasons, they are attractive model systems for biophysical investigations. In particular, there exists a well-developed mathematical and computational framework to understand their shape and mechanics [54] and well as their deformations and movement in shear flow [55]. The mechanics of their plasma membrane is dominated by local bending energy

$$\mathcal{H}_{bend} = 2\kappa \int dA \left(H - c_0\right)^2 \tag{1}$$

where $dA(u, v) = g(u, v)^{1/2} du dv$ is the integral measure of the surface area, g(u, v) the determinant of the metric tensor, and u and v are the internal coordinates of the surface. H(u, v) is the local mean curvature of the surface, c_0 is the spontaneous curvature and κ is the bending rigidity. In addition one has to take into account the area-difference-elasticity (ADE) arising from a difference in surface areas between the two leaflets

$$\mathcal{H}_{ADE} = \frac{\bar{\kappa}\pi}{2AD^2} \left(\Delta A - \Delta A_0\right)^2 \tag{2}$$

where D is the distance between the two leaflets (typical value 2 nm), $\bar{\kappa}$ the ADE-modulus and the difference in surface area can be calculated as

$$\Delta A = 2D \int dA \ H \ . \tag{3}$$

This shows that spontaneous curvature c_0 and area difference ΔA_0 have a similar effect and cannot be extracted independently of each other.

The mechanics of the spectrin-actin cytoskeleton underlying the plasma membrane is described by thin shell elasticity

$$\mathcal{H}_{elast} = \int dA \left(\frac{K_{\alpha}}{2} \alpha^2 + \mu \beta \right) \tag{4}$$

with the 2D stretch modulus K_{α} and the 2D shear modulus μ . Area and shear strain, respectively, follow from the principal extension ratios λ_1 and λ_2 of a deformed ellipse as

$$\alpha = \lambda_1 \lambda_2 - 1, \ \beta = \frac{1}{2} \left(\frac{\lambda_1}{\lambda_2} + \frac{\lambda_2}{\lambda_1} - 2 \right) \ . \tag{5}$$

To evalute the elastic energy, one has to define a reference shape. Additionally one can consider higher order terms, which become important at large deformations, where typically strain stiffening occurs due to the polymeric nature of the spectrin network. Finally bending and thin shell elasticity energies have to be complemented by Lagrange parameters to enforce constant area and volume.

Wildtype RBCs have a typical surface area of $A = 140 \ \mu m^2$ and a typical volume of $V = 100 \ \mu m^3 = 100 \ \text{fl}$ [57, 56]. The surface area of a sphere with the same volume would be $A = 104 \ \mu m^2$, thus the RBC has an excess area over the equivalent sphere of around 40%. Another way to express this important relation is to define the reduced volume

$$v = \frac{V}{V_0} = \frac{6\sqrt{\pi V}}{A^{3/2}}$$
(6)

which is the real volume V in relation to the volume of a sphere with the same area A. The RBC has v = 0.64, again indicating the high degree of excess area. The classical values for the elastic parameters are $\kappa \approx 50 \ k_B T$, $K_{\alpha} \approx 25 \ \kappa/\mu m^2$ and $\mu \approx K_{\alpha}/2 = 2.5 \ \mu N/m$ [58].



Fig. 5: (A) Shape of iRBC and parasite mass inside over the complete time of the 48 h asexual cycle as extracted by image processing from confocal stacks. One clearly sees that the parasite rounds up and moves towards the center of the iRBC, and that the iRBC itself also becomes round. (B-D) Time course of surface area and volume of iRBCs extracted from the image processed data. While surface area stays roughly constant, volume goes up by 60%. Solid lines are the predictions of the colloid-osmotic model. Adapted from [56].

Impressively, this theoretical framework gives rise to a complete understanding of the large zoo of RBC-shapes, including the stomatocyte-discocyte-echinocyte sequence arising from changing differential area [58] and the echinocytic shapes that arise as bilayer budding effected by large spontaneous curvature or differential area, but suppressed by elasticity [59]. For the wild-type RBC with the values reported above, the discocyte arises as the stable solution mainly due to the bending energy at reduced volume v = 0.64. The interface Hamiltonian for RBCs also suggests that the echinocytosis observed after merozoite invasion is related to some change in membrane composition, which has this global effect on RBC-shape.

It is a long-standing question how RBC standard shape is changed during the course of a malaria-infection. As shown in Fig. 5, the time course of the shape of iRBCs (A) and from this also the time courses of area (B) and volume (C,D) as a function of developmental time recently have been measured with high resolution [56]. In general, these measurements confirm earlier results that the iRBC starts to round up at around 20 h post invasions, at the same time when the parasite mass starts to grow and to move into the center. In regard to area and volume, it was found that surface area A is relatively constant, but that volume V increased by 60% (that is to $V = 160 \ \mu m^3$) from late ring to schizont, in very good agreement with the predictions of the colloid-osmotic model [26], but in contrast to earlier work that reported a reduction of A at relatively constant V [60, 61]. In particular, these data imply that the schizont has a reduced volume v close to 1 and thus can be modeled as a round cell in regard to its movement in hydrodynamic flow.

Another central issue are the values of the elastic parameters defined by the composite interface Hamiltonian introduced above. By fitting a multiscale model similar to the above continuum model to experimental deformation data from optical tweezer experiments, it was found that the


Fig. 6: (A) Representation of a biconcave wildtype RBC as a triangulated surface in a multiscale model that incorporates both bending and elastic energies. Taken from [55]. (B) Experimental and simulated force-deformation curves for stretched RBCs. One clearly sees that iRBCs become much stiffer as they develop from wildtype through trophozoite to schizont. Taken from [62].

wildtype 2D shear modulus should be $\mu = 8.3 \,\mu N/m$ and that the bending modulus κ should be larger than 50 k_BT [63]. A similar result, $\mu = 4.73 \,\mu N/m$ and κ in the range of 100 k_BT , was found by a multiscale model that also incorporates dynamical effects [64]. Fitting stretching data for iRBCs, it was found that the shear modulus μ increases to $14.5 \,\mu N/m$, $29 \,\mu N/m$ and $40 \,\mu N/m$ for ring, trophozoite and schizont stages [62], compare Fig. 6. This dramatic increase in stiffness underscores the fact that the parasite is under large evolutionary pressure to avoid passage through the spleen, where stiff RBCs are sorted out. Based on the AFM imaging data of the changes in the spectrin network over the 48 h asexual cycle [33], a multiscale model has been parametrized that suggests that the stiffness increase results mainly from the vertical links between the knobs and the spectrin network [34].

For the movement of RBCs in shear flow, we also have to know its viscoelastic properties. Their mean hemoglobin concentration is 33 g/dl and leads to an intracellular viscosity of $6 \cdot 10^{-3}$ Pa s, five times higher than the viscosity of the surrounding blood plasma. Higher hemoglobin concentrations would lead to a strong increase in viscosity. Due to ageing, RBCs loose surface area and volume, but not hemoglobin, leading to a strong reduction in cell deformability and removal in the spleen. These observations might explain why the malaria parasite seems to digest more hemoglobin than needed for its own metabolism.

5 Cytoadhesion of infected red blood cells

Malaria parasites induce cytoadherence of iRBCs in order to increase residency time in the vasculature and to avoid clearance by the spleen. To this end, they export proteins like KAHRP and PfEMP1 that self-assembly into thousands of adhesive knobs on the surface of the iRBC. While the diameter of these knobs becomes smaller during the asexual cycle, its height stays constant at a value of 10-20 nm, as measured with SEM and AFM [65, 66, 67]. Strikingly,

this design is similar to the one evolved by leukocytes, which use rolling adhesion to scan the endothelium for signals of inflammation [68]. To this end, they localize adhesion molecules from the selectin family to the tips of hundreds of microvilli covering their surfaces. Therefore it has been suggested that this common design of multiple adhesive protrusions is the result of an optimalisation process for cell capture and adhesion under flow conditions [69, 70, 71]. Another striking observation is the observation that iRBC-adhesion is flow-enhanced [72], as also known from leukocytes [73] and bacteria [74]. Often such behaviour results from molecular catch bonds (whose lifetime increases with force, in contrast to a decrease for the usual slip bonds) and indeed a very recent study has reported that the PfEMP1:ICAM1 bond has this property [75].

To model cytoadhesion of iRBCs in shear flow, one first needs to choose a suitable method to describe movement of the cell in hydrodynamic flow . For spherical cells, such a method has been introduced by Hammer and coworkers, who simulated the phase diagram of rolling adhesions for leukocytes [76, 77]. In later studies, this method has been extended to also resolve the receptors on the cell surface and the ligands on the substrates [70, 78, 79]. Recently this approach has also been applied to cytoadhesion of schizonts [80], because these can be considered to be round cells, compare Fig. 5. In essence, adhesive dynamics for round cells is the simulation of a Langevin equation

$$\partial_{t}X(t) = u^{\infty} + M\{F_{\rm S} + F_{\rm D}\} + k_{\rm B}T\nabla M + \xi(t),\tag{7}$$

where X(t) is a six-dimensional vector describing translation and rotation of the spherical cell. M is a mobility matrix that can be calculated semi-analytically from the solution of the Stokes equation for a sphere above a wall. u^{∞} is the imposed linear shear flow and $F_{\rm S}$ and $F_{\rm D}$ are shear and direct forces, respectively, with the former also resulting from the Stokes equation and the latter including the adhesion forces. $\xi(t)$ is the usual random force/torque and the term with ∇M arises due to the multiplicative noise. Additional model parameters are the rules for bond association and dissociation. Usually one assumes a constant on-rate below a typical encounter distance and an off-rate that depends exponentially on force according to the Bell-Evans-model for slip bonds. Finally one has to define the exact distribution of receptors and ligands. As shown in Fig. 7A, for schizonts one can model the clustering of PfEMP1-molecules into knobs, while the corresponding ligands are distributed with a typical distance on the substrate. In Fig. 7B it is shown that with this model, one can achieve good agreement with flow chamber experiments for rolling velocity as a function of shear rate. In particular, this work predicts that a typical number of PfEMP1-molecules per knob (*multiplicity*) should be six, in agreement with earlier estimates [81].

In order to also describe the other stages of the asexual cycle, when the iRBCs are not spherical, one has to implement a hydrodynamic method that can also deal with deformable cells. The same challenge in fact arises also for wildtype RBCs. One approach often applied to simulate blood flow is the Lattice Boltzmann Method (LBM) [82]. More recently, however, the movement of RBCs in shear flow has been simulated also with other methods, in particular with Multiparticle Collision Dynamics (MPCD) [83] and Dissipative Particle Dynamics (DPD) [84, 64]. These hydrodynamic methods are then coupled to the elasticity of the RBC, compare Fig. 6A, often implementing a multiscale model that in the continuum limit becomes the interface Hamiltonian described above. These approaches can predict many effects observed in blood flow, in particular the parachute shape of single RBCs in channel flow, the Fahraeus-Lindqvist effect (decrease of apparent viscosity with decreasing channel diameter) and the margination of white blood cells [55]. Cytoadhesion of trophozoites has been simulated with the DPD-approach and





Fig. 7: (A) Model for a round cell adhering in linear shear flow through adhesive knobs on its surface. Shear flow generates both a translational force F_s and a torque T_s . A fluorescent parasite mass (P) can be used to characterize rotation because it will oscillate as it goes in and out of the focal plane (FP). (B) Comparison of experimental and simulation data for rolling velocity as a function of shear rate with reasonable corridors for the model parameters (light shaded: 100 to 400 nm in ligand distance; dark shaded: 0.1 to 10 Hz in on-rate). Taken from [80].

revealed that these should flip rather than roll in shear flow due to their biconcave shape and smaller stiffness [62, 85], as recently indeed confirmed by flow chamber experiments [80].

6 Conclusions and outlook

As it is true for all pathogens, the malaria parasite is both frightening and fascinating to any researcher who studies the ways in which it interacts with its hosts. One of the most surprising aspects of the lifecycle is the observation how strongly the organization of the parasite depends on the environment with which it interacts. For example, when one compares the architecture of the sporozoite from Fig. 2A with the one of the merozoite from Fig. 2B, one cannot help to be surprised by the much larger size (10 versus 1 μ m) and completely different shape (crescent versus egg) of the two variants, which both arise from the same genome. Obviously these different architectures correspond to very different functions, namely fast motility in the skin versus invasion of RBCs in the blood. Evolutionary adaptability also becomes apparent in many other cellular functions described here. In particular, we have seen that the motion patterns of sporozoites are strongly shaped by their environment, as revealed by the pillar assays, and that the cytoadhesion of iRBCs mimics the way leukocytes interact adhesively with the endothelium, as confirmed by the adhesive dynamics simulations.

Although the malaria parasite is very special due to its unique lifecycle, the main biophysical questions that arise during its investigation are very similar to the ones that one also studies for other cell types. Here we have focused on two main stages in the human host, namely the skin and blood stages. For the sporozoite, which seems to be optimised for fast motion through the skin in search for blood vessels, we have seen that fast motility is achieved by retrograde flow of surface-anchored adhesins, and that simple physics models for sliding friction and self-propelled particles can explain some of its peculiar motion features. One open question in this context is the question how retrograde flow is accomplished by the interplay between the

myosin motor and actin, which in contrast to the lamellipodium of migrating animal cells does not seem to form long and branched or crosslinked filaments. Another open question in the sporozoite field is the question why the parasite has evolved very specific adhesion molecules when *in vitro* it does not need any special ligands to achieve its high level of motility.

Merozoites seem to be optimised for efficient invasion of RBCs and lead to a complete remodelling of the iRBC. Because RBCs are an extremely well-studied subject in biophysics, including mature theories for their shape, mechanics and movement in hydrodynamic flow, the biophysics of iRBCs has developed as a very fruitful sub-area of this large field. One additional aspect not present in healthy RBCs is cytoadherence of iRBCs, which has led to the development of adhesive dynamics simulations of both round and deformable cells. At the current stage, one of the main challenges is to establish a tighter connection to the underlying molecular processes through multiscale approaches. For example, it is still unclear how exactly the parasite remodels the spectrin network, how it controls transport through the cytoplasm and the different membranes, which variants of the adhesion molecule PfEMP1 are expressed in which context, how the endothelium reacts to different variants of iRBCs, and if and how iRBCs interact with other cells in the blood flow.

Finally we note that the liver stage, the sexual blood stage and the mosquito stages as described above are still largely *terra incognita* from the viewpoint of biophysics. Without doubt, the parasite has evolved unexpected solutions also for these stages. As is the case with all pathogens, one can only hope that understanding these mechanisms in more details also will provide better ways to fight this deadly disease.

Acknowledgments: I would like to thank all past and present members of the Heidelberg community working on the biophysics of malaria for helpful discussions and enjoyable collaborations, in particular Friedrich Frischknecht, Joachim Spatz, Sylvia Münter, Benedikt Sabass, Christine Selhuber, Anna Battista, Michael Lanzer, Anil Kumar Dasanna, Motomu Tanaka and Julia Jäger. This work was supported by the DFG Collaborative Research Center 1129 on Integrative analysis of pathogen replication and spread at Heidelberg.

References

- [1] L. H. Miller, D. I. Baruch, K. Marsh, and O. K. Doumbo, Nature 415(6872), 673 (2002).
- [2] World health organization (WHO), World malaria report 2016.
- [3] B. M. Cooke, N. Mohandas, and R. L. Coppel, Seminars in Hematology 41(2), 173 (2004).
- [4] M. Cyrklaff, C. P. Sanchez, F. Frischknecht, and M. Lanzer, Trends in Parasitology 28(11), 479 (2012).
- [5] S. Suresh, J. Spatz, J. P. Mills, A. Micoulet, M. Dao, C. T. Lim, M. Beil, and T. Seufferlein, Acta Biomaterialia 1, 15 (2005).
- [6] G. Y. H. Lee and C. T. Lim, Trends in Biotechnology 25(3), 111 (2007).
- [7] D. A. Fedosov, Drug Discovery Today: Disease Models 16, 17 (2015).
- [8] U. S. Schwarz, Seminars in Cell & Developmental Biology 46(Supplement C), 82 (2015).
- [9] R. E. Mebius and G. Kraal, Nature Reviews Immunology 5(8), 606 (2005).
- [10] I. V. Pivkin, Z. Peng, G. E. Karniadakis, P. A. Buffet, M. Dao, and S. Suresh, Proceedings of the National Academy of Sciences 113(28), 7804 (2016).

- [11] J. Baum and F. Frischknecht, Seminars in Cell & Developmental Biology 46(Supplement C), 78 (2015).
- [12] A. Hochstetter and T. Pfohl, Trends in Parasitology $\mathbf{0}(0)$ (2016).
- [13] M. De Niz, P.-C. Burda, G. Kaiser, H. A. del Portillo, T. Spielmann, F. Frischknecht, and V. T. Heussler, Nature Reviews Microbiology 15(1), 37 (2017).
- [14] R. Amino, S. Thiberge, B. Martin, S. Celli, S. Shorte, F. Frischknecht, and R. Mnard, Nature Medicine 12(2), 220 (2006).
- [15] A. Battista, F. Frischknecht, and U. S. Schwarz, Physical Review E 90(4), 042720 (2014).
- [16] J. K. Hellmann, S. Mnter, M. Kudryashev, S. Schulz, K. Heiss, A.-K. Mller, K. Matuschewski, J. P. Spatz, U. S. Schwarz, and F. Frischknecht, PLoS Pathog 7(6), e1002080 (2011).
- [17] M. J. Muthinja, J. Ripp, J. K. Hellmann, T. Haraszti, N. Dahan, L. Lemgruber, A. Battista, L. Schtz, O. T. Fackler, U. S. Schwarz, J. P. Spatz, and F. Frischknecht, Advanced Healthcare Materials pp. n/a–n/a (2017).
- [18] K. Keren, Z. Pincus, G. M. Allen, E. L. Barnhart, G. Marriott, A. Mogilner, and J. A. Theriot, Nature 453(7194), 475 (2008).
- [19] M. Dembo and Y.-L. Wang, Biophysical Journal 76(4), 2307 (1999).
- [20] A. F. Cowman and B. S. Crabb, Cell **124**(4), 755 (2006).
- [21] A. F. Cowman, D. Berry, and J. Baum, J Cell Biol 198(6), 961 (2012).
- [22] I. Tardieux and J. Baum, The Journal of Cell Biology 214(5), 507 (2016).
- [23] P. R. Gilson and B. S. Crabb, International Journal for Parasitology 39(1), 91 (2009).
- [24] S. Dasgupta, T. Auth, N. Gov, T. Satchwell, E. Hanssen, E. Zuccala, D. Riglar, A. Toye, T. Betz, J. Baum, and G. Gompper, Biophysical Journal 107(1), 43 (2014).
- [25] A. Crick, M. Theron, T. Tiffert, V. Lew, P. Cicuta, and J. Rayner, Biophysical Journal 107(4), 846 (2014).
- [26] J. M. A. Mauritz, A. Esposito, H. Ginsburg, C. F. Kaminski, T. Tiffert, and V. L. Lew, PLoS Comput Biol 5(4), e1000339 (2009).
- [27] L. G. Pologe, A. Pavlovec, H. Shio, and J. V. Ravetch, Proceedings of the National Academy of Sciences 84(20), 7139 (1987).
- [28] B. S. Crabb, B. M. Cooke, J. C. Reeder, R. F. Waller, S. R. Caruana, K. M. Davern, M. E. Wickham, G. V. Brown, R. L. Coppel, and A. F. Cowman, Cell 89(2), 287 (1997).
- [29] H. Weng, X. Guo, J. Papoin, J. Wang, R. Coppel, N. Mohandas, and X. An, Biochimica et Biophysica Acta (BBA) - Biomembranes 1838(1, Part B), 185 (2014).
- [30] J. M. Watermeyer, V. L. Hale, F. Hackett, D. K. Clare, E. E. Cutts, I. Vakonakis, R. A. Fleck, M. J. Blackman, and H. R. Saibil, Blood pp. blood-2015-10-674002 (2015).
- [31] D. I. Baruch, B. L. Pasloske, H. B. Singh, X. Bi, X. C. Ma, M. Feldman, T. F. Taraschi, and R. J. Howard, Cell 82(1), 77 (1995).
- [32] R. M. Fairhurst, D. I. Baruch, N. J. Brittain, G. R. Ostera, J. S. Wallach, H. L. Hoang, K. Hayton, A. Guindo, M. O. Makobongo, O. M. Schwartz, A. Tounkara, O. K. Doumbo, *et al.*, Nature **435**(7045), 1117 (2005).
- [33] H. Shi, Z. Liu, A. Li, J. Yin, A. G. L. Chong, K. S. W. Tan, Y. Zhang, and C. T. Lim, PLOS ONE 8(4), e61170 (2013).
- [34] Y. Zhang, C. Huang, S. Kim, M. Golkaram, M. W. A. Dixon, L. Tilley, J. Li, S. Zhang, and S. Suresh, Proceedings of the National Academy of Sciences 112(19), 6068 (2015).
- [35] M. Cyrklaff, C. P. Sanchez, N. Kilian, C. Bisseye, J. Simpore, F. Frischknecht, and M. Lanzer, Science 334(6060), 1283 (2011).
- [36] R. M. Fairhurst, C. D. Bess, and M. A. Krause, Microbes and Infection 14(10), 851 (2012).

- [37] L. Bannister, J. Hopkins, R. Fowler, S. Krishna, and G. Mitchell, Parasitology Today 16(10), 427 (2000).
- [38] M. Abkarian, G. Massiera, L. Berry, M. Roques, and C. Braun-Breton, Blood 117(15), 4118 (2011).
- [39] A. Callan-Jones, O. AlbarranArriagada, G. Massiera, V. Lorman, and M. Abkarian, Biophysical Journal 103(12), 2475 (2012).
- [40] M. Aingaran, R. Zhang, S. K. Law, Z. Peng, A. Undisz, E. Meyer, M. Diez-Silva, T. A. Burke, T. Spielmann, C. T. Lim, S. Suresh, M. Dao, *et al.*, Cellular Microbiology 14(7), 983 (2012).
- [41] M. Dearnley, T. Chu, Y. Zhang, O. Looker, C. Huang, N. Klonis, J. Yeoman, S. Kenny, M. Arora, J. M. Osborne, R. Chandramohanadas, S. Zhang, *et al.*, Proceedings of the National Academy of Sciences **113**(17), 4800 (2016).
- [42] D. Klug and F. Frischknecht, eLife 6, e19157 (2017).
- [43] L. D. Sibley, Science 304(5668), 248 (2004).
- [44] S. Muenter, B. Sabass, C. Selhuber-Unkel, M. Kudryashev, S. Hegge, U. Engel, J. P. Spatz, K. Matuschewski, U. S. Schwarz, and F. Frischknecht, Cell Host & Microbe 6(6), 551 (2009).
- [45] M. B. Heintzelman, Seminars in Cell & Developmental Biology 46(Supplement C), 135 (2015).
- [46] T. Mignot, J. W. Shaevitz, P. L. Hartzell, and D. R. Zusman, Science 315(5813), 853 (2007).
- [47] R. Balagam, D. B. Litwin, F. Czerwinski, M. Sun, H. B. Kaplan, J. W. Shaevitz, and O. A. Igoshin, PLOS Computational Biology 10(5), e1003619 (2014).
- [48] S. T. Islam and T. Mignot, Seminars in Cell & Developmental Biology 46(Supplement C), 143 (2015).
- [49] K. A. Quadt, M. Streichfuss, C. A. Moreau, J. P. Spatz, and F. Frischknecht, ACS Nano 10(2), 2091 (2016).
- [50] B. Sabass and U. S. Schwarz, Journal of Physics: Condensed Matter 22, 194112 (2010).
- [51] R. Milo and R. Phillips, *Cell Biology by the Numbers* (Garland Science, 2015), google-Books-ID: 9NPRCgAAQBAJ.
- [52] R. Sender, S. Fuchs, and R. Milo, PLOS Biology 14(8), e1002533 (2016).
- [53] S. E. Lux, Blood 127, 187 (2015).
- [54] Lim, H.W.G., Wortis, M., and Mukhopadhyay. R., in *Soft Matter* (Wiley-VCH Verlag GmbH & Co. KGaA, 2008), vol. 4, pp. 83–249.
- [55] D. A. Fedosov, H. Noguchi, and G. Gompper, Biomechanics and Modeling in Mechanobiology 13(2), 239 (2014).
- [56] M. Waldecker, A. K. Dasanna, C. Lansche, M. Linke, S. Srismith, M. Cyrklaff, C. P. Sanchez, U. S. Schwarz, and M. Lanzer, Cellular Microbiology 19(2), 1 (2017).
- [57] N. Mohandas and P. G. Gallagher, Blood 112(10), 3939 (2008).
- [58] G. L. H. W, M. Wortis, and R. Mukhopadhyay, Proceedings of the National Academy of Sciences 99(26), 16766 (2002).
- [59] R. Mukhopadhyay, H. W. Gerald Lim, and M. Wortis, Biophysical Journal 82(4), 1756 (2002).
- [60] A. Esposito, J.-B. Choimet, J. N. Skepper, J. M. A. Mauritz, V. L. Lew, C. F. Kaminski, and T. Tiffert, Biophysical Journal 99(3), 953 (2010).
- [61] I. Safeukui, P. A. Buffet, S. Perrot, A. Sauvanet, B. Aussilhou, S. Dokmak, A. Couvelard, D. C. Hatem, N. Mohandas, P. H. David, O. Mercereau-Puijalon, and G. Milon, PLOS

ONE 8(3), e60150 (2013).

- [62] D. A. Fedosov, B. Caswell, S. Suresh, and G. E. Karniadakis, Proceedings of the National Academy of Sciences 108(1), 35 (2011).
- [63] J. Li, M. Dao, C. T. Lim, and S. Suresh, Biophysical Journal 88(5), 3707 (2005).
- [64] D. A. Fedosov, B. Caswell, and G. E. Karniadakis, Biophysical Journal 98(10), 2215 (2010).
- [65] J. Gruenberg, D. R. Allred, and I. W. Sherman, The Journal of Cell Biology 97(3), 795 (1983).
- [66] E. Nagao, O. Kaneko, and J. A. Dvorak, Journal of Structural Biology 130(1), 34 (2000).
- [67] K. A. Quadt, L. Barfod, D. Andersen, J. Bruun, B. Gyan, T. Hassenkam, M. F. Ofori, and L. Hviid, PLoS ONE 7(9), e45658 (2012).
- [68] T. Springer, Cell **76**(2), 301 (1994).
- [69] M. Ho, M. J. Hickey, A. G. Murray, G. Andonegui, and P. Kubes, The Journal of Experimental Medicine 192(8), 1205 (2000).
- [70] C. Korn and U. S. Schwarz, Physical Review Letters 97(13), 138103 (2006).
- [71] G. Helms, A. K. Dasanna, U. S. Schwarz, and M. Lanzer, FEBS Letters 590(13), 1955 (2016).
- [72] H. Rieger, H. Y. Yoshikawa, K. Quadt, M. A. Nielsen, C. P. Sanchez, A. Salanti, M. Tanaka, and M. Lanzer, Blood 125(2), 383 (2015).
- [73] R. Alon, D. A. Hammer, and T. A. Springer, Published online: 06 April 1995; doi:10.1038/374539a0 374(6522), 539 (1995).
- [74] W. E. Thomas, E. Trintchina, M. Forero, V. Vogel, and E. V. Sokurenko, Cell 109(7), 913 (2002).
- [75] Y. B. Lim, J. Thingna, J. Cao, and C. T. Lim, Scientific Reports 7(1), 4208 (2017).
- [76] D. A. Hammer and S. M. Apte, Biophysical Journal 63(1), 35 (1992).
- [77] K.-C. Chang, D. F. J. Tees, and D. A. Hammer, Proceedings of the National Academy of Sciences 97(21), 11262 (2000).
- [78] C. B. Korn and U. S. Schwarz, The Journal of Chemical Physics 126(9), 095103 (2007).
- [79] C. B. Korn and U. S. Schwarz, Physical Review E 77(4), 041904 (2008).
- [80] A. K. Dasanna, C. Lansche, M. Lanzer, and U. S. Schwarz, Biophysical Journal 112(9), 1908 (2017).
- [81] X. Xu, A. K. Efremov, A. Li, L. Lai, M. Dao, C. T. Lim, and J. Cao, PLoS ONE 8(5), e64763 (2013).
- [82] C. Sun, C. Migliorini, and L. L. Munn, Biophysical Journal 85(1), 208 (2003).
- [83] H. Noguchi and G. Gompper, Proceedings of the National Academy of Sciences of the United States of America 102(40), 14159 (2005).
- [84] I. V. Pivkin and G. E. Karniadakis, Physical Review Letters 101(11), 118105 (2008).
- [85] D. Fedosov, B. Caswell, and G. Karniadakis, Biophysical Journal 100(9), 2084 (2011).

F 6 Alzheimer's disease

J. Kutzsche, D. Willbold Institute of Complex Systems (ICS-6) Forschungszentrum Jülich GmbH

Contents

1	Introduction	2
2	Neuropathology of AD	2
3	Amyloid cascade hypothesis	3
4	Treatment strategies	4
5	D-enantiomeric peptides	5
Refe	erences	11

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

1 Introduction

Alzheimer's disease (AD) is the most common type of dementia and currently more than 24 million people are affected worldwide. AD was first described from the German neuropathologist and psychiatrist Alois Alzheimer in 1907¹. The disease is characterized by a progressive decline in multiple areas of function, including memory, thinking, reasoning, communication, behaviour and the skills needed to carry out daily activities. AD leads to the progressive neuronal cell death and loss of tissue throughout the brain. The neurodegeneration results in a dramatically shrinkage of the brain over time and affects nearly all its functions. The cortex shrivels up, which leads to damage of areas, which are involved in thinking, planning and remembering. The shrinkage is especially severe in the hippocampus, an area of the cortex that plays a key role in formation of new memories.



Figure 1: Comparison of an healthy brain with the brain of an advanced Alzheimer's patient. ©2017 Alzheimer's Association. <u>www.alz.org</u>. All rights reserved. Illustrations by Stacy Jannis.

2 Neuropathology of AD

The key pathological hallmarks of AD are extracellular senile plaques (SPs) and intraneuronal neurofibrillary tangles (NFTs). The SPs are mainly composed of insoluble aggregated β -amyloid peptide (A β)². They occur in different morphological forms including diffuse, neuritic and dense-cored plaques³ and consist of different forms of A β peptide. A β is produced by the sequential processing of the amyloid precursor protein (APP) by β and γ -secretases. The cleavage of APP by the γ -secretase is somewhat imprecise, resulting in a C-terminal heterogeneity of the resulting peptide population. Therefore, multiple different A β species exist. A β 40 (ending at position 40) is the most abundant A β species, followed by A β 42. A β 42, which is more hydrophobic and fibrillogenic than A β 40 and is the principal species deposited in AD brains⁴.

Neurofibrillary tangles are composed of hyperphosphorylated forms of the microtubuleassociated protein (MAP) tau. Phosphorylated tau proteins accumulate early in neurons, even before formation of neurofibrillary tangles. The accumulation of neurofibrillary tangles and phosphorylated tau species is associated with disturbances of the microtubule network and, as a consequence of the axoplasmic transport⁵ (the transport of cargo to the nerve endings through cytoskeleton). The number of neurofibrillary tangles is tightly linked to the degree of dementia⁵.



Figure 2: Postmortem tissue sample from an AD patient brain reveals AD pathology including amyloid-beta plaques and Tau tangles. Taken from petridishtalk.com

3 Amyloid cascade hypothesis

Two observations resulted in the original formulation of the "amyloid cascade" hypothesis. First the detection of A β peptide as a main constituent of senile plaques⁶ and second the observation that humans with Down syndrome invariably develop neuropathological typical AD, due to the presence of three copies of the wild-type APP gene on chromosome 21. On the basis of these observations Hardy and Higgins postulated that the deposition of A β is the causative agent of AD pathology and that the neurofibrillary tangles, cell loss, vascular damage and dementia follow as a direct result of this deposition⁷. This hypothesis was further supported by different genetic studies of families with early-onset AD, in which mutations in the APP gene⁸ (substrate) or the genes for presenting 1 and 2 $^{9-11}$ ENREF 4 (essential components of the γ -secretase complexes¹²) were identified, which lead to an altered proteolytic processing of APP and subsequently to the accumulation of AB. Additionally apolipoprotein E4 (ApoE4), which is a predisposition in more than 40 % of the AD cases, has been described to impair the clearance of A β from the brain¹³ and increases the oligomerization of A β , which leads to increased A β oligomer levels in the brain¹⁴. So, mounting evidence from genetic, pathological, and functional studies have shown that an imbalance between the production and clearance of $A\beta$ in the brain results in accumulation and aggregation of A β . The resulting A β assemblies include soluble low molecular weight oligomers, protofibrils, and insoluble, fibrillar aggregates. Due to the fact that A β deposition is only weakly related to the degree of dementia¹⁵, the amyloid cascade hypothesis was controversially discussed for some time. However, studies over the more recent years have strengthened the hypothesis that soluble and freely diffusible AB oligomers are the major neurotoxic agent responsible for disease development and progression. It was shown that soluble Aβ42 oligomers isolated from AD patients' brains decreased the number of synapses¹⁶, inhibited long-term potentiation, and enhanced long-term synaptic depression in rodent hippocampus. Furthermore, injection of A β 42 oligomers into healthy rat brains lead to an impairment of memory. The human oligomers also induce hyperphosphorylation of tau at AD-relevant epitopes¹⁷ and cause neuritic dystrophy in cultured neurons. Crossing human APP with human tau transgenic mice enhances tau-positive neurotoxicity. Furthermore, there is more and more evidence that $A\beta$ oligomers or at least sub-fractions of them, are able to

replicate in a prion-like fashion¹⁸. It was shown that the inoculation of A β aggregatecontaining brain homogenates into aggregate-free tg AD mice brains overexpressing wildtype A β resulted in substantially increased A β deposition. This suggests an involvement of a seed-induced deposition mechanism¹⁹.

4 Treatment strategies

Despite intensive efforts in drug development, no preventive or curative treatment for AD has been achieved vet²⁰. Currently only two types of medications for the symptomatic treatment of AD are approved, cholinesterase inhibitors (Donepezil, Rivastigmine and Galantamine) and one N-methyl-D-aspartate (NMDA) receptor antagonist (Memantine). Treatment with these compounds does not slow down disease progression. Instead, only symptoms are treated with quite some risk for unpleasant side activities. Since 2003, there have not been any new drugs approved. Over the last decades different therapeutic approaches were pursued to establish a causative treatment against AD. Agents acting upon A β , such as vaccines, antibodies and inhibitors or modulators of γ - and β -secretase; agents directed against the tau protein as well as compounds acting as antagonists of neurotransmitter systems (serotoninergic 5-HT6 and histaminergic H3) were developed. During this time we have witnessed many failures of approaches that were aiming to reduce A β formation (e.g. secretase inhibitors or modulators) or increase AB clearance (e.g. active or passive immunization against A β species). Decreasing of total A β may possibly be suitable to slow down A β oligomer formation and thus be potentially helpful to prevent disease development. Once A β oligometric are formed, however, the reduction of total A β might not be enough anymore to slow down disease progression. Therefore, it is essential to eliminate already existing AB oligomers by converting them into AB species that are not toxic and cannot support prion-like replication of $A\beta$ oligomers anymore. Passive immunization with antibodies that are described as AB oligomer specific is one option that relies on the assumption that AB oligomers, which are bound by those antibodies, are subsequently degraded by microglia and astrocytes or any other component of the immune system. It is important to keep in mind that this is not necessarily true. A poligomers, with or without antibodies bound to them, are still cytotoxic. Also, it is not clear, whether these antibodies bind the relevant A β oligomer species at all, because it is just not known yet, what the disease-relevant oligomers look like. Also, as it is known from prion protein prion strains, that there might be an issue with development of resistance to antibodies, which are specific for a certain AB oligomer species. Therefore, our approach aims to specifically eliminate all toxic oligomers by stabilizing A β monomers in an aggregation incompetent conformation. This shifts any equilibrium between A β species towards A β monomers and away from A β oligomers.



5 D-enantiomeric peptides

D-enantiomeric peptides are less immunogenic and proteolytically more stable than Lpeptides, due to the fact that proteases are stereoselective for L-amino acid peptide bonds. Because of these advantageous properties, D-peptides are candidates for oral treatments.

In our institute a mirror image phage display approach was used to identify novel and highly specific ligands for the stabilization of A β monomers. A randomized 12-mer peptide library presented on M13 phages was screened for peptides with binding affinity for the mirror image of A β 1-42²¹. In one of these selections we have identified the fully D-enantiomeric peptide D3 (amino acid sequence rprtrlhthrnr, Figure 3).



Figure 3: Lewis structure of the D-enantiomeric peptide D3.

For this lead compound "D3" we were able to show that it specifically eliminates $A\beta$ oligomers and converts them into non-toxic species²². *In vivo*, we have shown that D3 is able to enhance cognition and reduce plaque load in several transgenic mouse models even after oral application²²⁻²⁵.



Figure 4: D3 treatment effect on the cognitive impairment of APP PS1 mice in the Morris water maze (MWM). Four months old APP PS1 mice had an Alzet minipump implanted under the skin with tubing inserted into the stomach, delivering the D3 peptide for the treated group (n=8), and delivering saline solution for the control group (n=9). The escape latency is the time, which the D3 treated (open circles) and untreated (filled circles) mice needed to find the

hidden platform in the Morris water maze assay. Error bars indicate standard deviations, asterisks (*) indicate p<0.05. Taken from Funke et al., 2010^{23} .

D3 and the D3-derivative D3D3 additionally slowed down neurodegeneration in TBA2.1 mice Figure 5²⁶. TBA2.1 mice are transgenic mice which exhibit a motor neurodegenerative phenotype induced by the neuronal presence of N-terminal truncated and pyroglutamated A β (3 42) (pEA β (3-42))²⁷. For pEA β (3-42) it was shown that it exhibits an increased cytotoxicity, enhanced oligomerization tendency compared to A β (1-42)^{28,29}, and is able to impair long-term potentiation²⁷.



Figure 5: Assessment of TBA 2.1 phenotype using (A) part of the SHIRPA test battery and (B) motor phenotype using rotarod performance test. Homozygous (HOM) and wild type (WT) TBA2.1 mice were treated intraperitoneally over 4 weeks with phosphate buffered saline (placebo) (n = 7), or with 5 mg per kg body weight D3 (n = 8) or D3D3 (n = 8) per day. (A) SHIRPA test before (black) and after i.p. treatment (white). D3D3-treated mice did not show significant worsening of their phenotypes in the SHIRPA test, whereas the phenotypes of D3-treated mice and saline-treated mice worsened significantly or very significantly (Repeated measures ANOVA before vs. after p = 0.0027, paired t-test (hypothesized difference ≥ 0) before vs. after, *p = 0.035, **p = 0.010.) (B) Rotarod analysis of HOM demonstrates a worsening of the motor phenotype in placebo treated mice whereas D3 and D3D3 administration inhibited this process. Taken from Brener et al.,²⁶.

More recently, we developed derivatives of D3 with improved properties during a lead optimization strategy that focused primarily on the A β oligomer elimination efficiency. We used our newly developed A β -QIAD (quantitative determination of interference with A β aggregate size distribution) to quantitatively measure A β oligomer elimination efficiency and thus in vitro target engagement²⁶.

Our most advanced D3-derivative RD2 is in any regard more efficient than D3. *In vitro* RD2 is indeed eliminating toxic A β oligomers very efficiently Figure 7³⁰ and inhibits the formation of A β fibrils demonstrated in the ThT assay³⁰.



Figure 6: Lewis structure of the D3- derivative RD2.



Figure 7: $A\beta QIAD$ assay $A\beta(1-42)$ size distributions without d-peptide (black in A, B and C) and in the presence of compounds were analysed by density gradient centrifugation. $A\beta(1-42)$ concentrations for each fraction were determined by UV absorption during RP-HPLC. $A\beta(1-42)$ oligomers of interest are located in fractions 4 to 6 (A and B), or fraction 2 (C), which represents the fractions 4, 5 and 6 of the standard QIAD assay. Comparison of 20 μ M D3 (grey) and 20 μ M RD2 (blue) (A and B) revealed a higher oligomer elimination efficacy of RD2. In (C) a QIAD assay for $A\beta(1-42)$ size distributions in dependence of different RD2 concentration (0 μ M black, 1 μ M till 40 μ M, from light to dark blue) is shown. Graphical

representation of the measured decrease in oligomer concentration in % of the amount of oligomers found in the control with $0 \mu M$ RD2 in dependence of the different RD2 concentration. The curve was fitted according to a logistic fit function yielding an IC50 value (D). In each experiment, data is represented as mean \pm SD; A: one-way ANOVA, with Bonferroni post hoc analysis, $A\beta(1-42)$ vs. D3 or RD2 in fractions 4-6, ***p ≤ 0.001 , fraction 4 D3 vs RD2, *p = 0.019, fraction 5 D3 vs RD2, *p = 0.025, fraction 4 D3 vs RD2, *p = 0.045. Data of $A\beta(1-42)$ and D3 were taken from Brener et al. 2015^{26} . Taken from van Groen et al., 2017^{30} .



Figure 8: Aggregation inhibition assay: $A\beta(1-42)$ fibril formation was monitored by Thioflavin T (ThT) in the absence or presence of different RD2 concentrations (0 μ M black, 80 μ M red, 40 μ M light blue, 20 μ M pink, 10 μ M green, 5 μ M blue, 2.5 μ M lilac, 1.25 μ M purple) (A). Fibril mass was normalized to the A β control and the inhibition in % was calculated. The curve was fitted according to a logistic fit function yielding an IC₅₀ value (B). The data represent three replicates. Taken from van Groen et al 2017³⁰.

Additionally, it has been shown to neutralize $A\beta$ -induced neurotoxicity in cell cultures (MTT assay) and to reduce the catalytic effect of preformed seeds on $A\beta$ aggregation thus preventing prion-like propagation of $A\beta$ Figure 9.





Figure 9: *MTT- Cell* viability assay and $A\beta$ Seeding assay: (A,B) RD2 reduced the negative impact of $A\beta(1-42)$ on cell viability in PC-12 and SH-SY5Y cells. Cell viability was assessed by MTT after incubation of PC-12 (A) or SH-SY5Y (B) cells with 1 μ M $A\beta(1-42)$ (light grey) or 1 μ M $A\beta(1-42)$ co-incubated with 5 μ M RD2 (grey), respectively. Data confirms the efficacy of RD2 to significantly increase the cell viability after co-incubation with $A\beta(1-42)$ on both cell lines. Data is represented as mean \pm SD, one-way ANOVA with Bonferroni post hoc analysis, *** $p \le 0.001$. (C,D) RD2 reduced the catalytic effect of preformed seeds on $A\beta$ aggregation The aggregation of monomeric $A\beta$ (m, white) alone or together with preformed seeds (s, grey) incubated with or without RD2 (black) was monitored and the half-life ($t_{1/2}$) (C) as well as the amplitude of the aggregation (D) was determined. In (D) the amplitude of monomeric $A\beta$ without seeds was set to be 100% and the amplitudes of the other samples were normalized accordingly. Data is represented as mean \pm SD, one-way ANOVA with Bonferroni post hoc analysis, *** $p \le 0.001$. Taken from van Groen et al 2017³⁰.

Furthermore, it was shown that the spatial learning and cognition (Morris water maze test) was significantly improved after oral³¹ or intraperitoneal treatment³⁰ with RD2 in two different mouse models. These experiments have been carried out in three different laboratories. These studies also demonstrated that the general or anxiety-related behaviors of the mice were not affected in these studies (open field and zero maze tests – data not shown).



Figure 10: *RD2* treatment effect on the cognitive impairment of APPSL in the Morris water maze (MWM). APP_{SL} mice were treated orally for 6 weeks with either RD2 (n = 15) or 0.9 %

*NaCl as control group (placebo) (n=15). For the validation of the Morris water maze, a group of non-transgenic littermates was used as control (n=15). RD2 treated mice showed significantly improved learning behavior during the acquisition phase of the MWM compared to those of placebo treated mice (two-way RM ANOVA, Fisher post hoc analysis, *p <0.05). The performance RD2 treated mice was nearly the same to those of non-transgenic littermates (ntg), suggesting a reversal of the cognitive impaired phenotype of APP_{SL} mice. Taken from Kutzsche et al., 2017³¹*

Even transgenic mice at the advanced age of 18 months with full-blown pathology that were orally treated with RD2 for 12 weeks demonstrated significantly improved cognitive performance and significantly reduced plaque load in the cortex (unpublished data). Moreover RD2 was able to slow down the progression of the neurodegenerative phenotype of the pyroglutamate-modified A β (pEA β (3-42)) expressing mouse model TBA2.1, which shows hippocampal neuron loss and progressive motor neurodegeneration (unpublished data). This indicates a strong therapeutic potential of RD2 also against pEA β (3-42)-induced neurodegeneration.

Additionally, D3D has very suitable pharmacokinetic properties. In particular, D3D has a long half live in plasma and brain. Three to four hours after application, it reaches a brain-plasma ratio close to one³².



Figure 11: Pharmacokinetic properties of 3H-radioactively labelled RD2 in male C57BL/6 mice were studied using different administration routes. The applied amount contained 10 mg/kg for oral gavage (p.o.) and 3.3 mg/kg for intravenous (i.v.) injection. RD2 shows long terminal half-life in brain, very high bioavailability in the brain, and a high brain-plasma ratio, which demonstrates that RD2 reaches the target organ.

Taken together the advantages of our strategy and of our drug candidate RD2 are manifold and unique: RD2 can be applied orally. It does not recognize a specific A β oligomer species, but stabilizes A β monomers in an aggregation incompetent conformation, thus avoiding any danger of resistance development. We do not at all intend to decrease total A β levels in brain, plasma or CSF. RD2 prevents A β oligomer formation and eliminates oligomers that are already formed. Therefore, RD2 is a promising new drug candidate for the causative treatment of AD.

References

- 1 Alzheimer, A. Über eine eigenartige Erkrankung der Hirnrinde. *Allgemeine Zeitschrift fur Psychiatrie und Psychisch-gerichtliche Medizin.* **64**, 146-148 (1907).
- 2 Cras, P. *et al.* Senile plaque neurites in Alzheimer disease accumulate amyloid precursor protein. *Proceedings of the National Academy of Sciences of the United States of America* **88**, 7552-7556 (1991).
- 3 Armstrong, R. A. The molecular biology of senile plaques and neurofibrillary tangles in Alzheimer's disease. *Folia Neuropathologica* **47**, 289-299 (2009).
- Selkoe, D. J. Alzheimer's disease: genes, proteins, and therapy. *Physiological reviews* 81, 741-766 (2001).
- 5 Brion, J. P. Neurofibrillary tangles and Alzheimer's disease. *European neurology* **40**, 130-140 (1998).
- 6 Selkoe, D. J. The deposition of amyloid proteins in the aging mammalian brain: implications for Alzheimer's disease. *Annals of medicine* **21**, 73-76 (1989).
- 7 Hardy, J. A. & Higgins, G. A. Alzheimer's disease: the amyloid cascade hypothesis. *Science* **256**, 184-185 (1992).
- 8 Goate, A. *et al.* Segregation of a missense mutation in the amyloid precursor protein gene with familial Alzheimer's disease. *Nature* **349**, 704-706, doi:10.1038/349704a0 (1991).
- 9 Clark, R. F. *et al.* The structure of the presenilin 1 (S182) gene and identification of six novel mutations in early onset AD families. *Nature Genetics* 11, 219, doi:10.1038/ng1095-219 (1995).
- 10 Levy-Lahad, E. *et al.* Candidate gene for the chromosome 1 familial Alzheimer's disease locus. *Science* **269**, 973-977, doi:10.1126/science.7638622 (1995).
- 11 Rogaev, E. I. *et al.* Familial Alzheimer's disease in kindreds with missense mutations in a gene on chromosome 1 related to the Alzheimer's disease type 3 gene. *Nature* **376**, 775, doi:10.1038/376775a0 (1995).
- 12 De Strooper, B. *et al.* Deficiency of presenilin-1 inhibits the normal cleavage of amyloid precursor protein. *Nature* **391**, 387-390, doi:10.1038/34910 (1998).
- 13 Castellano, J. M. *et al.* Human apoE isoforms differentially regulate brain amyloidbeta peptide clearance. *Science translational medicine* **3**, 89ra57, doi:10.1126/scitranslmed.3002156 (2011).
- 14 Hashimoto, T. *et al.* Apolipoprotein E, especially apolipoprotein E4, increases the oligomerization of amyloid β peptide. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **32**, 15181-15192, doi:10.1523/JNEUROSCI.1542-12.2012 (2012).
- 15 Josephs, K. A. *et al.* Beta-amyloid burden is not associated with rates of brain atrophy. *Annals of neurology* **63**, 204-212, doi:10.1002/ana.21223 (2008).

- 16 Selkoe, D. J. Soluble Oligomers of the Amyloid β-Protein Impair Synaptic Plasticity and Behavior. *Behavioural brain research* **192**, 106-113, doi:10.1016/j.bbr.2008.02.016 (2008).
- 17 Jin, M. *et al.* Soluble amyloid β-protein dimers isolated from Alzheimer cortex directly induce Tau hyperphosphorylation and neuritic degeneration. *Proceedings of the National Academy of Sciences* **108**, 5819-5824, doi:10.1073/pnas.1017033108 (2011).
- 18 Jucker, M. & Walker, L. C. Pathogenic protein seeding in Alzheimer disease and other neurodegenerative disorders. *Annals of neurology* 70, 532-540, doi:10.1002/ana.22615 (2011).
- 19 Morales, R., Bravo-Alegria, J., Duran-Aniotz, C. & Soto, C. Titration of biologically active amyloid–β seeds in a transgenic mouse model of Alzheimer's disease. *Scientific reports* 5, 9349, doi:10.1038/srep09349
- https://www.nature.com/articles/srep09349#supplementary-information (2015).
- Huang, Y. & Mucke, L. Alzheimer Mechanisms and Therapeutic Strategies. *Cell* **148**, 1204-1222, doi:https://doi.org/10.1016/j.cell.2012.02.040 (2012).
- 21 Wiesehan, K. *et al.* Selection of D-amino-acid peptides that bind to Alzheimer's disease amyloid peptide abeta1-42 by mirror image phage display. *Chembiochem : a European journal of chemical biology* **4**, 748-753, doi:10.1002/cbic.200300631 (2003).
- 22 van Groen, T. *et al.* Reduction of Alzheimer's disease amyloid plaque load in transgenic mice by D3, A D-enantiomeric peptide identified by mirror image phage display. *ChemMedChem* **3**, 1848-1852, doi:10.1002/cmdc.200800273 (2008).
- 23 Aileen Funke, S. *et al.* Oral treatment with the d-enantiomeric peptide D3 improves the pathology and behavior of Alzheimer's Disease transgenic mice. ACS Chem Neurosci 1, 639-648, doi:10.1021/cn100057j (2010).
- 24 van Groen, T., Kadish, I., Funke, A., Bartnik, D. & Willbold, D. Treatment with Abeta42 binding D-amino acid peptides reduce amyloid deposition and inflammation in APP/PS1 double transgenic mice. *Advances in protein chemistry and structural biology* 88, 133-152, doi:10.1016/b978-0-12-398314-5.00005-2 (2012).
- 25 van Groen, T., Kadish, I., Funke, S. A., Bartnik, D. & Willbold, D. Treatment with D3 removes amyloid deposits, reduces inflammation, and improves cognition in aged AbetaPP/PS1 double transgenic mice. *Journal of Alzheimer's disease : JAD* 34, 609-620, doi:10.3233/JAD-121792 (2013).
- 26 Brener, O. *et al.* QIAD assay for quantitating a compound's efficacy in elimination of toxic Abeta oligomers. *Scientific reports* **5**, 13222, doi:10.1038/srep13222 (2015).
- 27 Alexandru, A. *et al.* Selective hippocampal neurodegeneration in transgenic mice expressing small amounts of truncated Abeta is induced by pyroglutamate-Abeta formation. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **31**, 12790-12801, doi:10.1523/jneurosci.1794-11.2011 (2011).
- 28 Gunn, A. P., Masters, C. L. & Cherny, R. A. Pyroglutamate-Abeta: role in the natural history of Alzheimer's disease. *The international journal of biochemistry & cell biology* 42, 1915-1918, doi:10.1016/j.biocel.2010.08.015 (2010).
- 29 Jawhar, S., Wirths, O. & Bayer, T. A. Pyroglutamate amyloid-beta (Abeta): a hatchet man in Alzheimer disease. *The Journal of biological chemistry* 286, 38825-38832, doi:10.1074/jbc.R111.288308 (2011).
- 30 van Groen T, Schemert S., Brener O, Gremer L, Ziehm T, Tusche M, Nagel-Steger L, Kadish I, Schartmann AE, Elfgen A, Jürgens D, Willuweit A, Kutzsche J, Willbold D. The Aβ oligomer eliminating D-enantiomeric peptide RD2 improves cognition

without changing plaque pathology. *Scientific reports*, doi:doi: 10.1038/s41598-017-16565-1 (2017).

- 31 Kutzsche, J. *et al.* Large-Scale Oral Treatment Study with the Four Most Promising D3-Derivatives for the Treatment of Alzheimer's Disease. **22**, doi:10.3390/molecules22101693 (2017).
- 32 Leithold, L. H. *et al.* Pharmacokinetic Properties of a Novel D-Peptide Developed to be Therapeutically Active Against Toxic beta-Amyloid Oligomers. *Pharmaceutical research* **33**, 328-336, doi:10.1007/s11095-015-1791-2 (2016).

F 7 Quantitative approaches to antibiotic resistance

Tobias Bollenbach Institute of Biological Physics University of Cologne

Contents

1	Intro	oduction	2
	1.1	Antibiotic resistance	2
	1.2	Evolution experiments	2
2	Inve	stigating antibiotic resistance using experimental evolution	4
	2.1	Evolution under increasing antibiotic concentration over time	4
	2.2	Evolution in spatial antibiotic concentration profiles	4
3	Рори	ulation genetics models	6
	3.1	Mutation, selection, and genetic drift	6
	3.2	Model for an asexual population of constant size	7
4	Gen	etic interactions of resistance mutations	7
	4.1	Distribution of fitness effects	7
	4.2	Epistasis and fitness landscapes	9
5	Dru	g combinations	11
	5.1	Drug interactions: Synergy and antagonism	12
	5.2	Drug combination that minimize resistance evolution	12
6	Con	clusion	14

Lecture Notes of the 49th IFF Spring School "Physics of Life" (Forschungszentrum Jülich, 2018). All rights reserved.

Parts of these lecture notes are based on a recent review article on the same topic [50].

1 Introduction

1.1 Antibiotic resistance

The high rate at which bacteria evolve resistance to new antibiotics is a serious public health problem. Infectious diseases caused by pathogenic bacteria are becoming increasingly hard to treat due to resistance [12, 61, 2]. To make things worse, the discovery rate of new antibiotics is in decline (despite some notable exceptions [48]). Spontaneous mutations which occur during DNA replication are an important source of resistance. Spontaneous resistance evolution is of fundamental interest since it enables us to observe in the laboratory how genes with diverse functions turn into resistance genes, often by just a few point mutations. A deeper understanding of the dynamics and the underlying mechanisms of resistance evolution is a crucial prerequisite for the long-term goal of predicting rates of resistance evolution and for developing strategies to minimize it. Progress in our understanding of the evolutionary dynamics leading to antibiotic resistance thus holds promise to help avert the looming resistance crisis [47].

The resistance of a bacterium to a drug is determined by measuring the dose-response curve, i.e. the specific growth rate as a function of the drug concentration, and the minimal inhibitory concentration (MIC), i.e. the lowest drug concentration that completely inhibits growth of a clonal population (Figure 1) [11]. An increase in resistance occurs when the population can grow at higher concentrations of antibiotic. Resistance is a genetically inherited trait; bacteria become resistant through one of two main processes: *de novo* mutations that occur during DNA replication and horizontal gene transfer where genes directly move from the genome of one species to another [28, 43]. For simplicity, we focus on *de novo* mutations here.

The investigation of antibiotic resistance has elucidated several distinct molecular mechanisms that enable bacteria to evade antibiotics. Common mechanisms of antibiotic resistance observed in the clinic or laboratory include enzymes that degrade antibiotics, efflux pumps that remove drugs from the cell interior, the prevention of drug uptake, and drug target modification [7, 35, 40, 74]. These mechanisms have been characterized in detail, culminating in databases of the "resistome", i.e. the collection of all genes known to confer antibiotic resistance [53, 76, 77]. It is of fundamental interest to understand the dynamics of resistance evolution in microbial populations, identify general principles that underlie this phenomenon, and to develop a theoretical framework that enables predicting future resistance evolution.

1.2 Evolution experiments

Spontaneous antibiotic resistance evolution can be studied systematically in controlled and reproducible experiments. The extremely short generation time of bacteria (e.g. less than 30 min for *Escherichia coli*) enables us to follow evolutionary adaptation in the laboratory within days [23]; in the presence of antibiotics, particularly fast adaptation with large, easily detectable changes in drug resistance occurs [6, 18, 50, 71]. Due to great advances in DNA sequencing technology, the entire time course of mutations that occur during an evolution experiment (the "evolutionary path") can be readily observed. Recently, biological physicists developed several new experimental evolution techniques and high-throughput methods which provide powerful tools for reconstructing and rationalizing mutational paths to drug resistance. These techniques



Fig. 1: Schematic of a dose-response curve which measures the growth rate (response) of a bacterium as function of the concentration (dose) of an antibiotic. Empirically, most dose-response curves can be well approximated with a Hill function with parameters IC_{50} (drug concentration at which growth is inhibited by 50%) and m (a parameter corresponding to the steepness of the curve). The minimum inhibitory concentration (MIC) is the lowest concentration at which no growth occurs. Figure from [56].

expose bacteria to antibiotic concentrations that increase in time or in space and thus enable following resistance evolution over multiple sequential mutation steps, each of which increases resistance by up to an order of magnitude. Total resistance increases of 100 to 1,000-fold are commonly observed within weeks in the lab [6, 18, 50, 71]. In addition, evolution experiments can be automated using dedicated liquid handling robots and thus be performed in many replicates. This approach enables rigorous statistical investigations of parallel evolution [22]. This is crucial since evolutionary dynamics is inherently stochastic: Even perfect replicates of an evolution experiment in the exact same conditions and starting from the same initial condition can generally have different outcomes as different resistance mutations may occur due to random chance.

One of the most common experimental evolution protocols is serial transfer of a microbial culture [23]. In this protocol, bacterial cultures grow in liquid nutrient medium in flasks (with typical volumes of 10-100 mL) or on microtiter plates which contain e.g. 96 wells (with typical volumes of 200 μ L per well); these cultures are diluted into fresh medium by a fixed dilution factor (e.g. 100-fold) at regular time intervals (e.g. every 24 hours). These experiments can be kept running almost indefinitely: A famous long-term evolution experiment by Richard Lenski [23] has by now exceeded a staggering 60,000 generations in 28 years and is still ongoing. One of the key challenges in such evolution experiments is avoiding contamination: The longer the experiment runs the more likely it becomes that a microbe from the environment invades the culture and, in the worst case, outcompetes the bacterium studied in the experiment. Due to the relative simplicity of the serial transfer protocol, one can run hundreds of evolution experiments in parallel. Together with increasingly inexpensive whole genome sequencing techniques [4] which can provide a comprehensive overview of all mutations that occurred during the experiment, this enables a statistical investigation of the intrinsically stochastic evolutionary dynamics and for identifying general principles governing microbial evolution [22, 44, ?, 73].

In principle, the serial transfer protocol can be used to investigate antibiotic resistance evolution [60, 33, 59]. However, a drawback of serial transfer protocols is their inability to control key

parameters that affect the evolutionary process: The population size changes considerably over each growth-and-dilution cycle; moreover, different cultures usually differ in their growth rates and in the time they spend in stationary phase (the state that is reached when a microbial population has run out of nutrients and stopped growing) between dilutions. These issues complicate the quantitative investigation of the evolutionary process and its comparison among different microbes and conditions. A particular problem with the basic serial transfer protocol is that it is not clear how the antibiotic concentration used in the experiment should be chosen to gain maximum insight into the process of resistance evolution: If it is too low, there is virtually no selection for antibiotic resistance — the bacteria will happily grow as in the absence of drug with no need to evolve drug resistance; if the antibiotic concentration is too high, the bacteria cannot grow at all, which largely prevents them from evolving at any significant rate because, in the absence of DNA replication, there are few mutations that could occur to increase resistance.

2 Investigating antibiotic resistance using experimental evolution

2.1 Evolution under increasing antibiotic concentration over time

Recently developed techniques for evolution experiments circumvent the limitations of choosing a single fixed drug concentration by exposing bacteria to increasing antibiotic concentrations. Before these techniques were developed, theoretical work using models similar to those introduced in section 3 suggested that temporally or spatially increasing drug concentration gradients can facilitate the sequential emergence and fixation of multiple resistance mutations and thus lead to increasingly higher resistance levels [31, 34]. Motivated by these theoretical results, advanced protocols that gradually increase antibiotic concentrations in time or space have been developed [5, 69, 71, 78].

A notable example is the "morbidostat": This feedback-controlled device keeps liquid bacterial cultures of fixed volume growing in exponential phase and automatically increases the antibiotic concentration during the course of the experiment such that they keep growing at a pre-defined rate despite their increasing resistance. As in a simpler chemostat, fresh nutrient medium is provided throughout the experiment; the rate at which this medium is added is also controlled to maintain a constant population size. Both the constant growth rate and the constant population size are continuously monitored via turbidity measurements (the turbidity of a bacterial culture is a well-established proxy for the number of bacteria per volume in the culture); throughout the experiment, a control software decides when fresh growth medium or more antibiotic needs to be added. In this way, strong selection pressure for resistance is constantly maintained and both the selection pressure and the population size are quantitatively controlled — an ideal situation for the comparison to theoretical models of evolution (see section 3). For some antibiotics, this protocol enables the highly reproducible evolution of a \sim 1,000-fold resistance increase in just a few weeks [71].

2.2 Evolution in spatial antibiotic concentration profiles

Spatial antibiotic concentration profiles may enable even faster resistance evolution. An example is given by a recently developed microfluidic device: Here, a concentration gradient of the antibiotic ciprofloxacin (which targets the bacterial DNA replication machinery) was



Fig. 2: Schematic representation and direct visualization of evolving populations on antibiotic landscapes. (a) A Muller plot showing the relative abundance of different mutants in an evolving population. Each colored region represents a subpopulation with a different fitness; darker green indicates higher fitness. (b) A picture of the MEGA-plate showing growth of E. coli over bands of increasing trimethoprim concentration (left to right). Snapshot after 220h of growth [5]. (c) The same evolving population as in a, represented on a discrete fitness landscape of four mutations: the circles represent all possible genotypes of four loci, each with two alleles. Two points are connected by an edge if the two genotypes are one mutational event apart. The arrows show the only path that has monotonically increasing fitness from the least fit to the fittest genotype. Figure from [50].

maintained across a hexagonal device of \sim 2-3 cm diameter; the device consists of over 100 connected micro-compartments, so that bacteria could move between different concentrations. From an initial population size of 10⁶ bacteria, a surprisingly strong over 200-fold increase in resistance occurred; this increase resulted from multiple reproducible mutations that had occurred as early as 10 hours after inoculation [78].

In contrast, a recent study used a very large setup and followed evolution on a huge, meter-scale agar plate (the "MEGA plate") [5]. The size of the plate allows for large bacterial population sizes which should accelerate the occurrence of resistance mutations in the population (see section 3). In contrast to normal small agar plates (diameter ~ 10 cm) where rapid diffusion of the antibiotic in the agar quickly destroys any spatial concentration profiles, they remain relatively stable on this larger plate over the time scale of the experiment (weeks). In this experiment, a front of bacterial growth moves across the plate and resistance evolution to extremely high levels is observed for different antibiotics within weeks. By using concentration gradients with different slopes, it was found that smaller increments in drug resistance over the same distance enable the multi-step evolution of high resistance levels that are practically impossible to reach in a single mutation step [5]. A fascinating aspect of this experiment is that the front of bacteria that grows across the plate can be viewed as a living "Muller diagram" (Figure 2a) that shows the emergence, competition, and fate of different mutants that occur during the experiment. The MEGA plate thus directly visualizes the evolutionary record and the key role of stochastic events in this process (Figure 2b): Some of the most highly resistant lineages ultimately stalled in this assay because they emerged in an unfavorable location too far away from the growth front [5], illustrating the stochasticity of the process. Together, these results highlight

the potential of new assays with well-defined spatial drug landscapes as tools for investigating resistance evolution.

Overall, there are powerful new tools for observing resistance evolution at different levels. Beyond these experimental techniques, theoretical descriptions of evolution are needed to interpret the data and ultimately make predictions for the outcomes of evolution experiments.

3 Population genetics models

Population genetics provides the theoretical framework for describing evolving populations [32]. A typical population genetics model describes a population of N individuals which can have different genotypes. Here, we focus on microbial populations, in particular those in an experimental evolution device without spatial structure like the morbidostat [71]. The population is haploid, i.e. each individual carries only one copy of its chromosome, and asexual, i.e. there is no recombination among the chromosomes of different individuals; each individual simply produces two offspring cells, in most cases with the exact same genotype. The key processes that determine the evolutionary dynamics of such a population are mutation, selection, and genetic drift.

3.1 Mutation, selection, and genetic drift

Mutations are rare errors that occur when the DNA is copied before cell division. Mutation events include point mutations where a single nucleobase is replaced with a different one, deletions in which longer DNA sequences are lost, duplications where a DNA sequence is by accident copied twice, and insertions where certain DNA sequences ("insertion elements") are introduced at a random location in the genome. Mutations are rare: A typical estimate for the mutation rate is 10⁻¹⁰ per base pair and generation, i.e. the cellular DNA replication machinery on average makes its first error only after it has copied 10^{10} base pairs. For a bacterium like E. coli with a genome size of 4.5×10^6 bp, this implies that, on average, it can copy its DNA and divide over 1,000 times before the first mutation occurs. This extremely low error rate well below what equilibrium thermodynamics would allow – is achieved by a mechanism called "kinetic proofreading" [36] in which energy is dissipated to introduce irreversible steps in the DNA replication process and thus increase fidelity. Nevertheless, mutations can occur rapidly in typical microbial populations because their population size is huge: Typical microbial populations consist of $N \sim 10^9$ individuals. This implies that, within a single generation time, $\sim 10^6$ mutations occur in the population. Thus, for any given position in the genome, there will be an individual in the population with a mutation in that position after just a few generations.

Most mutations are neutral or deleterious, i.e. they have no effect on fitness or lead to a decrease in fitness. Here, "fitness" quantifies the reproductive success of an individual (or genotype) in the population; for microbes in evolution experiments, we can often approximate fitness with the growth rate (which can be measured easily). Deleterious mutations are not expected to affect the evolutionary dynamics of large populations because they should die out quickly. More relevant are beneficial mutations, i.e. those that increase fitness in the current environment. Thus, population genetics models often capture only beneficial mutations.

Selection results from the fact that the individuals in the population in general have different fitness. The selection coefficient *s* measures the relative fitness advantage of an individual with a beneficial mutation compared to its ancestor. For example, a mutant with selection coefficient

s = 0.05 produces viable offspring at a rate that is 5% greater than that of its ancestor; for a bacterium in an exponentially growing population that just experienced a beneficial mutation with selection coefficient s = 0.05, the specific growth rate is simply 5% greater than for the rest of the population.

Genetic drift results from random sampling of the surviving individuals, which inevitably needs to happen to maintain a finite population size: In a population of fixed size, half of the individuals must disappear after each individual has produced on average one offspring. Thus, even if an individual has a fitness advantage compared to the others and produces more offspring, its lineage may die out due to genetic drift. Drift plays an important role in evolutionary dynamics since individuals with a new beneficial mutation are initially always present in extremely low numbers and can thus easily die out. Another consequence of drift is that, in the absence of new mutations, after long times the entire population will consist of descendants of a single individual. In large populations, individuals that do not have a fitness advantage compared to others in the population are virtually certain to die out.

3.2 Model for an asexual population of constant size

To illustrate the essence of a populations genetics model, we focus on a simple model of an asexual population of size N that can easily be implemented and solved numerically [27]. For simplicity, we assume discrete non-overlapping generations. In the first step, we randomly select the individuals that will survive to the next generation. Here, the survival probability of each individual is proportional to its fitness and the overall survival probabilities are normalized such that on average only half the population will survive to the next generation. In step two, each surviving individual duplicates. In the last step, each individual can acquire a new beneficial mutation with probability U_b (beneficial mutation rate). When a beneficial mutation has occurred, its fitness effect is drawn from a probability distribution, the Distribution of Fitness Effects (DFE [24], see below). This procedure is then repeated for each generation. Multiple beneficial mutations can be present in different individuals of the population – this effect, called "clonal interference", plays an important role for the evolutionary dynamics of large populations. From such simulations, it is straightforward to determine quantities such as the rate at which the average fitness of the population increases or the fixation rate of beneficial mutations.

The shape of the DFE is usually unknown. It is common to assume that this distribution (or its extreme value distribution) is exponential or has some other defined shape. However, it is an active field of research to measure the DFE (or at least proxies thereof) for different situations. For evolution experiments in the presence of antibiotics in a concentration regime where bacterial growth is strongly inhibited by the antibiotic, recent work has shown that the shape of the DFE of the most common mutations can be approximated by assuming a log-normal distribution of changes in drug resistance resulting from mutations [18]. These resistance changes can then be translated into fitness changes using the inverse function of the dose-response curve of the antibiotic (Figure 3).



Fig. 3: The Distribution of Fitness Effects (DFE) of mutations under antibiotics is strongly affected by the shape of the dose-response curve. Schematics of two different dose response curves: the left curve is steep, i.e. the growth rate is sensitive to small changes in drug concentration; the right curve is shallow. Mutations cause changes in drug resistance which can be interpreted as shifts in the effective drug concentration a bacterium experiences (e.g. a 2-fold increase in resistance leads to a 2-fold lower effective drug concentration); the typical magnitude of these shifts is surprisingly similar for diverse antibiotics [18]. The distribution of effective drug concentrations resulting from mutations is shown in gray. These mutations produce distributions of growth rates (fitness) that are wide for the steep dose-response curve (left) and narrow for the shallow dose-response curve (right). Figure from [50].

4 Genetic interactions of resistance mutations

4.1 Distribution of fitness effects

Key parameters in theoretical descriptions of evolution are mutation rates and selection coefficients. While mutation rates can be measured or at least estimated from experimental data [38], it is challenging to measure the DFE of mutations [24] (which determines the selection coefficients): In principle, the fitness effects of all possible mutations would need to be measured to map this distribution. In practice, it is feasible to quantify the fitness effects of large numbers of mutations which can be generated using different techniques. Recent studies have measured proxies of the DFE in the absence [?] or presence of antibiotics and revealed general quantitative relations that partly explain the shape of this distribution [18, 51, 68, 72].

In particular, the width of the DFE in the presence of antibiotics depends strongly on the shape of the dose-response curve of the drug. The DFEs for 8 antibiotics representing diverse modes of action were approximated by measuring growth rates in the presence of fixed concentrations of the antibiotics for all 4,000 strains of the genome-wide *E. coli* gene deletion library [18, 3]. In each of these strains, one gene is deleted; in this way, a spontaneous loss-of-function mutation in the gene is mimicked. Since loss-of-function is the most likely (non-neutral) outcome of spontaneous mutations, these strains provide a relevant proxy for the primary effect of spontaneous mutations. Interestingly, the widths of the observed DFEs vary drastically across antibiotics. These differences are almost completely explained by the shape of their dose-response curves (Figure 4): When the growth rate of a bacterium is sensitive to small differences in the concentration of a particular antibiotic, the corresponding DFE is wide; conversely, when the



Fig. 4: The width of the DFE in the presence of antibiotics correlates with the steepness of the dose-response curve. Scatterplot of dose-sensitivity n (a measure for the steepness of the dose-response curve) and DFE width (interquartile range, IQR); Pearsons $\rho = 0.96$, $p = 1.3 \times 10^{-4}$; n error bars show the standard deviation of replicate measurements; DFE width error bars show bootstrap 95% confidence interval. Horizontal dashed line shows DFE width in the absence of antibiotics. Gray line shows a linear relation as a guide to the eye. Figure modified from [18].

growth rate depends only weakly on such small dose differences, the corresponding DFE is narrow (Figure 3).

A population genetics model of the type described in section 3.2 predicted that the rate of resistance evolution should increase with increasing DFE width (Figure 5). An additional prediction from this model was that more diverse resistance mutations should be observed for increasing DFE width. These predictions were confirmed in evolution experiments using the morbidostat [18]. These results highlight the potential of identifying general principles and quantitative determinants that shape the DFE and thus affect resistance evolution.

4.2 Epistasis and fitness landscapes

The DFE generally depends on the genetic background, since the same mutations can have different effects on fitness dependent on the presence or absence of other mutations in the genome. The DFE can thus change as soon as the first mutation has occurred. The phenomenon where the effect of mutations depends on the presence of other mutations is called "epistasis" [64]. Measuring the extent of epistasis is important for evolutionary predictions. The reason is that epistasis can lead to multiple fitness peaks, i.e. multiple different genotypes that represent a local maximum of fitness in the sense that any single mutation generally leads to genotypes with lower fitness. This situation can prevent a population from reaching the global fitness maximum [21]; in particular, this is the case for "reciprocal sign epistasis" where the fitness effect of a mutation changes sign (e.g. from positive to negative) depending on the genetic background [65].

Epistasis can be visualized using discrete fitness landscapes [62, 75]. A discrete fitness landscape is a graph where the vertices are genotypes; two genotypes are connected by an edge if they are a single mutational event apart; a fitness value is assigned to all genotypes (Figure 2c). Paths on the landscape are accessible by evolution if they represent sequences of genotypes with monotonically increasing fitness, i.e. all mutations along the path are beneficial. This reflects that it is usually highly unlikely that a mutant with lower fitness can survive long enough in a population to acquire a second mutation. Strictly, however, this is possible for very large popu-





Fig. 5: Theoretical predictions for the rate of resistance evolution for two different antibiotics and experimental validation. (A) Simulation results from a theoretical model of resistance evolution in a morbidostat: IC_{50} increase over time for a drug with narrow distribution of resistance (IC_{50}) changes and two available large-effect mutations that can occur only once (magenta) or none (gray); light lines are sample runs; dark lines are mean of 200 runs; inset: distribution of relative IC_{50} changes used in the simulations. IC_{50} is the drug concentration at which growth is inhibited by 50% (Figure 1). (B) Same as in panel A for a wider distribution of resistance changes. (C, D) Results from morbidostat laboratory evolution experiments: IC_{50} increase over time for the antibiotics nitrofurantoin (C) and chloramphenicol (D) which have different modes of action but similar dose-sensitivity n; light lines are individual runs; dark lines are mean, error bars standard deviation; shaded region in C indicates early phase during which large-effect mutations fix. Figure from [18].

lations or high mutation rates. The assumption that only individual mutation steps can be taken on the landscape corresponds to the "strong selection, weak mutation" (SSWM) limit in which mutations are sufficiently rare that they occur and fix in the population one-by-one.

If only few mutational paths on the landscape are accessible, evolutionary trajectories become more constrained and predictable. The fitness of all mutations and their combinations is experimentally inaccessible due to the astronomically large number of possible genotypes, even if mutations are restricted to short DNA sequences (e.g. a single gene). Therefore, recent studies have focused on full landscape reconstructions of just a few mutations relevant for drug resistance [62, 75, 29] and proposed biophysical models to predict epistatic interactions from protein structure and function [26, 66].

The metabolic enzyme dihydrofolate reductase (DHFR) has served as a key model for higherorder epistasis and the biophysical constraints that shape fitness landscapes. DHFR is the target of the antibiotic trimethoprim. Resistance to this antibiotic can evolve via a few point mutations in the DHFR gene and its promoter. These mutations lead to the overexpression of DHFR which can partly compensate its inhibition by the antibiotic; in addition, they can reduce the affinity of the antibiotic to its target. A recent study identified seven different mutations in DHFR that occurred in replicate evolution experiments; the entire fitness landscape of these mutations was then determined by constructing mutants with all possible pairwise and higher order combinations of these seven mutations and measuring their fitness [62]. This approach revealed that pervasive epistasis may increase the accessibility of all peaks on the fitness landscape, provided that indirect mutational paths (in which the same locus mutates more than once) are allowed [62]. Furthermore, the effects of three resistance mutations in DHFR and their combinations on biophysical parameters of DHFR (enzyme efficiency, stability, and ability to bind trimethoprim) were measured [66]. This work revealed a biophysical trade-off between affinity to trimethoprim, enzyme efficiency, and stability which shapes the epistatic interactions in the fitness landscape. This underlines the importance of genetic interactions across different cellular mechanisms for predicting evolution.

A broader investigation of interactions between drug resistance and other cellular functions can uncover genes that accelerate resistance evolution. Notable examples are mechanisms that increase the mutation rate in response to an antibiotic challenge ("stress-induced mutagenesis"): An increased mutation rate can accelerate resistance evolution because the most beneficial resistance mutations occur faster in the population. In particular, stress-induced mutagenesis can occur by upregulation of the mutation-inducing DNA damage response ("SOS response") [19, 30, 52], induction of mutagenic oxidative damage [42, 54], or by DNA uptake from the environment [16]. Such mechanisms that directly affect the ability of cells to evolve adaptively ("evolvability") are a potential target for new drugs that could be combined with established antibiotics to hamper spontaneous resistance evolution — an idea that has triggered efforts to develop inhibitors of the SOS response [19, 1, 57].

5 Drug combinations

Combining multiple antibiotics is a promising possibility for slowing down resistance evolution [14, 33, 55, 41, 58, 37, 45, 59, 70]. Drug combinations are further used in basic research as a means of perturbing multiple cell functions to reveal relationships in cell physiology [46, 25], analogous to genetic epistasis measurements [20]. At the same time, drug combinations may increase the efficacy of the drugs. While combination treatments are still relatively rare for



Fig. 6: Drug interactions are defined by the shape of lines of equal effect in two-drug concentration space. Schematics showing growth rate (grayscale) and minimal inhibitory concentration (MIC) line (black, line of zero growth) in the two-dimensional concentration space of drugs A and B. The additive reference is given by linear interpolation of the MICs of the individual drugs [49]. For synergistic and antagonistic drug interactions the MIC line lies below or above this additive expectation, respectively. Suppression is a hyper-antagonistic case in which drug A alleviates the effect of drug B. Insets: growth rates in the absence of drugs ('0'), and at fixed concentrations of drugs A and B individually and combined ('A+B'). The dashed horizontal line in insets indicates the additive expectation [8]. Figure from [9].

antibiotics, they are commonly used in the treatment of many conditions and diseases, including tuberculosis and cancer [13, 39]. One of the central questions is which drugs should be combined, and at which concentration ratio, to maximize the success of a treatment while minimizing the emergence of drug resistance and side effects.

5.1 Drug interactions: Synergy and antagonism

The interaction between two drugs is synergistic if the joint effect of the drugs is stronger than an additive expectation and it is antagonistic if it is weaker. These effects can be defined based on the shape of the response surface, i.e. the growth rate g as a function of the two drug concentrations c_A , c_B [49, 8] (Figure 6). Suppression is an extreme kind of antagonism in which one drug alleviates the effect of the other (Figure 6). A few synergistic antibiotic pairs have been applied for decades as they can reduce side effects in treatments and increase the potency of drugs that are ineffective alone [63]. Despite their growing relevance, fundamental questions about drug interactions remain unanswered. In particular, little is known about the underlying mechanisms of most drug interactions [9, 10, 17, 56]. Still, our understanding of how antibiotic combinations affect microbes has advanced considerably in recent years. Specifically, networks of pairwise interactions among large numbers of drugs were quantified and several studies provided insights into the effects of drug combinations on resistance evolution.

5.2 Drug combination that minimize resistance evolution

In some cases, resistance evolution can be slowed down simply because several independent mutations are required to become resistant against a combination of drugs with different cellular targets. As discussed in section 2, microbes offer unique possibilities for studying resistance evolution in well-controlled evolution experiments. Rates of resistance evolution can vary considerably between antibiotics for reasons that are largely unknown [61, 71, 18]. Several studies used evolution experiments to investigate the effects of antibiotic combinations with different



Fig. 7: Rescaling of response surface in resistant mutants. (a-d) Lines of zero growth (green, sensitive strain; red, doxycycline-resistant strain) and other growth rate isoboles (grey scale) in synergistic (a, b) and suppressive (c, d) drug pairs. The shape of the lines of the sensitive and resistant strains is similar; rescaling the doxycycline concentrations maps lines of zero growth almost on top of each other (b, d). In the synergistic case (a), this scaling leaves the growth region of the sensitive strain fully enclosed by that of the resistant strain. In contrast, in the suppressive case, the scaling generates a region in which only the sensitive strain grows (c, asterisk). Figure modified from [14].

drug interactions on spontaneous resistance evolution and revealed a general trend that antagonistic drug combinations lead to slower resistance evolution than synergistic ones [33, 55]. However, recent work on *S. aureus* suggested that this trend may not hold generally when bacteria evolve higher levels of resistance as the drug interactions themselves might change due to resistance mutations [67].

A particularly promising approach for identifying drug combinations that minimize resistance evolution is to use competition experiments between drug-sensitive and resistant strains. The potential of this approach is highlighted by the observation that suppressive drug interactions (Figure 6) can invert selection: A drug-sensitive *E. coli* strain was shown to rapidly outcompete a doxycycline-resistant strain under the suppressive combination of doxycycline and ciprofloxacin. Hence, these suppressive drug combinations can in principle be used to select against drug resistance. This effect occurred for different resistance mechanisms and should hold more generally for suppressive drug interactions [14].

The reason that suppressive interactions can select against resistance is that, in most cases,

resistance mutations lead to a simple, approximately linear rescaling of the response surface (Figure 7): If the response surface of the wild type for two drugs is $g_{WT}(c_A, c_B)$, the response surface of a mutant is

 $g_{\text{mut}}(c_A, c_B) = \lambda_g g_{\text{WT}}(\lambda_A c_A, \lambda_B c_B)$

with suitable scaling factors λ_g , λ_A , $\lambda_B > 0$ [14, 17]. As a result of this rescaling, for suppressive drug pairs, there exists a region in 2d drug concentration space where the sensitive wild type strain grows but a mutant that is resistant to one of the drugs does not (Figure 7c). These observations motivated the development of an innovative screening technique, using neutral labeling of sensitive and resistant bacteria with different fluorescent proteins to identify natural product molecules that can select against drug resistance [15].

6 Conclusion

Antibiotic resistance evolution offers a rare opportunity where evolution can be observed on relatively short time scales and thus quantitatively investigated in well-controlled laboratory settings. New technology for evolution experiments together with improved mapping of mutational fitness effects and epistatic interactions will soon enable the field to statistically test predictions from theoretical models of evolution in various simple and structured environments. A key technical challenge is to scale recently developed, precisely controlled lab evolution protocols [71] to higher throughput so that many antibiotics and various strains can be tested in parallel at high replication. This technology would enable a more systematic approach to exploiting the vast potential of drug combinations for preventing, slowing, circumventing, or even countering resistance. It would further enable a systematic statistical investigation of mutations and other factors that affect resistance evolution. A deeper understanding of the general principles and dynamics of resistance evolution would result. Crucially, this would enable the development of a theoretical framework that can make quantitative predictions of resistance evolution. In the long run, it will be crucial to translate the fundamental insights on resistance evolution from basic research into specific intervention strategies that are effective against pathogenic bacteria in an infected host but, unlike current treatments, keep resistance in check.

References

- M. K. Alam, A. Alhhazmi, J. F. DeCoteau, Y. Luo, and C. R. Geyer. RecA Inhibitors Potentiate Antibiotic Activity and Block Evolution of Antibiotic Resistance. *Cell Chem. Biol.*, 23(3):381–91, mar 2016.
- [2] D. I. Andersson and D. Hughes. Microbiological effects of sublethal levels of antibiotics. *Nat. Rev. Microbiol.*, 12(7):465–78, jul 2014.
- [3] T. Baba, T. Ara, M. Hasegawa, Y. Takai, Y. Okumura, M. Baba, K. A. Datsenko, M. Tomita, B. L. Wanner, and H. Mori. Construction of Escherichia coli K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol. Syst. Biol.*, 2(1):20060008, jan 2006.
- [4] M. Baym, S. Kryazhimskiy, T. D. Lieberman, H. Chung, M. M. Desai, and R. Kishony. Inexpensive multiplexed library preparation for megabase-sized genomes. *PLoS One*, 10(5):e0128036, 2015.
- [5] M. Baym, T. D. Lieberman, E. D. Kelsic, R. Chait, R. Gross, I. Yelin, and R. Kishony. Spatiotemporal microbial evolution on antibiotic landscapes. *Science*, 353(6304):1147– 51, 2016.
- [6] M. Baym, L. K. Stone, and R. Kishony. Multidrug evolutionary strategies to reverse antibiotic resistance. *Science*, 351(6268):aad3292, jan 2016.
- [7] J. M. A. Blair, M. A. Webber, A. J. Baylay, D. O. Ogbolu, and L. J. V. Piddock. Molecular mechanisms of antibiotic resistance. *Nat. Rev. Microbiol.*, 13(1):42–51, jan 2015.
- [8] C. I. Bliss. The toxitcity of poisons applied jointly. Ann. Appl. Biol., 26(3):585–615, aug 1939.
- [9] T. Bollenbach. Antimicrobial interactions: mechanisms and implications for drug discovery and resistance evolution. *Curr. Opin. Microbiol.*, 27:1–9, may 2015.
- [10] T. Bollenbach, S. Quan, R. Chait, and R. Kishony. Nonoptimal microbial response to antibiotics underlies suppressive drug interactions. *Cell*, 139(4):707–18, nov 2009.
- [11] A. Brauner, O. Fridman, O. Gefen, and N. Q. Balaban. Distinguishing between resistance, tolerance and persistence to antibiotic treatment. *Nat. Rev. Microbiol.*, 14(5):320–30, apr 2016.
- [12] K. Bush, P. Courvalin, G. Dantas, J. Davies, B. Eisenstein, P. Huovinen, G. A. Jacoby, R. Kishony, B. N. Kreiswirth, E. Kutter, S. A. Lerner, S. Levy, K. Lewis, O. Lomovskaya, J. H. Miller, S. Mobashery, L. J. V. Piddock, S. Projan, C. M. Thomas, A. Tomasz, P. M. Tulkens, T. R. Walsh, J. D. Watson, J. Witkowski, W. Witte, G. Wright, P. Yeh, and H. I. Zgurskaya. Tackling antibiotic resistance. *Nat. Rev. Microbiol.*, 9(12):894–6, dec 2011.
- [13] J. A. Caminero, G. Sotgiu, A. Zumla, and G. B. Migliori. Best drug treatment for multidrug-resistant and extensively drug-resistant tuberculosis. *Lancet. Infect. Dis.*, 10(9):621–9, sep 2010.
- [14] R. Chait, A. Craney, and R. Kishony. Antibiotic interactions that select against resistance. *Nature*, 446(7136):668–71, apr 2007.
- [15] R. Chait, S. Shrestha, A. K. Shah, J.-B. Michel, and R. Kishony. A Differential Drug Screen for Compounds That Select Against Antibiotic Resistance. *PLoS One*, 5(12):e15179, jan 2010.
- [16] X. Charpentier, P. Polard, and J.-P. Claverys. Induction of competence for genetic transformation by antibiotics: convergent evolution of stress responses in distant bacterial species lacking SOS? *Curr. Opin. Microbiol.*, 15(5):570–6, oct 2012.
- [17] G. Chevereau and T. Bollenbach. Systematic discovery of drug interaction mechanisms. *Mol. Syst. Biol.*, 11(4):807, apr 2015.
- [18] G. Chevereau, M. Dravecká, T. Batur, A. Guvenek, D. H. Ayhan, E. Toprak, and T. Bollenbach. Quantifying the Determinants of Evolutionary Dynamics Leading to Drug Resistance. *PLoS Biol.*, 13(11):e1002299, nov 2015.
- [19] R. T. Cirz, J. K. Chin, D. R. Andes, V. de Crécy-Lagard, W. A. Craig, and F. E. Romesberg. Inhibition of mutation and combating the evolution of antibiotic resistance. *PLoS Biol.*, 3(6):e176, jun 2005.
- [20] M. Costanzo, A. Baryshnikova, J. Bellay, Y. Kim, E. D. Spear, C. S. Sevier, H. Ding, J. L. Y. Koh, K. Toufighi, S. Mostafavi, J. Prinz, R. P. St Onge, B. VanderSluis, T. Makhnevych, F. J. Vizeacoumar, S. Alizadeh, S. Bahr, R. L. Brost, Y. Chen, M. Cokol, R. Deshpande, Z. Li, Z.-Y. Lin, W. Liang, M. Marback, J. Paw, B.-J. San Luis, E. Shuteriqi, A. H. Y. Tong, N. van Dyk, I. M. Wallace, J. A. Whitney, M. T. Weirauch, G. Zhong, H. Zhu, W. A. Houry, M. Brudno, S. Ragibizadeh, B. Papp, C. Pál, F. P. Roth, G. Giaever, C. Nislow, O. G. Troyanskaya, H. Bussey, G. D. Bader, A.-C. Gingras, Q. D. Morris, P. M.

Kim, C. A. Kaiser, C. L. Myers, B. J. Andrews, and C. Boone. The genetic landscape of a cell. *Science*, 327(5964):425–31, jan 2010.

- [21] J. A. G. M. de Visser and J. Krug. Empirical fitness landscapes and the predictability of evolution. *Nat. Rev. Genet.*, 15(7):480–90, jul 2014.
- [22] M. M. Desai. Statistical questions in experimental evolution. J. Stat. Mech. Theory Exp., 2013(01):P01003, jan 2013.
- [23] S. F. Elena and R. E. Lenski. Evolution experiments with microorganisms: the dynamics and genetic bases of adaptation. *Nat. Rev. Genet.*, 4(6):457–69, jun 2003.
- [24] A. Eyre-Walker and P. D. Keightley. The distribution of fitness effects of new mutations. *Nat. Rev. Genet.*, 8(8):610–8, aug 2007.
- [25] S. B. Falconer, T. L. Czarny, and E. D. Brown. Antibiotics as probes of biological complexity. *Nat. Chem. Biol.*, 7(7):415–423, jul 2011.
- [26] M. Figliuzzi, H. Jacquier, A. Schug, O. Tenaillon, and M. Weigt. Coevolutionary Landscape Inference and the Context-Dependence of Mutations in Beta-Lactamase TEM-1. *Mol. Biol. Evol.*, 33(1):268–80, jan 2016.
- [27] C. A. Fogle, J. L. Nagle, and M. M. Desai. Clonal interference, multiple mutations and adaptation in large asexual populations. *Genetics*, 180(4):2163–73, dec 2008.
- [28] L. S. Frost, R. Leplae, A. O. Summers, and A. Toussaint. Mobile genetic elements: the agents of open source evolution. *Nat. Rev. Microbiol.*, 3(9):722–32, sep 2005.
- [29] S. J. Gabryszewski, C. Modchang, L. Musset, T. Chookajorn, and D. A. Fidock. Combinatorial Genetic Modeling of pfcrt-Mediated Drug Resistance Evolution in Plasmodium falciparum. *Mol. Biol. Evol.*, 33(6):1554–70, jun 2016.
- [30] R. S. Galhardo, P. J. Hastings, and S. M. Rosenberg. Mutation as a stress response and the regulation of evolvability. *Crit. Rev. Biochem. Mol. Biol.*, 42(5):399–435, 2007.
- [31] P. Greulich, B. Waclaw, and R. J. Allen. Mutational pathway determines whether drug gradients accelerate evolution of drug-resistant cells. *Phys. Rev. Lett.*, 109(8):088101, aug 2012.
- [32] D. L. Hartl and A. G. Clark. Principles of population genetics. Sinauer Associates, 2007.
- [33] M. Hegreness, N. Shoresh, D. Damian, D. Hartl, and R. Kishony. Accelerated evolution of resistance in multidrug environments. *Proc. Natl. Acad. Sci. U. S. A.*, 105(37):13977–81, sep 2008.
- [34] R. Hermsen, J. B. Deris, and T. Hwa. On the rapidity of antibiotic resistance evolution facilitated by a concentration gradient. *Proc. Natl. Acad. Sci. U. S. A.*, 109(27):10775–80, jul 2012.
- [35] A. H. Holmes, L. S. P. Moore, A. Sundsfjord, M. Steinbakk, S. Regmi, A. Karkey, P. J. Guerin, and L. J. V. Piddock. Understanding the mechanisms and drivers of antimicrobial resistance. *Lancet (London, England)*, 387(10014):176–87, jan 2016.
- [36] J. J. Hopfield. Kinetic proofreading: a new mechanism for reducing errors in biosynthetic processes requiring high specificity. *Proc. Natl. Acad. Sci. U. S. A.*, 71(10):4135–4139, 1974.
- [37] L. Imamovic and M. O. A. Sommer. Use of Collateral Sensitivity Networks to Design Drug Cycling Protocols That Avoid Resistance Development. *Sci. Transl. Med.*, 5(204):204ra132–204ra132, sep 2013.
- [38] J. Jee, A. Rasouly, I. Shamovsky, Y. Akivis, S. R. Steinman, B. Mishra, and E. Nudler. Rates and mechanisms of bacterial mutagenesis from maximum-depth sequencing. *Nature*, 534(7609):693–6, jun 2016.
- [39] C. T. Keith, A. A. Borisy, and B. R. Stockwell. Multicomponent therapeutics for net-

worked systems. Nat. Rev. Drug Discov., 4(1):71-8, jan 2005.

- [40] B. Khameneh, R. Diab, K. Ghazvini, and B. S. Fazly Bazzaz. Breakthroughs in bacterial resistance mechanisms and the potential ways to combat them. *Microb. Pathog.*, 95:32–42, jun 2016.
- [41] S. Kim, T. D. Lieberman, and R. Kishony. Alternating antibiotic treatments constrain evolutionary paths to multidrug resistance. *Proc. Natl. Acad. Sci. U. S. A.*, 111(40):14494– 9, oct 2014.
- [42] M. A. Kohanski, M. A. DePristo, and J. J. Collins. Sublethal antibiotic treatment leads to multidrug resistance via radical-induced mutagenesis. *Mol. Cell*, 37(3):311–20, feb 2010.
- [43] E. V. Koonin, K. S. Makarova, and L. Aravind. Horizontal gene transfer in prokaryotes: quantification and classification. *Annu. Rev. Microbiol.*, 55:709–42, 2001.
- [44] S. Kryazhimskiy, D. P. Rice, E. R. Jerison, and M. M. Desai. Global epistasis makes adaptation predictable despite sequence-level stochasticity. *Science*, 344(6191):1519–1522, jun 2014.
- [45] V. Lazar, G. Pal Singh, R. Spohn, I. Nagy, B. Horvath, M. Hrtyan, R. Busa-Fekete, B. Bogos, O. Mehi, B. Csorgo, G. Posfai, G. Fekete, B. Szappanos, B. Kegl, B. Papp, and C. Pal. Bacterial evolution of antibiotic hypersensitivity. *Mol. Syst. Biol.*, 9(1):700–700, apr 2014.
- [46] J. Lehár, G. R. Zimmermann, A. S. Krueger, R. A. Molnar, J. T. Ledell, A. M. Heilbut, G. F. Short, L. C. Giusti, G. P. Nolan, O. A. Magid, M. S. Lee, A. A. Borisy, B. R. Stockwell, and C. T. Keith. Chemical combination effects predict connectivity in biological systems. *Mol. Syst. Biol.*, 3(80):80, 2007.
- [47] S. B. Levy and B. Marshall. Antibacterial resistance worldwide: causes, challenges and responses. *Nat. Med.*, 10(12s):S122–S129, dec 2004.
- [48] L. L. Ling, T. Schneider, A. J. Peoples, A. L. Spoering, I. Engels, B. P. Conlon, A. Mueller, D. E. Hughes, S. Epstein, M. Jones, L. Lazarides, V. a. Steadman, D. R. Cohen, C. R. Felix, K. A. Fetterman, W. P. Millett, A. G. Nitti, A. M. Zullo, C. Chen, and K. Lewis. A new antibiotic kills pathogens without detectable resistance. *Nature*, 517(7535):455–459, jan 2015.
- [49] S. Loewe. Die quantitativen Probleme der Pharmakologie. *Ergebnisse der Physiol.*, 27(1):47–187, dec 1928.
- [50] M. Lukačišinová and T. Bollenbach. Toward a quantitative understanding of antibiotic resistance evolution. *Curr. Opin. Biotechnol.*, 46:90–97, mar 2017.
- [51] R. C. MacLean and A. Buckling. The distribution of fitness effects of beneficial mutations in Pseudomonas aeruginosa. *PLoS Genet.*, 5(3):e1000406, mar 2009.
- [52] I. Matic, F. Taddei, and M. Radman. Survival versus maintenance of genetic stability: a conflict of priorities during stress. *Res. Microbiol.*, 155(5):337–41, jun 2004.
- [53] A. G. McArthur, N. Waglechner, F. Nizam, A. Yan, M. A. Azad, A. J. Baylay, K. Bhullar, M. J. Canova, G. De Pascale, L. Ejim, L. Kalan, A. M. King, K. Koteva, M. Morar, M. R. Mulvey, J. S. O'Brien, A. C. Pawlowski, L. J. V. Piddock, P. Spanogiannopoulos, A. D. Sutherland, I. Tang, P. L. Taylor, M. Thaker, W. Wang, M. Yan, T. Yu, and G. D. Wright. The comprehensive antibiotic resistance database. *Antimicrob. Agents Chemother*, 57(7):3348–57, jul 2013.
- [54] O. Méhi, B. Bogos, B. Csörg, F. Pál, A. Nyerges, B. Papp, and C. Pál. Perturbation of iron homeostasis promotes the evolution of antibiotic resistance. *Mol. Biol. Evol.*, 31(10):2793–804, oct 2014.
- [55] J.-B. Michel, P. J. Yeh, R. Chait, R. C. Moellering, and R. Kishony. Drug interactions modulate the potential for evolution of resistance. *Proc. Natl. Acad. Sci. U. S. A.*,

105(39):14918-14923, sep 2008.

- [56] K. Mitosch and T. Bollenbach. Bacterial responses to antibiotics and their combinations. *Environ. Microbiol. Rep.*, 6(6):545–57, dec 2014.
- [57] C. Y. Mo, S. A. Manning, M. Roggiani, M. J. Culyba, A. N. Samuels, P. D. Sniegowski, M. Goulian, and R. M. Kohli. Systematically Altering Bacterial SOS Activity under Stress Reveals Therapeutic Strategies for Potentiating Antibiotics. *mSphere*, 1(4):1–15, 2016.
- [58] C. Munck, H. K. Gumpert, A. I. N. Wallin, H. H. Wang, and M. O. A. Sommer. Prediction of resistance development against drug combinations by collateral responses to component drugs. *Sci. Transl. Med.*, 6(262):262ra156–262ra156, nov 2014.
- [59] T. Oz, A. Guvenek, S. Yildiz, E. Karaboga, Y. T. Tamer, N. Mumcuyan, V. B. Ozan, G. H. Senturk, M. Cokol, P. Yeh, and E. Toprak. Strength of selection pressure is an important parameter contributing to the complexity of antibiotic resistance evolution. *Mol. Biol. Evol.*, 31(9):2387–401, sep 2014.
- [60] C. Pál, B. Papp, and V. Lázár. Collateral sensitivity of antibiotic-resistant microbes. *Trends Microbiol.*, 23(7):401–407, 2015.
- [61] A. C. Palmer and R. Kishony. Understanding, predicting and manipulating the genotypic evolution of antibiotic resistance. *Nat. Rev. Genet.*, 14(4):243–8, apr 2013.
- [62] A. C. Palmer, E. Toprak, M. Baym, S. Kim, A. Veres, S. Bershtein, and R. Kishony. Delayed commitment to evolutionary fate in antibiotic resistance fitness landscapes. *Nat. Commun.*, 6(May 2014):7385, jun 2015.
- [63] S. K. Pillai, R. C. Moellering Jr., and G. M. Eliopoulos. Antimicrobial Combinations. In V. Lorian, editor, *Antibiot. Lab. Med.*, volume 5, pages 365–440. Lippincott Williams and Wilkins, Philadelphia, 2005.
- [64] F. J. Poelwijk, V. Krishna, and R. Ranganathan. The Context-Dependence of Mutations: A Linkage of Formalisms. *PLoS Comput. Biol.*, 12(6):e1004771, 2016.
- [65] F. J. Poelwijk, S. Tanase-Nicola, D. J. Kiviet, and S. J. Tans. Reciprocal sign epistasis is a necessary condition for multi-peaked fitness landscapes. *J. Theor. Biol.*, 272(1):141–4, mar 2011.
- [66] J. V. Rodrigues, S. Bershtein, A. Li, E. R. Lozovsky, D. L. Hartl, and E. I. Shakhnovich. Biophysical principles predict fitness landscapes of drug resistance. *Proc. Natl. Acad. Sci.* U. S. A., 113(11):E1470–8, mar 2016.
- [67] M. Rodriguez de Evgrafov, H. Gumpert, C. Munck, T. T. Thomsen, and M. O. A. Sommer. Collateral Resistance and Sensitivity Modulate Evolution of High-Level Resistance to Drug Combination Treatment in Staphylococcus aureus. *Mol. Biol. Evol.*, 32(5):1175– 85, may 2015.
- [68] A. Sousa, S. Magalhães, and I. Gordo. Cost of antibiotic resistance and the geometry of adaptation. *Mol. Biol. Evol.*, 29(5):1417–28, may 2012.
- [69] F. Spagnolo, C. Rinaldi, D. R. Sajorda, and D. E. Dykhuizen. Evolution of Resistance to Continuously Increasing Streptomycin Concentrations in Populations of Escherichia coli. *Antimicrob. Agents Chemother.*, 60(3):1336–42, dec 2015.
- [70] S. Suzuki, T. Horinouchi, and C. Furusawa. Prediction of antibiotic resistance by gene expression profiles. *Nat. Commun.*, 5:5792, dec 2014.
- [71] E. Toprak, A. Veres, J.-B. Michel, R. Chait, D. L. Hartl, and R. Kishony. Evolutionary paths to antibiotic resistance under dynamically sustained drug selection. *Nat. Genet.*, 44(1):101–105, jan 2012.
- [72] S. Trindade, A. Sousa, and I. Gordo. Antibiotic resistance and stress in the light of Fisher's model. *Evolution*, 66(12):3815–24, dec 2012.

- [73] S. Venkataram, B. Dunn, Y. Li, A. Agarwala, J. Chang, E. R. Ebel, K. Geiler-Samerotte, L. Hérissant, J. R. Blundell, S. F. Levy, D. S. Fisher, G. Sherlock, and D. A. Petrov. Development of a Comprehensive Genotype-to-Fitness Map of Adaptation-Driving Mutations in Yeast. *Cell*, 166(6):1585–1596.e22, sep 2016.
- [74] C. Walsh. Antibiotics. American Society of Microbiology, jan 2003.
- [75] D. M. Weinreich, N. F. Delaney, M. A. Depristo, and D. L. Hartl. Darwinian evolution can follow only very few mutational paths to fitter proteins. *Science*, 312(5770):111–4, apr 2006.
- [76] J. D. Winkler, A. L. Halweg-Edwards, K. E. Erickson, A. Choudhury, G. Pines, and R. T. Gill. The Resistome: A Comprehensive Database of Escherichia coli Resistance Pheno-types. ACS Synth. Biol., 5(12):1566–1577, 2016.
- [77] G. D. Wright. The antibiotic resistome: the nexus of chemical and genetic diversity. *Nat. Rev. Microbiol.*, 5(3):175–86, mar 2007.
- [78] Q. Zhang, G. Lambert, D. Liao, H. Kim, K. Robin, C.-k. Tung, N. Pourmand, and R. H. Austin. Acceleration of emergence of bacterial antibiotic resistance in connected microenvironments. *Science*, 333(6050):1764–7, sep 2011.

Index

A

 α - β -transition, C3.13 accumulation, E3.6 actin, C3.5, C5.5, D6.5 actin filament, C3.2, D7.3, F3.5, F3.12 actin polymerization, E4.9 action potential, B4.9, C7.4, D3.9, D9.3, D9.9, D9.13, F4.3 active fluid, F5.7 active gel, C5.9, C6.2, C6.3, C6.8 active matter, E3.2, E4.2, I.2 active Ornstein-Uhlenbeck particle (AOUP), E3.5 active particle, composite, E4.11 active pressure, E3.7 active stress, C6.2, C6.5, C6.10 active Brownian particle (ABP), E3.4 adenylate kinase, B2.8 adherens junction, D7.4 adhesion bond model, E2.9 agent-based model, F5.8 aggregation inhibition assay, F6.8 Airy pattern, A1.2 algae, E4.5 allele, F2.14 alternating-access mechanism, B4.4 Alzheimer's disease, B7.2, F6.2 amide I absorption, B2.3 AMPA receptor, C7.7 amyloid, B7.2, F6.3 amyloid cascade, F6.3 amyloid oligomers, B7.3 amyloid- β , B7.3 animal cells, D6.3 anion channel, F4.5 anisotropic filter, A8.7 anomalous dispersion, A5.9 anomalous subdiffusion, D1.4 anomaly exponent, D1.4, D1.11 antagonism, F7.12 antibiotic resistance, F2.7, F2.15, F7.2 antibody fragment motion, B2.11 antigen, F3.3 antigen presenting cell, F3.3

antiporter, B4.3 apoferritin, B3.3 apparent viscosity, E2.7 arrestin, F1.7 arhythmogenic right ventricular cardiomyopathy, D7.11 artificial biological devices, I.3 artificial biological organisms, I.3 artificial cell, B9.17 Asakura Osawa Vrii, D2.3 astrocyte, C7.3, C7.10, C7.12 ataxia, F4.5 atomic force microscopy, E1.19, E5.12 atomic vibration, B2.3 ATP synthase, D4.10 ATPase, B4.4 ATPase motor, B2.14 axon, D3.2 axon hillock, D3.9 axoneme, E4.6

B

backbone dynamics, B2.8 bacteria, E4.4, E5.2, F7.2 bacteria, swimming near surfaces, E4.20 base-pair, B5.2 Bayes, A9.6 Becker-type myotonia, F4.2 Beer-Lambert's law, F1.3 bending energy, C1.4, F5.9 bending potential, A10.4 bending rigidity, F5.9 Berg-Purcell limit, D5.9, F1.10 bilateral filter, A8.7 binding equilibria, force dependence, E5.2 binding force, E5.3 binodal, B3.21 bioelectronics, D9.15, D9.17, I.4 biofilm, E5.2 biofluidmechanics. I.2 biofluids, I.2 biofunctionalization, D4.1, D4.6 biological function, I.2 biomacromolecule, I.2

biomechanics, L4 biomimetic microswimmers, E4.7 bistable gene circuit, B9.2 bivalent ions, D4.8, D4.10 blood, E2.2 blood flow, E2.3 blood flow in microcirculation, E2.11 blood flow resistance, E2.7 Boltzmann distribution, F1.6 Boltzmann factor, A10.9 bond potential, A10.4 bovine serum albumin, B3.3 Bragg reflection, A5.8, A9.7 brain, E7.2 Brownian dynamics (BD), A10.11 Brownian motion, D1.2, F1.2, A3.2 bystander cells, F3.5

С

calcium dynamics, F3.8 calorimetry, B1.12 cancer. E6.2 Canham-Helfrich energy, C1.5 catalytic effect, F6.9 catastrophe, C3.7 cell crawling, E4.2 cell mechanics, C3.2, E6.2 cell migration, E4.2, E4.10 cell motility, E4.2, E4.10 cell shape, E5.5, E5.8 cell sorting, E5.7 cell viability, F6.9 cell-cell interaction, E5.3-E5.5 cell-electrode cleft, D9.3, D9.5, D9.7 cell-free gene expression, B9.16 cell-surface interaction, E5.3, E5.4 centrosome, F3.7 cerebral cortex, E7.3 channel dysfunction, F4.2 channel rhodopsin, B8.2, B8.4 chemical shift, B2.5 chemical synapse, C7.6 chemoattractant, F1.9 chemoreception, F1.2 Chevron-plot, B1.7 Chlamydomonas reinhardtii, E4.5 chloride pump, B8.9

chromatin, B5.2 chromosome, B5.2 cilia, E4.5 cilia, synchronization, E4.22 circle-swimming bacteria, E4.20 circuit-chassis interaction, B9.9 classification, A6.4 ClC-K channel, F4.4 clonal interference, F2.9, F7.7 cluster, E3.14 cluster size distribution, B3.15 cluster-fluid phase, B3.16 CNG channels, F1.5 cognitive impairment, F6.5 cognitive performance, F6.10 coherent scattering, A2.9 collagen, D6.5 collective behavior, E3.9 collective diffusion, B3.9 colloid-osmotic model, F5.10 colloidal dispersion, A10.5 colloidal system, E5.5 colocalization analysis, A8.13 communication, E7.9 compartmentalization, B9.10, B9.15, B9.17, D4.3 competition, E5.9, E6.7, F7.13 complex systems, I.3 complexity, I.3 composite active particle, E4.11 computational electrophysiology, B4.11 computational neuroscience, E7.2 confocal microscopy, E1.16, A1.13 conformational entropy, B1.5, B1.12 connected component analysis, A8.7 connectivity, E7.2 constitutive relation, E1.14 contrast transfer function (CTF), A6.4 convolution, A8.5 convolutional neuronal network, A8.11 cooperative rotation, B2.14 correlation, E7.7 corresponding states, A10.7 cotransporter, B4.3 Coulomb potential, A10.5 **CRISPR** interference, B9.6 crowding, B2.14, B3.2, D1.2, D1.10, D2.2

cryo-electron microscopy, A9.3 cryo-EM, A9.3, A9.4 crystal structure, B8.6 crystallography, B2.4 CUSUM test, D5.10 cytoadherence, F5.3, F5.11 cytoskeletal filament , D7.3 cytoskeleton, C3.2, C5.2, E4.8, F3.7, F3.11 cytoskeleton-membrane interaction, E4.10

D

D-enantiomeric peptide, F6.5 D3-derivative, F6.6 deep learning, A8.11 dementia, F6.2 dendrite. D3.2 dendritic cells, F3.3 density-based spatial clustering, A8.14 depletion, D2.3 depletion force, E5.8 depolarization, D3.3 desmosome, D7.4 developmental biology, I.3 differential adhesion, E5.6 differential cross-section, A2.3 diffraction, A1.3 diffusion, F3.5, F3.11 diffusion coefficient, D1.3, A10.7 diffusion function, B3.10 diffusion limited reaction, D2.13 diffusion potential, B4.8 dihydrofolate reductase, F7.11 dimeric motor, C5.10 disease, F4.2, F5.2, F5.14, F6.2 dissipative particle dynamics (DPD), A10.12 distance transform, A8.10 distribution of fitness effects, F7.7 DNA, B5.2, E1.19 DNA origami, B9.11 domain motion, B2.9 dose-response curve, F7.3, F7.8 double differential cross-section, A2.7 droplet-supported GUV, D4.3, D4.4 Drosophila, C5.2 drug combinations, F7.11 drug development, F6.4 drug interactions, F7.12

ductile, E1.11 dynamic frequency sweeps, E1.7 dynamic scattering, A2.7 dynamic strain sweep, E1.11 dynamic structure factor, B3.9, A2.9 dynein, C3.4, F3.7

Е

E. coli, E3.2, E4.4, F7.2 elastic normal modes, B2.10 elastic scattering, A2.3 electrical synapse, C7.5 electro-chemical equilibrium, D3.4 electrochemical gradient, B4.3 electrode-3D, D9.5, D9.8, D9.10 electron beam lithography, D9.6 electron microscope, A6.3 electroporation, C1.3 elevator transporter, B4.4, B4.6 emergent behavior, I.3 endoplasmatic reticulum, C1.2, C1.14, F3.9 endothelial cells, D6.4 energy consumption, E4.2, E4.8 energy landscape, B1.5 energy transfer, A2.7 entropy of protein and hydration water, B1.13 environmental signals, I.4 enzyme cascade, F1.5 enzyme clustering, B9.15 epidermis, D7.8 epidermolysis bullosa simplex, D7.9 epigenetics, B5.13 epilepsy, F4.5 epistasis, F2.13, F7.9 epithelial cells, D6.5 epithelium, D7.6 equilibrium thermodynamics, B1.5 Ermak-McCammon equation, D1.6 Escherichia coli, F7.2 eukaryotes, B5.2 evolution, E5.8, F2.2, F7.2 evolutionary biology, I.3 evolutionary trajectory, F7.11 Ewald sphere, A5.8 exchanger, B4.3 excitatory amino acid transporter, F4.5 excluded volume, D2.2

experiment, E7.2 extracellular gate, B8.6 extracellular matrix, D6.3, E5.7, E5.9 extracellular signal, D9.3, D9.12

F

F-actin, E1.9 Förster Resonance Energy Transfer (FRET), A3.2, A3.23 Fahraeus effect, E2.7 Fahraeus-Lindqvist effect, E2.7, F5.12 fibrillar aggregate, F6.3 fibrils, B7.2 fibroblasts, D6.3 field-effect transistor (FET), D9.2, D9.10 filopodium, D7.4 Fisher's fundamental theorem, F2.4 fitness, E5.8, F2.3 fitness graph, F2.14 fitness landscape, F2.14, F7.9 fixation, F2.5 flagella, E4.4, E5.4 flagellar beat, 2nd harmonic, E4.18 flexible electronics, D9.14, D9.17 flock of birds, E3.2 fluctuation dissipation relations, A10.14 fluctuation spectrum, C1.5 fluid-fluid phase coexistence, D2.9 fluidification, E6.10 fluorescence, A1.5, A3.6, D4.6- D4.12 fluorescent correlation spectroscopy (FCS), D1.8, E1.19, A3.2, A3.7 fluorescent recovery after photobleaching (FRAP), A3.2, A3.16, D1.9 focal adhesion, D7.4 focused ion-beam (FIB), D9.7, D9.16 Fokker-Planck equation, E3.16, A10.13 folding funnel, B1.6 folding intermediate, B1.3 folding kinetics, B1.7 folding mechanobiology, D6.3 force dipole, E3.11 force sensing, D5.2 force spectroscopy, E5.12 forward flux sampling, B7.10 free energy barrier, B1.5 free volume theory, D2.4

friction anisotropy, E4.16 friction force, A10.10 FRM II reactor, A5.7 functional amyloid, B7.5

G

GABA, C7.8 GABA receptor, C7.9 GABA transporter, F4.6 Galilean invariance, A10.13 gap junction, C7.5, C7.13 gating, B4.9 Gauss-Bonnet theorem, C1.5, C1.11 Gaussian blur, A8.6 Gaussian curvature, C1.4, C1.16 gene expression, B9.3 gene regulation, B9.5 generalized hydrodynamics, C6.2 generalized Stokes-Einstein relations, B3.12 generalized transport coefficient, D1.4 genetic drift, F2.5, F7.7 genetic factor, F4.2 genetic logic gate, B9.12 genetic oscillator, B9.2 genomics, I.3 GFP protein, B8.5 giant unilamellar vesicles (GUV), D4.3 glass transition, A3.22 gliding, F5.3, F5.6 globular proteins, B3.3 GltPh, B4.6, B4.12 glutamate, C7.7 glutamate receptor, C7.7 glutamate transporter, B4.6, B4.12, C7.15, F4.5 gold nanoparticles, D4.5, D4.16 gold-linked surfactants, D4.5-D4.7 Goldman equation, D3.5 Goldman-Hodgkin-Katz equation, D3.5 graphene, D9.10 graphene FETs (GFETs), D9.10 graphene MEAs (GMEAs), D9.12 green fluorescent protein (GFP), D1.9, D1.10 growth, mechanics of, E6.2 GYG motif, B4.5

H

hard sphere crystal, D2.8

hard sphere fluid, D2.6 hard-sphere colloids, E5.5 Helfrich energy, C1.5 hematocrit, E2.3 hemidesmosome, D7.4 hemoglobin, F5.5, F5.11 high-resolution structure, B8.7 hinge, B2.8, B2.10 Hinshelwood distribution, F1.6 histone, B5.2 homeostatic pressure, E6.4 hybrid modeling, A9.2, A9.7 hydrodynamic flow, F5.12 hydrodynamic interactions, B3.4, D1.7, A10.10 hydrodynamic interactions, far-field approximation, E4.15 hydrodynamic radius, B1.9 hydrodynamic synchronization, E4.22 hydrodynamics, E3.10, A10.1 hydrodynamics interactions, squirmers, E4.16 hydrodynamics, dipole near wall, E4.27 hydrodynamics, dipole tensor, E4.27 hydrodynamics, Oseen tensor, E4.26 hydrophobic collape, B1.9 hyperpolarization, D3.3

I

image features, A8.7 image segmentation, A8.7 image stitching, A8.12 immune repertoire, F3.3 immune system, F3.2 immunological synapse, F3.2, F3.3, F3.11, F3.12 in-vivo neural probes, D9.17 incoherent scattering, A2.8 induced-fit binding, B4.7 inelastic scattering, A2.7 integrin, D4.10- D4.12, D4.16 intensity threshold, A8.7 intercalated disc, D7.11 intermediate filament, C3.2, D7.3 intermediate scattering function, A2.9 intermittent search, F3.11 internal dynamics (proteins), A2.18 intestinal epithelium, D7.7 intracellular diffusion, D1.9- D1.12

intrinsic shape, B5.6 intrinsically disordered protein, B2.12 ion channel, B4.3, D3.3, D6.8 ion pump, B8.4 ion selectivity, B8.6 isomerization, F1.4 isomorphous replacement, A5.9

J

junction, D7.4

K

K+ buffering, C7.14 KcsA, B4.5 Kern-Frenkel potential, B3.18 kinesin, C3.4, C5.3 kinetic descriptions, C6.2

L

lamellipodium, C6.5, D7.4 Langevin equation, D1.5, E3.4, E4.16, F5.12, A10.11, A10.13 large scale dynamics, B2.13 laser tweezers, E5.12 lattice Boltzmann (LB), A10.11 lattice gas automata (LGA), A10.11 leaky integrate-and-fire, E7.4 Lennard Jones (LJ) potential, A10.4 lifecycle, F5.2 light sensor, B8.2 light-driven pump, B8.3, B8.8 light-energy transducer, B8.2 line tension, C1.9 lipid, C1.2, C1.9, A10.5 lipid bilayer, C1.3 lipid phases, C1.9 liposomes, D4.7, D4.10 liquid-crystalline order, E5.5 living matter, E4.2, I.2 localization, A1.4 locomotion, E3.13, E4.2 loss modules, D6.3 low-resolution data, A9.2 low-resolution structure determination, A2.17 low-Reynolds-number hydrodynamics, E4.25 LTP, long term potentiation, C7.9 lysozyme, B3.3 lytic granule, F3.2, F3.7, F3.12

Lévy-walk, F3.5

Μ

macaque, E7.10 macromolecule, I.2 major histocompatibility complex, F3.2 malaria, F5.2 margination, E2.8 Markov state model, B7.11 mass attenuation coefficients, A6.4 material property, I.2 materials, I.2 maximum entropy, A9.6 maximum likelihood, A6.4 Maxwell-Boltzmann distribution, A10.14 mean curvature, C1.4, C1.16 mean first passage time, F3.12 mean square displacement (MSD), A3.2, B3.6, D1.3, E3.5 mean-field theory, E7.8 mechanical evolution, B5.10 mechanical genome, B5.8 mechanosensing, D7.5 mechanotransduction, D5.2 MEGA-plate, F7.5 membrane-mediated interactions, C1.12 membrane (axisymmetric), C1.8 membrane (cylindrical), C1.6, C1.14 membrane (spherical), C1.6, C1.14 membrane model, E2.4 membrane potential, C7.4, D3.3 membrane shape equation, C1.8, C1.13 memory consumption, E7.8 metabolic engineering, B9.15 metachronal wave, E4.5, E4.22 metachronal wave, transport properties, E4.23 metadynamics, B7.10 micelle, C1.2 micro-organisms, E4.2 microbot, E4.2 microcircuit, E7.3 microcirculation, E2.3 microelectrode array (MEA), D9.2, D9.5 microfluidics, A10.8, D4.3, D4.4, D4.7, D4.8, D4.11. D4.15 microglia, C7.10 microswimmer, E3.2, E4.2, F5.4

microswimmer, biomimetic, E4.7 microswimmer, chiral, E4.18 microswimmer, force dipole, E4.14 microswimmer, puller, E4.14 microswimmer, pusher, E4.14 microswimmer, synthetic, E4.7 microtubule, C3.2, C5.6, D7.3, E3.3, F3.7, F3.12 microtubule organizing center, F3.7, F3.12 migration, D7.2 Mikado model, C5.2 mitochondria, C1.14, F3.7, F3.9 modeling, E7.2 modularity, B9.7, B9.10 Moiré pattern, A1.9 molecular clock, F2.8 molecular dynamics, A10.1, B4.11, B7.9 molecular mechanics energy, A9.2 molecular motor, C5.2 molecular replacement, A5.10 Monge gauge, C1.5, C1.13 monogenetic disease, F4.2 Monte Carlo, A10.1, A10.8 Monte-Carlo sampling, A9.7 Moran model, F2.3 morbidostat, F7.4 morphogenesis, I.3 morphological operations, A8.10 motility, E4.2, E5.3, E5.4 motility-induced phase separation (MIPS), E3.8 Muller diagram, F7.5 Muller's ratchet, F2.11 multiparticle collision dynamics (MPC), A10.1, A10.14, B3.16 multiplexing, B5.2 multiscale modeling, I.4 muscle chloride channel, F4.2, F4.4 mutation, E5.9, F4.4, F7.6 Mutation Monte Carlo (MMC), B5.5 mutation rate, F7.8 mutational load, F2.11 myelin, D3.10 myofibroblasts, D6.3 myoglobin, A5.12 myoglobin as model system, B1.10 myosin, C3.4 myotonia congenita, F4.2

Ν

nanocavities, D9.6 nanostructured droplets, D4.7, D4.12 narrow escape problem, F3.12 natural killer cells, F3.2 natural selection, F2.2 Navier slip length, E4.21 Navier-Stokes equation, E4.13 negative stain, A6.2 Nernst equation, B4.8, D3.4 NEST. E7.9 neurodegeneration, F6.2 neurofibrillary tangle, F6.2 neurometabolic coupling, C7.15 neuropathology, F6.2 neuroprosthetics, D9.17, D9.20 neuroscience, D9.2 neurotransmitter. D3.2 neurovascular coupling, C7.15 neutral evolution. F2.6 neutron cross-section, A6.4 neutron imaging, A6.2 neutron protein crystallography, A5.5 neutron spectroscopy, B1.11 neutron spin echo spectroscopy, B2.9, A2.18 NMDA receptor, C7.7 NMR-2D, B2.6 noise, F1.2 non-affine deformation, E1.12 non-equilibrium systems, I.2 non-local mean filter, A8.7 normal mode deformations, B2.11 nuclear magnetic resonance, B2.5 nucleosome, B5.2 nucleosome positioning rules, B5.8 Nyquist criterion, A1.4

0

oligodendrocyte, C7.12 oligomer, F6.4 Onsager transition, C5.7 Opalina, E4.5 optical microscopy, A1.2 optogenetic instrument, B8.6 optogenetics, B8.2, I.4 optogenetics tools, B8.3 Orai channel, F3.9 Oseen tensor (hydrodynamics), A10.10, E4.26 osmotic pressure, C1.3 Otsu method, A8.7 out-of-equilibrium state, E4.8 out-of-equilibrium system, E4.2

P

Péclet number, C5.8, E1.5 pair correlation function, A2.5, A10.6, B3.8 parallel evolution, F2.7 Paramecium, E4.5 particle tracking, A8.12 partition function, A10.9 patchy interactions, B3.17 pathwalker approach, A9.4 PDMS, D4.5 periodic selection, F2.9 perisynaptic process, C7.13 perivascular process, C7.13 persistence length, B5.3, C3.10 petascale, E7.9 pharmacokinetic property, F6.10 phase behavior, B3.14, D2.6 phase transition, E3.10 pheromone, F1.9 phosphoglycerate kinase, B2.10 phosphorescence, A3.7 photo-bleaching, A3.7 photoactivated localization microscopy (PALM), A1.7 photoreception, F1.2 pico-injection, D4.5 piezo channels, D5.6 pilus, E5.3- E5.5 plasma membrane, C1.2 platelet, E2.3 point contact model, D9.4 point spread function (PSF), A1.2 Poisson distribution, F1.10 Poiseuille flow, A10.8 polyelectrolyte, C3.9 polymer, A10.5 polymerization force, E4.9 population genetics, F7.6 poroelasticity, C6.8 postsynapse, C7.4 potassium conductance, F4.5

predictive neuroinformatics, E7.10 presynaptic terminal, C7.4 primary hemostasis, E2.4 principal curvature, C1.4 printed electronics, D9.15, D9.17 projection matching, A6.4 protease ClpP, B2.7 protein aggregation, B7.3 protein conformation, A9.3 protein crystallography, A2.14 protein data bank, A5.3 protein dynamics, B1.11 protein reconstitutions, D4.1, D4.10 proteoliposomes, D4.10, D4.12 protocell, B9.17 proton pump, B8.2 puller (microswimmer), E4.14 pump, B4.3 pumping activity, B8.5 pusher (microswimmer), E4.14

Q

quantum yield, F1.4 quorum sensing, B9.13

R

R-factor, A5.10 radiography, A6.2 random walk, F3.5, F3.11 rare event, B7.10 rate constant, D2.12 reaction diffusion model, F3.10 receptor potential, D3.7 reciprocal sign epistasis, F7.9 red blood cell, C1.7, C1.8, E1.16, E2.2, F5.2, F5.8, F5.10 reduced units, A10.7 reduced volume, F5.9 reflection interference contrast microscopy, F5.7 self-diffusion, D1.2, D1.4 refractory phase, D3.9 region labeling, A8.7 replica exchange molecular dynamics, B7.9 replicator equation, F2.4 repolarization, D3.9 resistive-force theory, E4.16 resolution, A1.2, A5.11 response surface, F7.14 resting potential, B4.8

retrograde flow, F5.6 reversible adhesion, E2.11 Reynolds number, E3.10, E4.13 RGD , D4.12, D4.16 rheology, I.2 rheology of tissues, E6.10 rhodopsin, B8.2 rhodopsin kinase, F1.7 riboregulator, B9.6 rigid base-pair model, B5.5 rod-like colloids, A10.5 rolling adhesion, F5.12 Rotne-Prager-Yamakawa (RPY) tensor, D1.7 rouleaux, E1.16 run-and-tumble motion, E4.4

S

Salmonella, E4.3 SALR solutions, B3.12 saltatory conduction, D3.10 sample preparation, A1.12 scalability, E7.8 scaled particle theory, D2.5 scallop theorem, E3.10, E4.14 scanning, A1.13 scanning electron microscopy, D9.7, D9.16 scattering, A2.2 scattering function, A2.9 scattering length, A2.3, A2.5 scattering vector, A2.3 screened Coulomb potential, B3.4 seal resistance, D9.6, D9.12 search strategy, F3.11 selection. F7.6 selection coefficient, F2.3, F7.8 selectivity filter, B4.5 self propulsion, E3.2 self-propelled particle, E4.2, F5.8 self-propelled rods, C5.7 sender-receiver circuit, B9.2 sequence space, B5.5 serial transfer, F7.3 sex, F2.11 shear flow, A10.8 shear modulus, F5.9 shear thickening, E1.14

shear thinning, E1.14 shear-activated polymer, E2.6 SHIRPA test. F6.6 sidechain movement, B2.4 single molecule localization microscopy, A1.6 single-particle tracking (SPT), D1.8 sliding friction, F5.7 slip length, hydrodynamic, E4.21 small unilamellar vesicles, D4.1, D4.8 small-angle neutron scattering (SANS), A2.16, E1.16 small-angle x-ray scattering (SAXS), A2.16 smoothed dissipative particle dynamics (SDPD), suppression, F7.12 E2.13 smoothed particle hydrodynamics (SPH), A10.11surface accumulation, E3.6 solution structure determination, A2.16 sparse ensemble selection, A9.6 spectral precision distance microscopy (SPDM), A1.6 spectrin, F5.5 sperm, E4.4 sperm navigation, E4.18 sperm steering, E4.18 sperm, chiral, E4.18 sperm, swimming velocity, E4.18 spiking, E7.5 spinodal, B3.21 spontaneous curvature, C1.4 squirmer, E3.12 squirmer interactions, E4.16 static scattering, A2.3 stick-slip boundary conditions, A10.7 stiffness matrix, B5.6 STIM protein, F3.9 stimulated emission depletion microscopy (STED), A1.8 stimulation, I.4 stochastic optical reconstruction microscopy (STORM), A1.7 stochastic simulations, C6.2 Stokes equation, E3.10, E4.13 strain, D6.4 strain amplitude, D6.5, E1.4 strain hardening, E1.11 strain softening, E1.11 stress, D6.4, E1.6 stress fibers, C3.2

stress-induced mutagenesis, F7.11 stress-strain dependence, C3.13 strong selection, F2.5 structural biology, A9.2 structural model, B8.9 structure factor, A2.5 structure prediction, A9.3 structured illumination microscopy (SIM), A1.9 substitution rate, F2.8 superpixels, A8.9 superresolution microscopy, A1.2, A8.14 supported lipid bilayer, D4.8, D4.9 supramolecular activation cluster, F3.4 surface tension, C1.3 surfactants, D4.4-D4.6, D4.12-D4.14 swarming, E3.2 swimming, helical, E4.18 symporter, B4.3 synapse, D3.6, E7.3 synaptic potential, D3.7 synchronization, hydrodynamic, E4.22 synchrotron radiation, A5.5 synergism, F7.12 synthetic biology, I.3, D4.3, D4.13 synthetic cells, D4.3, D4.15, D4.16

Т

T cell, F3.2 T cell motility, F3.5 T cell receptor, F3.3 talin, unfolding of, D5.6 tension sensing, D5.2 tethered membrane channels, D5.4 tethered sperm, E4.19 tetrameric motor, C5.10 thermal fluctuations, C1.5, A10.10, A10.12 thermal ratchet, E4.9 thermal stability, F1.6 thermodynamic perturbation theory, B3.18 thin shell elasticity, F5.9 threshold detection, D5.10 tisssue growth, I.3, E6.1 tissue competition, E6.7 tissue spheroids, E6.12 tobacco mosaic virus (TMV), A6.2

tomography, A6.2 torsion potential, A10.4 traction force microscopy, F5.7 trainable segmentation, A8.10 transition path sampling, B7.10 transition state, B1.5 transition-state controlled reaction, D2.12 transporter, B4.2 treadmilling, C3.7, E4.9 tripartite synapse, C7.13 triple point, D2.9 tubular membrane, F4.3 tubulin, C3.2, C3.4 two-state force sensor, D5.9

U

uncertainty, A9.4, A9.5 unfolding, D6.8 unfolding transition, B1.10 uniporter, B4.3 UV absorption, F6.7

V

van Hove correlation function, A2.10 vesicle, C7.7 Vicsek model, E3.9 viscoelastic, D6.3, D6.5, E1.4 viscosity, B3.11 visual cortex, E7.10 Volvox, E4.7 von Willebrand factor (VWF), E2.4 VWF adhesion, E2.9 VWF model, E2.6 VWF-platelet aggregate, E2.11

W

water-in-oil droplets , D4.3, D4.5, D4.7, D4.8 watershed algorithm, A8.8 WCA potential, A10.5 weak selection, F2.5 white blood cell, E2.3 Wiener-Khintchine theorem, A2.2 worm-like chain, C3.11 Wright-Fisher model, F2.3

Х

X-ray crystallography, A9.4 X-ray diffraction, A9.2 x-ray scattering, B3.7 xenorhodopsin, B8.7

Y

Young's modulus, D9.17

The IFF Spring School in Jülich, the Peter Grünberg Institute, and JARA-FIT

PGI/JCNS-TA, 52425 Jülich, GermanyPhone:++49 2461 61-4750Email:springschool@fz-juelich.deweb:www.iff-springschool.de

The annual IFF Spring School is a long-standing tradition of the Institut für Festkörperforschung (IFF) which was founded in 1969. The institute's research topics ranged from electronic and structural properties of solids and nanoelectronics, to the thermal and dynamical behaviour of soft matter. The IFF has organized the Spring School for over 40 years. Since the restructuring in 2011, research in the area of electronic systems, their phenomena, as well as their applications in information technology, became part of the Peter Grünberg Institute (PGI) named after the IFF scientist who received the Nobel Prize in Physics in 2007. Biophysics and soft matter research is now located at the Institute of Complex Systems (ICS). These institutes are linked together and supported by the Institute for Advanced Simulation (IAS), which focuses on developing and applying high-performance computing to understand complex systems, and the Jülich Centre for Neutron Science (JCNS), which is dedicated to the operation of neutron scattering instruments at national and international neutron sources. The IFF Spring School is now organized in turns by PGI and ICS.

The PGI consists of several departments: Quantum Theory of Materials, Theoretical Nanoelectronics, Functional Nanostructures at Surfaces, Scattering Methods, Microstructure Research, Electronic Properties, Electronic Materials, Bioelectronics, Semiconductor Nanoelectronics, and the JARA Institutes for Green Information Technology and for Quantum Information. We operate the Helmholtz Nanoelectronics Facility (HNF) and, together with the Central Facility for Electron Microscopy at the RWTH Aachen University, the Ernst Ruska-Centre (ER-C) for Microscopy and Spectroscopy with Electrons. In addition, our departments participate in the operation of synchrotron and neutron beam lines as well as the Jülich supercomputers. We are part of the Jülich Aachen Research Alliance within the section Fundamentals of Future Information Technology (JARA-FIT) in which we collaborate with the physicists, chemists, electrical engineers, material scientists, and biologists of the RWTH Aachen University. In a concerted way, we conduct exploratory research in nanoelectronics and quantum phenomena with an emphasis on potential long-term applications in information technology and beyond.

The Institute of Complex Systems (ICS) consists of following departments: neutron scattering, theoretical soft matter and biophysics, soft matter, cellular biophysics, molecular biophysics, structural bio-chemistry, biomechanics, and bioelectronics. A major objective of biophysics research is to understand processes far from equilibrium, which distinguishes dead from living matter, and thus to elucidate the structure and function of biological matter and living systems. Examples include biomolecules, such as DNA and various proteins, cells with their complex machinery and functions, tissues which represent a collective organization of cells, and systems biology which concerns intricate inter-actions between different biological entities starting from single molecules to cells and tissues, right up to organs and whole organisms.

Physics of Life

Lecture Notes 49th IFF Spring School **2018** 26 February – 09 March 2018 G. Gompper, J. Dhont, J. Elgeti, C. Fahlke, D. Fedosov, S. Förster, P. Lettinga, A. Offenhäusser ISBN: 978-3-95806-286-3

Topological Matter - Topological

Insulators, Skyrmions and Majoranas Lecture Notes 48th IFF Spring School **2017** 27 March – 07 April 2017 S. Blügel, Y. Mokrousov, T. Schäpers, Y. Ando ISBN: 978-3-95806-202-3

Memristive Phenomena - From Fundamental Physics to Neuromorphic Computing

Lecture Notes 47th IFF Spring School **2016** 22 February - 04 March 2016 R. Waser and M. Wuttig ISBN: 978-3-95806-091-3

Functional Soft Matter

Lecture Notes 46th IFF Spring School **2015** 23 February - 06 March 2015 J. Dhont, G. Gompper, G. Meier, D. Richter, G. Vliegenthart and R. Zorn ISBN: 978-3-89336-999-7

Computing Solids - Models, ab-initio methods and supercomputing

Lecture Notes 45th IFF Spring School **2014** 10 - 21 March 2014 S. Blügel, N. Helbig, V. Meden and D. Wortmann ISBN: 978-3-89336-912-6

Quantum Information Processing

Lecture Notes 44th IFF Spring School **2013** 25 February - 08 March 2013 D. DiVincenzo ISBN: 978-3-89336-833-4

Scattering Methods for Condensed Matter Research: Towards Novel Applications at Future Sources

Lecture Notes 43rd IFF Spring School **2012** 05 - 16 March 2012 M. Angst, T. Brückel, D. Richter and R. Zorn ISBN: 978-3-89336-759-7

Macromolecular Systems in Soft- and Living-Matter

Lecture Notes 42nd IFF Spring School **2011** 14 – 25 February 2011 J. K.G. Dhont, G. Gompper, P. R. Lang, D. Richter, M. Ripoll, D. Willbold and R. Zorn ISBN: 978-3-89336-688-0

Electronic Oxides - Correlation Phenomena, Exotic Phases and Novel

Functionalities Lecture Notes 41st IFF Spring School **2010** 08 – 19 March 2010 S. Blügel, T. Brückel, R. Waser and C.M. Schneider ISBN: 978-3-89336-609-5

Spintronics – From GMR to Quantum Information

Lecture Notes 40th IFF Spring School **2009** 09 – 20 March 2009 S. Blügel, D. Bürgler, M. Morgenstern, C. M. Schneider and R. Waser ISBN: 978-3-89336-559-3

Soft Matter - From Synthetic to Biological Materials

Lecture Notes 39th IFF Spring School **2008** 03 – 14 March 2008 J.K.G. Dhont, G. Gompper, G. Nägele, D. Richter and R.G. Winkler ISBN: 978-3-89336-517-3

Probing the Nanoworld - Microscopies, Scattering and Spectroscopies of the Solid State

Lecture Notes 38th IFF Spring School **2007** 12 - 23 March 2007 K. Urban, C. M. Schneider, T. Brückel, S. Blügel, K. Tillmann, W. Schweika, M. Lentzen and L. Baumgarten ISBN: 978-3-89336-462-6

Computational Condensed Matter Physics

Lecture Notes 37th IFF Spring School **2006** 06 - 17 March 2006 S. Blügel, G. Gompper, E. Koch, H. Müller-Krumbhaar, R. Spatschek and R. G. Winkler ISBN: 978-3-89336-430-5

Magnetism goes Nano - Electron Correlations, Spin Transport, Molecular Magnetism

Lecture Notes 36th IFF Spring School **2005** 14 - 25 February 2005 S. Blügel, T. Brückel and C. M. Schneider ISBN: 3-89336-381-5

Physics meets Biology - From Soft Matter of Cell Biology

Lecture Notes 35th IFF Spring School **2004** 22 March - 02 April 2004 G. Gompper, U. B. Kaupp, J. K. G. Dhont, D. Richter and R. G. Winkler ISBN: 3-89336-348-3

Fundamentals of Nanoelectronics

Lecture Notes 34th IFF Spring School **2003** 10 - 21 March 2003 S. Blügel, M. Luysberg, K. Urban and R. Waser ISBN: 3-89336-319-X

Band / Volume 10 **Soft Matter - Complex Materials on Mesoscopic Scales** Lecture Notes 33rd IFF Spring School **2002** 04 - 15 March 2002 J.K.G. Dhont, G. Gompper and D. Richter ISBN: 3-89336-297-5

Neue Materialien für die Informationstechnik

Manuskripte des 32. IFF-Ferienkurses **2001** 05. - 16. März 2001 R. Waser (Editor) ISBN: 3-89336-279-7

Schriften des Forschungszentrums Jülich Reihe Schlüsseltechnologien / Key Technologies

Band / Volume 147 Neutron Scattering

Lectures of the JCNS Laboratory Course held at Forschungszentrum Jülich and at the Heinz-Maier-Leibnitz Zentrum Garching edited by T. Brückel, S. Förster, G. Roth, and R. Zorn (Eds.) (2017), ca 400 pp ISBN: 978-3-95806-243-6

Band / Volume 148 Neutron scattering

Experimental Manuals of the JCNS Laboratory Course held at Forschungszentrum Jülich and at the Heinz-Maier-Leibnitz Zentrum Garching edited by T. Brückel, S. Förster, G. Roth, and R. Zorn (Eds.) (2017), ca 200 pp ISBN: 978-3-95806-244-3

Band / Volume 149 **Kinetic and thermodynamic considerations on the formation of heteromolecular layers on metal surfaces** C. Henneke (2017), vii, 157, XIV pp ISBN: 978-3-95806-245-0

Band / Volume 150
Spectroscopic characterization of local valence change processes in resistively switching complex oxides
C. Bäumer (2017), x, 206 pp
ISBN: 978-3-95806-246-7

Band / Volume 151 **Magnetic structure in relation to the magnetic field induced ferroelectricity in Y-type hexaferrite Ba** _{2-x}Sr_xZn₂Fe₁₂O₂₂ P. Thakuria (2017), 17, 180 pp ISBN: 978-3-95806-250-4

Band / Volume 152 **Statistical analysis tools for assessing the functional relevance of higher-order correlations in massively parallel spike trains** V. Rostami (2017), x, 176 pp ISBN: 978-3-95806-251-1

Band / Volume 153 **The influence of the substrate on the structure and electronic properties of carbon-based 2D materials** J. Sforzini (2017), XIII, 145 pp ISBN: 978-3-95806-255-9 Band / Volume 154 Gate-All-Around Silicon Nanowire Tunnel FETs for Low Power Applications G. V. Luong (2017), ii, 136 pp ISBN: 978-3-95806-259-7

Band / Volume 155 Graphene Devices for Extracellular Measurements D. Kireev (2017), ix, 169 pp ISBN: 978-3-95806-265-8

Band / Volume 156 Nanoscale 3D structures towards improved cell-chip coupling on microelectrode arrays S. D. Weidlich (2017), II, 154 pp ISBN: 978-3-95806-278-8

Band / Volume 157 Interface phenomena in La_{1/3}Sr_{2/3}FeO₃ / La_{2/3}Sr_{1/3}MnO₃ heterostructures and a quest for p-electron magnetism M. Waschk (2017), ix, 205 pp ISBN: 978-3-95806-281-8

Band / Volume 158 **Physics of Life** Lecture Notes of the 49th IFF Spring School 2018 26 February – 09 March 2018, Jülich, Germany ed. by G. Gompper, J. Dhont, J. Elgeti, C. Fahlke, D. Fedosov, S. Förster, P. Lettinga, A. Offenhäusser (2018), ca 1000 pp ISBN: 978-3-95806-286-3

Weitere *Schriften des Verlags im Forschungszentrum Jülich* unter <u>http://wwwzb1.fz-juelich.de/verlagextern1/index.asp</u>
