

Weiterentwicklung und  
Ausbau numerischer  
Strukturen in den  
AC<sup>2</sup>-Programmen  
ATHLET und COCOSYS

**Weiterentwicklung und  
Ausbau numerischer  
Strukturen in den  
AC<sup>2</sup>-Programmen  
ATHLET und COCOSYS**

Tim Steinhoff  
Volker Jacht

August 2020

**Anmerkung:**

Das diesem Bericht zugrunde liegen-  
de Forschungsvorhaben wurde mit  
Mitteln des Bundesministeriums für  
Wirtschaft und Energie (BMWi) unter  
dem Kennzeichen RS1558 durchge-  
führt.

Die Verantwortung für den Inhalt die-  
ser Veröffentlichung liegt beim Auf-  
tragnehmer.

Der Bericht gibt die Auffassung und  
Meinung des Auftragnehmers wieder  
und muss nicht mit der Meinung des  
Auftraggebers übereinstimmen.

## **Deskriptoren**

AC<sup>2</sup>, Differentialgleichungen, MPI, Numerical Toolkit, NuT, Numerik, PETSc, Softwarearchitektur

## Kurzfassung

Thermohydraulische Vorgänge werden im GRS-Programmpaket AC<sup>2</sup> auf Basis von klassischen Erhaltungsgleichungen und einem Finite-Volumen-Ansatz mittels der Zeitintegration des resultierenden Anfangswertproblems simuliert. Diese Probleme sind in der Regel als *steif* zu betrachten, was seitens der Numerik eine implizite Behandlung notwendig macht. Hiermit einher geht das Erfordernis, lineare Gleichungssysteme zu lösen, in welche die Jacobimatrix der rechten Seite des Anfangswertproblems eingeht. Die Dimension der Probleme ist im Bereich von  $\mathcal{O}(10^3)$  bis gelegentlich  $\mathcal{O}(10^5)$  Unbekannte angesiedelt. Die Jacobimatrix ist stets dünnbesetzt.

In diesem Vorhaben wurden numerisch und software-technisch geprägte Verbesserungen in AC<sup>2</sup> eingebracht. Der Fokus lag hierbei auf der Komponente NuT (Numerical Toolkit) und dessen Interaktion mit den AC<sup>2</sup>-Komponenten ATHLET / ATHLET-CD (ATHLET/CD) sowie COCOSYS.

Das Numerical Toolkit bietet Zugriff auf dedizierte Numerik-Bibliotheken wie PETSc und MUMPS, wodurch skalierbare und moderne Algorithmen zur Lösung dünnbesetzter linearerer Gleichungssysteme zur Verfügung stehen. Die Grundform von NuT wurde bereits im Vorgängerprojekt RS1530 etabliert. Im Rahmen der Projektarbeiten wurde NuT nun erstmalig in der Version 1.0 und als Teil von AC<sup>2</sup> 2019 auch für den Endnutzer einsetzbar gemacht.

Der funktionale Umfang von NuT ist im Vorhaben in verschiedene Richtungen ausgebaut worden. Zum einen wurde die bereits vorhandenen Möglichkeit, die lineare Algebra der ATHLET/CD-Transientenrechnung in NuT anzugehen, mittels weiterer Löserangebote verfeinert und durch NuT-Unterstützungen in der ATHLET-Startrechnung ergänzt. So mit steht jetzt auch in der Startrechnung ein standardisierter Löser bereit, welcher die Dünnbesetzung der zugrundeliegenden Systeme berücksichtigt. Des Weiteren ist es im Restart-Szenario nun möglich, mittels NuT Jacobimatrizen an Restart-Punkten zu speichern. Folglich entfällt die sonst erzwungene Neuberechnung der Jacobimatrix, wenn ein Restart durchgeführt wird. Ebenso besteht mittels PETSc-Skripten die Option, die Matrix z. B. in Python zu laden, um weiterführende Untersuchungen wie eine Spektralanalyse durchzuführen.

Eine wesentliche Neuerung ergab sich durch eine NuT-interne Modernisierung, indem der Code komplett von FORTRAN auf C/C++ umgestellt wurde. Hierdurch ist das Toolkit robuster, übersichtlicher und leistungsstärker geworden. Auf Basis dieser 2.0-Version

von NuT wurden Arbeiten hinsichtlich eines allgemeinen AC<sup>2</sup>-Kommunikationsmodells durchgeführt. Sehr nützlich war hierbei der Einsatz der Kommunikations-Bibliothek MMA. Sowohl NuT als auch ATHLET/CD sind bereits auf das Kommunikationsmodell ausgerichtet worden. Bezuglich COCOSYS sind grundlegende Arbeiten in diese Richtung unternommen worden. Eine komplette Integration steht jedoch noch aus.

Im Thermohydraulik-Modul von COCOSYS wurden zusätzlich funktionale Anpassungen am Code vorgenommen, um ebenfalls auf NuT-Unterstützung zugreifen zu können. Ein derartiger Zugriff sollte möglich sein, sobald das Kommunikationsmodell entsprechend erweitert werden wird. Auch hier ist der Zweck die Auslagerung der linearen Algebra ins Numerical Toolkit.

Der linearen Algebra ist die numerische Differentialgleichungslogik übergeordnet, welche Steuer- und Methodenlogik umfasst. Die Standardmethode in AC<sup>2</sup> ist zurzeit ein Extrapolationsverfahren basierend auf dem expliziten und linear-impliziten Euler-Verfahren. Um hier eine Alternative mit verbesserten Eigenschaften zu bieten, wurden theoretische Arbeiten durchgeführt, die viel Potential im Sinne konkreter Methodenimplementierungen liefern. Arbeiten diesbezüglich können direkt an dem bisher Geleisteten angeknüpft werden.

Im Rahmen der allgemeinen Validierungsarbeiten wurde ein Kohärenz-Test für die Jacobimatrix erstellt, welcher das aktuelle Systemverhalten mit der Netzwerkinformation vergleicht, welche zum Bau der Jacobimatrix genutzt wird. Hiermit konnten bereits einige Diskrepanzen behoben werden. Auch für zukünftige Arbeiten sollte sich der Test als wertvolle Ergänzung zum *Continuous Integration* zeigen.

Das Projekt führt konsequent die Maxime des Vorgängervorhabens fort, sinnvolle und benötigte Numerik innerhalb des AC<sup>2</sup>-Kontextes zentral zu kapseln und den Anwendungen über einfach zu handhabende Schnittstellen zur Verfügung zu stellen. Die geleisteten Arbeiten tragen zur Vereinheitlichung bei oder bereiten hierauf gut vor. Dies bildet eine hervorragende Basis für zukünftige Anpassungen oder Erweiterungen der ACII-Numerik.

## Abstract

The simulation of thermohydraulic processes in GRS's software suite AC<sup>2</sup> is based on classical conservation laws for mass, energy, and momentum of a two-phase flow. For the sake of actual numerical computations spatial discretization by some specific finite volume approach is applied yielding an initial value problem (IVP) composed of a system of ordinary differential equations. Generally, this IVP comes with the attribute of being *stiff*, which makes it necessary to utilize implicit methods in order to ensure an efficient and robust numerical treatment. Implicit methods require solving linear systems that involve the Jacobian of the underlying IVP. Due to the network structure of the problem the Jacobian is sparse. The problem size is usually in the range of  $\mathcal{O}(10^3)$  unknowns but may occasionally go up to  $\mathcal{O}(10^5)$  degrees of freedom as well.

In this project several enhancements and new features related to numerics and software engineering were introduced into AC<sup>2</sup>. The focus of the project was on AC<sup>2</sup>'s component NuT (Numerical Toolkit) and its interaction with other AC<sup>2</sup> components like ATHLET / ATHLET-CD (ATHLET/CD) as well as COCOSYS.

The Numerical Toolkit provides easy access to dedicated numerical libraries like PETSc and MUMPS. This allows for the use of modern and scalable algorithms to handle the linear systems mentioned above in an efficient way. A basic version of NuT was already achieved in the predecessor project RS1530. As one of the main results of the work in this project, NuT 1.0 is part of AC<sup>2</sup>'s release version 2019 and, hence, is available to all AC<sup>2</sup> users.

Several work packages of this project focused on functional upgrades of the Numerical Toolkit within the overall software architecture of AC<sup>2</sup>:

New and refined solver presets were implemented to tackle the aforementioned linear systems. This provides the user with an extended variety of easy to use presets to choose from. As a suitable addendum the steady state calculation of ATHLET was also improved by making use of the chosen preset. A recommendation for when to apply a certain preset is given in NuT's user's manual, which was established as part of the release-related work.

Additional enhancements related to the Jacobian were made in the context of restart calculations. While NuT is activated and making use of appropriate PETSc techniques to save matrix data to disk, it is no longer necessary to calculate a new Jacobian every time a simulation run is initiated at a restart point. This comes with the further benefit

to have Jacobian data available to load into, e.g., Python-tools in order to investigate stability-related attributes of the Jacobian like its spectrum.

A significant improvement towards NuT's software foundation was established via a complete overhaul of its code basis. Except from some necessary interface routines, no FORTRAN code is included in NuT anymore – all functionality is provided by C/C++ code. This results in a more reliable, robust, and powerful code. Based on this 2.0 version of the Numerical Toolkit, a general communication model for AC<sup>2</sup> was developed. This model introduces the open-source communication library MMA to the AC<sup>2</sup>-context. Its main purpose is to support the establishment of communication between heterogeneous processes in an elegant and easy-to-use way. So far ATHLET/CD and NuT have been made compatible to this communication model. Regarding COCOSYS basic steps were implemented, but full compatibility is yet to be achieved.

COCOSYS's thermohydraulic software module comes with the same requirements as ATHLET/CD to solve linear systems during the time integration of a simulation run. Several modifications of the module were made to provide the necessary data to involve the Numerical Toolkit in the solving process. As soon as COCOSYS is fully integrated in the aforementioned communication model, accessing NuT on the part of COCOSYS shall be possible without much additional work.

Solving linear systems arises as a sub-problem during numerical integration by means of an implicit method. AC<sup>2</sup>'s current default method is an extrapolation ansatz based on a combination of the explicit and linearly-implicit Euler's method. To provide an alternative approach with more amiable features as the default several theoretical results were established. These can directly be used to construct concrete methods to be implemented in the Numerical Toolkit.

As part of the general validation process for the modifications listed above, a so-called coherency test for the Jacobian was created. By means of this test it is possible to check whether the given network dependencies and the current system behavior are coherent in terms of the Jacobian's non-zero pattern. The test already helped to eliminate some discrepancies and is a valuable tool for future developments as well as for any reasonable *Continuous Integration* concept in the context of AC<sup>2</sup>.

This project has resolutely continued the effort of its predecessor to centralize required and useful numerics within AC<sup>2</sup> and to ensure easy access to it. The results achieved in this project directly contribute to this maxim, and give a stable platform for further work in this direction.

# Inhaltsverzeichnis

<b>Kurzfassung</b> .....	I
<b>Abstract</b> .....	III
<b>1 Einleitung</b> .....	1
1.1 Anwendungsszenarien von NuT .....	2
1.2 Weitere allgemeine Anmerkungen .....	3
<b>2 AP1 – Ausbau und Verbesserung der Parallelisierungs- und Funktionsmächtigkeit innerhalb des Numerical Toolkit</b> .....	5
2.1 Untersuchungen zur Einbringung von OpenMP in NuT .....	5
2.1.1 Test und Resultate .....	6
2.1.2 Schlussfolgerung .....	10
2.2 Weitere Optionen zur Präkonditionierung iterativer Verfahren .....	10
2.3 Umschreiben des Numerical Toolkit-Codes von FORTRAN nach C++ ...	12
<b>3 AP2 – Verbesserung und Ausbau der Zusammenarbeit von ATHLET mit dem Numerical Toolkit</b> .....	17
3.1 Aufbau eines Doppelpuffers zur effizienteren Nutzung des ATHLET/NuT-Kommunikationskanals .....	17
3.2 Entwicklung einer effizienten Dimensionierungsstrategie zu den linearen Gleichungssystemen .....	19
3.2.1 Effiziente lineare Algebra .....	19
3.2.2 TOP-basierte Dimensionierungsstrategie .....	20
3.3 Nutzung von NuT zur Behandlung der linearen Algebra während der Startrechnung .....	23
3.4 Weiterentwicklung und Ausbau der Differentialgleichungsnumerik im Numerical Toolkit .....	25
3.4.1 ATHLET-spezifische Anforderungen an die Zeitintegration .....	25
3.4.2 Der FiterRK-Ansatz .....	28
3.4.3 Anforderung an die technische Umsetzung .....	34
3.5 Einbringen der Kommunikations-Bibliothek MMA .....	35
3.5.1 Einbinden von MMA .....	36
3.5.2 Funktionsweise von MMA .....	37
3.6 Speichern und Laden der Jacobimatrix im Restart-Szenario mittels NuT	40

<b>4</b>	<b>AP3 – Aktualisierung der DGL-Numerik in der AC<sup>2</sup>-Komponente</b>	
	<b>COCOSYS .....</b>	<b>45</b>
4.1	Analyse des Ist-Zustandes der Differentialgleichungsnumerik im Thermo-hydraulik-Modul von COCOSYS .....	45
4.2	Zugang von COCOSYS zur AC <sup>2</sup> -einheitlichen Handhabung der DGL-Numerik über die NuT-Architektur .....	47
4.2.1	Funktionsbasierte Modifikationen am Code des Thermohydraulik-Moduls	47
4.2.2	Funktionsbasierte Modifikationen des anwendungsseitigen NuT-Codes .	50
4.2.3	Ausstehendes Bindeglied .....	51
4.3	Ausarbeitung einer Kommunikations- und Parallelitätsstrategie im Arbeitsverbund der AC <sup>2</sup> -Komponenten COCOSYS und ATHLET mit dem Numerical Toolkit .....	52
<b>5</b>	<b>AP4 – Validierung der Modifikationen.....</b>	<b>55</b>
5.1	Release und QS zu NuT 1.0 .....	55
5.2	Kohärenz-Test zur Berechnung der Jacobimatrix.....	56
5.2.1	Motivation .....	57
5.2.2	Mögliche Auswirkungen einer Kohärenz-Diskrepanz .....	57
5.2.3	Funktionsweise und Anwendungsbereich des Kohärenz-Tests .....	61
5.2.4	Empfehlungen zum Einsatz des Tests .....	61
5.3	Entwicklung dedizierter Programme zum Austesten einzelner Funktionalitäten .....	62
<b>6</b>	<b>Fazit und Ausblick .....</b>	<b>63</b>
	<b>Literaturverzeichnis.....</b>	<b>65</b>
	<b>Abbildungsverzeichnis .....</b>	<b>70</b>
	<b>Tabellenverzeichnis .....</b>	<b>71</b>
	<b>Verzeichnis von Codeabschnitten.....</b>	<b>72</b>
<b>A</b>	<b>Anhang .....</b>	<b>73</b>
A.1	Masterarbeit Ravil Dorozhinskii .....	73
A.1.1	Effiziente lineare Algebra .....	73
A.1.2	Implementation eines Doppelpuffers .....	73
A.1.3	Betreuung .....	73

# 1 Einleitung

Die GRS entwickelt und validiert das Programmsystem AC<sup>2</sup>, welches die Programm-Module ATHLET und ATHLET-CD (ATHLET/CD) sowie COCOSYS und NuT umfasst. Hierbei dient ATHLET (Analyse der Thermo-Hydraulik von LEcks und Transienten) zur Simulation des Kreislaufverhaltens von Transienten und Störfällen. Mittels der Erweiterung ATHLET-CD (Core Degradation) besteht zusätzlich die Möglichkeit, Unfälle mit Kernzerstörung zu simulieren. Der Einsatz von COCOSYS (COntainment-COdesystem) erlaubt die Berücksichtigung von Stör- und Unfallabläufen in Containments.

Als jüngstes Mitglied der AC<sup>2</sup>-Code-Familie ist NuT (Numerical Toolkit) mit dem Erstrelease von AC<sup>2</sup> in der Version 2019 dem Programm-Kanon hinzugefügt worden. Das Toolkit bietet einen einfach zu bedienenden Zugriff auf dedizierte numerische Bibliotheken, welche skalierbare Lösungsalgorithmen zur linearen Algebra bereitstellen. Siehe Abschnitt 1.1 hinsichtlich des Anwendungskontextes im Rahmen von AC<sup>2</sup>. Die Schnittstelle zu NuT wird über den Kommunikationsstandard MPI (Message Passing Interface) realisiert.

Bereits im Rahmen des Vorgängervorhabens RS1530 wurde ein Plugin für ATHLET/CD erstellt, um die lineare Algebra im Zeitintegrationsprozess der Simulationsläufe optional mittels NuT-Algorithmen anzugehen. Mitunter konnten deutliche Performance-Gewinne erzielt werden. Siehe hierzu den Abschlussbericht zu RS1530 /STE 17b/.

Der Fokus dieses Projektes liegt auf der stetigen Verbesserung des Numerical Toolkit selbst sowie der Interaktion mit anderen AC<sup>2</sup>-Komponenten. Entsprechend ergeben sich die Arbeitspakete, die je einem Kapitel zugeordnet sind

- Kap. 2 – AP1 *Ausbau und Verbesserung der Parallelisierungs- und Funktionsmächtigkeit innerhalb des Numerical Toolkit*
- Kap. 3 – AP2 *Verbesserung und Ausbau der Zusammenarbeit von ATHLET mit dem Numerical Toolkit*
- Kap. 4 – AP3 *Aktualisierung der DGL-Numerik in der AC<sup>2</sup>-Komponente COCOSYS*
- Kap. 5 – AP4 *Validierung der Modifikationen*

Zusätzlich ergaben sich die Arbeiten zum Release, welche Kapitel 5 zugeordnet sind. Die in diesem Bericht protokollierten Arbeiten berücksichtigen bereits die Anpassungen des Arbeitsplanes, welche durch den Änderungsantrag vom Juni 2019 eingehen.

## 1.1 Anwendungsszenarien von NuT

Ausgehend von den Erhaltungsgleichungen für Masse, Energie und Impuls bezüglich einer flüssigen sowie gasförmigen Phase werden für die thermohydraulischen Berechnungen in AC<sup>2</sup> finite Volumenansätze bezüglich der Raumdimension genutzt. Diese Ansätze führen zu einem Netzwerk aus *control volumes*, welche durch ein hierzu verschobenes Gitter aus *junctions* miteinander verknüpft sind. Auf den *junctions* werden die Masseströme simuliert, während in den Volumen Zustandsgrößen wie Temperatur, Druck und relative Aufteilung von Flüssig- und Gasphase das System beschreiben. Trotz grundsätzlich ähnlicher Strukturen und Konzepte finden in ATHLET und dem Thermohydraulik-Modul von COCOSYS verschiedene Modelle Anwendung. Details hierzu finden sich in /AUS 19, Kap. 2/ beziehungsweise /ARN 19, Absch. 2.1/.

Gleich ist jedoch, dass das Systemverhalten letztendlich mittels eines Anfangswertproblems beschrieben wird:

$$y' = f(t, y), \quad y(t_0) = y_0. \quad (1.1)$$

Dieses ist in der Regel als *steif* anzusehen, was bedeutet, dass es impliziter Verfahren bedarf, um es effizient numerisch zu lösen. Sowohl ATHLET als auch das Thermohydraulik-Modul von COCOSYS bedienen sich des gleichen linear-impliziten Verfahrens, siehe Unterabschnitt 3.4.1.1 beziehungsweise Abschnitt 4.1.

Als Teilproblem in der Zeitintegration ergibt sich durch den impliziten Charakter des Verfahrens das Lösen von linearen Gleichungssystemen, in welche eine Approximation  $J$  der Jacobimatrix zu (1.1) eingeht:

$$(I - h\gamma J)\Delta y = b. \quad (1.2)$$

Hierbei ist  $h$  eine gegebene Zeitschrittweite und  $\gamma$  ein methodenspezifischer Parameter. Das Numerical Toolkit setzt genau in diesem Kontext der Differentialgleichungsnumerik und numerischen linearen Algebra an. Die einzelnen Arbeitspunkte dieses Projektes zielen darauf ab, das Numerical Toolkit weiter hinsichtlich Robustheit, Bedienbarkeit und Funktionsumfang im Rahmen des AC<sup>2</sup>-Kontextes zu verbessern.

## 1.2 Weitere allgemeine Anmerkungen

- Viele Arbeiten zu diesem Projekt zeigen einen starken software-technischen Charakter. Für eine klare Darstellung wird auf eine Übersetzung etablierter englischer Fachbegriffe verzichtet. Entweder ist der deutsche Bezeichner sehr unüblich oder gar nicht erst kreiert worden.
- Stehende Begriffe innerhalb der jeweiligen Kontexte werden in der Regel nicht gesondert erläutert. Zur Klärung kann auf ein breites Spektrum an Basisliteratur zugegriffen werden. Nachfolgend einige Empfehlungen:
  - Softwareentwicklung in C++ und zu MPI  
/LIP 12/,/MAR 08/, /MES 15/
  - Numerik gewöhnlicher Differentialgleichungen  
/HAI 93/,/HAI 96/,/BUT 16/
  - Numerik linearer Gleichungssysteme  
/MEI 11/,/SAA 03/,/DAV 06/
- Mehrere Tausend Zeilen Code wurden im Rahmen der Arbeiten für das Projekt geschrieben. Diese sind dem Bericht nicht beigefügt. Vereinzelt kann es im Sinne der besseren Erläuterung zu der Darstellung einzelner Codezeilen kommen.
- Mit den abgeschlossenen Arbeiten aus Abschnitt 2.3 liegt das Numerical Toolkit in der Version 2.0.0 vor. Es ist nicht mehr kompatibel zu der aktuellen Release-Version 1.0.1. Es ist geplant, die Neuerungen in der nächsten Major-Version von AC<sup>2</sup> Berücksichtigung finden zu lassen und somit der allgemeinen Nutzerschaft zur Verfügung stellen zu können.



## **2 AP1 – Ausbau und Verbesserung der Parallelisierungs- und Funktionsmächtigkeit innerhalb des Numerical Toolkit**

### **2.1 Untersuchungen zur Einbringung von OpenMP in NuT**

OpenMP-Techniken bieten einfach zu handhabende Mechanismen, um einzelne Code-Abschnitte in der Ausführung zu parallelisieren und dadurch ggf. einen Performancegewinn zu erzielen. Der Einsatz ist hierbei sehr vom jeweiligen Kontext abhängig. Insofern wird zu diesem Arbeitspunkt zunächst auf das strukturelle NuT-Umfeld eingegangen. Konkrete Tests werden in Abschnitt 2.1.1 thematisiert.

Das Numerical Toolkit liefert leicht zu bedienende High-Level-Funktionen, um über eine MPI-Schnittstelle Host-Applikationen bei Aufgaben der numerischen linearen Algebra zu unterstützen. Im AC<sup>2</sup>-Kontext wird dies für die linearen Systeme genutzt, welche während der (linear-)impliziten Zeitintegration der Modellgleichungen auftreten. Im Wesentlichen gliedert sich das Toolkit hierbei in die Ebenen

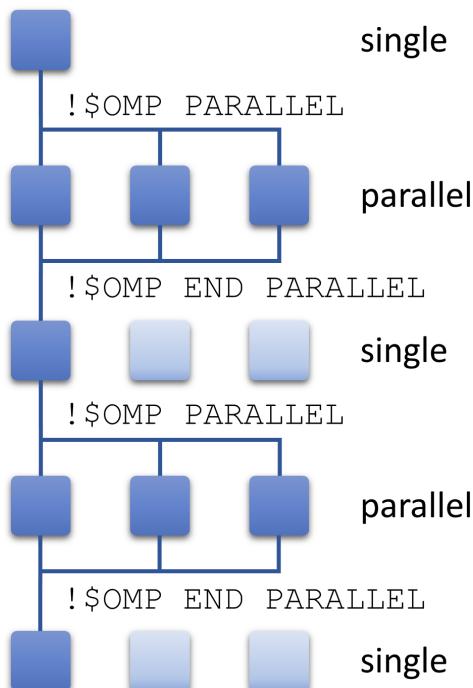
- Kommunikation mit Host-Applikationen via MPI-Direktiven,
- Verwalten von abstrakten Daten und Meta-Algorithmen,
- Zugriff auf dedizierte Numerik-Bibliotheken zur grundlegenden Datenhaltung und Berechnung.

**Bemerkung 2.1.** Die Untersuchungen zur etwaigen Einbindung von OpenMP-Mechanismen in NuT basieren auf der FORTRAN-Implementation des Toolkits, welche in den relevanten Codeabschnitten den Release-Versionen 1.0.0 und 1.0.1 entspricht, s. Abschnitt 5.1 für weitere Anmerkungen zu Releases. Die in Abschnitt 2.3 erläuterten Arbeiten zum Umschreiben des Toolkits nach C/C++ fanden nach den hier beschriebenen Tätigkeiten statt. Da die OpenMP-Konzepte die gleichen sind, die NuT-Funktionalitäten analog übernommen wurden und aufgrund der in Abschnitt 2.1.1 aufgezeigten Ergebnisse wurde davon abgesehen, sich des Themas auch für die C/C++-Implementation anzunehmen.

Eine Anwendung von OpenMP-Techniken setzt eine Shared-Memory-Architektur voraus, s. /OPE 18, Sec. 1.4/, wodurch die in oben stehender Aufzählung angeführte Kommunikationsebene aufgrund der MPI-Mechanismen außen vor ist. Analoges gilt für die genutzten dedizierten Numerik-Bibliotheken, da diese nicht Teil des NuT-Codes sind. Zwar besteht grundsätzlich die Möglichkeit zur Modifikation. Jedoch ist es sinnvoller, dies ggf. als Anregung an die eigentlichen Entwickler weiterzugeben. Es bleibt somit die

Ebene der Meta-Algorithmen hinsichtlich einer näheren Untersuchung zur Anwendung von OpenMP-Techniken.

OpenMP zeigt seine Stärken bei der Parallelisierung von Schleifen, indem möglichst disjunkte Abschnitte der Schleifeniteration auf verschiedene Threads aufgeteilt werden, welche wiederum dynamisch und automatisiert erzeugt werden, vgl. Abbildung 2.1. Zu diesen Anforderungen passende Algorithmen im Toolkit beziehen sich auf das Verwalten des Alloziierens und Setzens von Matrixelementen sowie auf die Organisation des Seeding-Prozesses. Zur Klärung des Seeding-Begriffs siehe Abschnitt 5.2.2.



**Abb. 2.1** Prinzip dynamisch erstellter Threads per OpenMP

### 2.1.1 Test und Resultate

Im Sinne des geforderten Austestens von OpenMP-Funktionalitäten innerhalb NuTs wurde sich für die Subroutine `nut_matCreateFromBlockCSR` entschieden, welche auf Basis von Strukturinformation im Format Compressed Sparse Row (CSR) eine Matrix in PETSc alloziert und mit 0-Werten vorbelegt, s. Code 2.1. Hierbei wird auf den anwendungsspezifischen Fall eingegangen, dass eine Block- und eine darin enthaltene Elementebene der Matrix existiert, was zu einer Komplexitätssteigerung der aufgerufenen Routinen `nut_matAllocateBlockCSR` und `nut_matSetBlockCSR` führt. Die OpenMP-Direktiven finden sich in `nut_matAllocateBlockCSR` wieder, da hier die Logik eine höhere Schachtelung erfährt, s. Code 2.2. Grundsätzlich gilt aber für beide Routinen, dass die Matrix

blockzeilenweise durchgegangen wird, und die Blockzeilen sind in Bezug auf die ausführte Logik unabhängig voneinander. In `nut_matAllocateBlockCSR` gibt es lediglich gemeinsame Zählvariablen (`d_nnz` sowie `o_nnz`), welche entsprechend in der OpenMP-Direktive berücksichtigt werden. Bezüglich der Disjunktheit der Daten bestehen somit ideale Voraussetzungen.

**Code 2.1** OpenMP-Performance-Messung I: Timer-Einbettung anhand der Beispielroutine `nut_matCreateFromBlockCSR`

```

...
! declaration of global variables
timer timer_1 ! abstract timer object
timer timer_2 ! abstract timer object
file output_file ! abstract file object
...
subroutine nut_matCreateFromBlockCSR(nut_A, nut_ePtr, nut_rowPtr,
    nut_colInd, activeRows_array, activeColumns_array, forColoring)
    ...
        timer_1.start = MPI_Wtime ()
        ...
        call nut_matAllocateBlockCSR(nut_A, nut_ePtr, nut_rowPtr,
            nut_colInd, activeRows_array, activeColumns_array)
        call nut_matSetBlockCSR(nut_A, nut_ePtr, nut_rowPtr,
            nut_colInd, activeRows_array, activeColumns_array)
        ...
        timer_1.end = MPI_Wtime ()
        timer_1.time = timer_1.end - timer_1.start
        output_file.append ( timer_1.time , timer_2.end )
    ...
end subroutine

```

**Code 2.2** OpenMP-Performance-Messung II: Anweisungsdefinitionen anhand der Beispielroutine `nut_matAllocateBlockCSR`

```

subroutine nut_matAllocateBlockCSR(nut_A, nut_ePtr, nut_rowPtr,
    nut_colInd, activeRows_array, activeColumns_array)
    ...
        timer_2.start = MPI_Wtime ()
        !$OMP PARALLEL PRIVATE ( block_y , block_ys , element_y ,
            block_x ,
            !+ block_xs , element_x ,i,j,x,y,y_s , y_e)
        ...
        !$OMP DO SCHEDULE ( DYNAMIC ) REDUCTION (+: d_nnz , o_nnz
            )

```

```

do i = 1, blocks                                8
    block_y = i                                  9
    block_ys = ePtr_array(block_y+1) - ePtr_array(block_y) 10
    element_y = ePtr_array(block_y)               11
    if (element_y .gt. high) then                12
        exit
    end if
    if (element_y + block_ys .lt. low) then      13
        cycle
    end if
    do j = rowPtr_array(i), rowPtr_array(i+1)-1 14
        block_x = colInd_array(j)
        block_xs = ePtr_array(block_x+1) - ePtr_array(block_x) 15
        element_x = ePtr_array(block_x)
        do x = element_x, element_x + block_xs -1 16
            y_s = max(element_y, low)
            y_e = min(element_y + block_ys -1, high) 17
            if (x .ge. low .and. x .le. high) then 18
                do y = y_s, y_e
                    if((activeRows_array(y) .and. 19
                        activeColumns_array(x)) .or. ( 20
                        includeDiagonals .and. (x .eq. y))) then 21
                        d_nnz(y-low) = d_nnz(y-low) + 1 22
                    end if
                end do
            else
                do y = y_s, y_e
                    if(activeRows_array(y) .and. 23
                        activeColumns_array(x)) then 24
                        o_nnz(y-low) = o_nnz(y-low) + 1 25
                    end if
                end do
            end if
        end do
    end do
end do
! $OMP END DO
! $OMP END PARALLEL
timer_2.end = MPI_Wtime ()
timer_2.time = timer_2.end - timer_2.start
...
end subroutine

```

**Tab. 2.1** OpenMP-Performance für K3-2 mit den OpenMP-Anweisungen aus Code 2.2

Iteration	1 Thread			4 Threads		
	timer_1	timer_2	Gesamtlaufzeit	timer_1	timer_2	Gesamtlaufzeit
1	1.76E-01	7.65E-03	–	1.76E-01	7.65E-03	–
2	5.44E-02	8.58E-03	–	5.43E-02	8.58E-03	–
3	5.19E-02	8.58E-03	–	5.17E-02	8.72E-03	–
4	2.38E-02	8.63E-03	–	2.36E-02	8.65E-03	–
5	5.36E-02	8.60E-03	–	5.34E-02	8.59E-03	–
6	5.17E-02	8.58E-03	–	5.17E-02	8.57E-03	–
7	2.39E-02	8.89E-03	–	2.38E-02	8.70E-03	–
8	5.34E-02	8.59E-03	–	4.74E-02	7.66E-03	–
9	2.22E-02	8.64E-03	–	1.99E-02	7.64E-03	–
10	2.26E-02	8.68E-03	–	3.11E-02	7.65E-03	–
Summe:	5.33E-01	8.54E-02	5.58E+03	5.33E-01	8.24E-02	5.31E+03

Es wurden Simulationsläufe mit verschiedenen Inputfiles gestartet. Ggf. wurde der Endzeitpunkt der Zeitintegration angepasst, da nicht der eigentliche Verlauf von Bedeutung ist, sondern nur der Einfluss der OpenMP-Direktiven auf die Laufzeit. Man beachte, dass keine wesentliche Floatingpoint-Arithmetik in `nut_matAllocateBlockCSR` auftritt. Die Modifikationen durch OpenMP haben somit keinerlei Einfluss auf die Zahlenwerte des Simulationsergebnis.

Da sich die Qualität der OpenMP-bezogenen Resultate de facto unabhängig von der Eingabe gezeigt hat, sind in Tabelle 2.1 nur die Zeitangaben für den Anwendungsfall K3-2 gegeben. Der Datensatz basiert auf den Informationen aus /TER 08/ mit einigen Modifikationen im Detail. Hierbei wird das Abschalten einer der vier Hauptkühlmittelpumpen in einem VVER-1000 Kern simuliert. Heißer Kanal und ein Nachbarelement sind stabweise nodalisiert. Entsprechend groß ist das System mit 128.673 Systemvariablen. Die zugehörige Jacobimatrix hat im Schnitt knapp drei Millionen Nicht-Null-Elemente und ist somit dünnbesetzt.

Die Rechnungen wurden auf einem Node des Linux-Clusters der GRS ausgeführt. Mit zehn verfügbaren physischen Kernen gab es keine Überbelegung. Auch im Falle eines Threads wurde OpenMP aktiviert. Ansonsten wurde das Numerical Toolkit grundlegend sequentiell betrieben, d. h. nur ein MPI-Prozess wurde angesetzt. Die Timer-Angaben `timer_1` und `timer_2` entsprechen den gleichnamigen Größen in Code 2.2.

Wie anhand der Zeitwerte in Tabelle 2.1 zu erkennen ist, nimmt die Zeit in `nut_mat-`

`AllocateBlockCSR` keinen maßgeblichen Einfluss auf die Gesamlaufzeit des Simulationslaufes. Selbst auf Ebene von `nut_matCreateFromBlockCSR`, welches neben dem Allozieren noch das Setzen von Werten in der Matrix umfasst, ergibt sich kein nennenswerter Beitrag zur Simulationslaufzeit. Entsprechend gering fällt auch die Beschleunigung durch den Einsatz von OpenMP aus.

### 2.1.2 Schlussfolgerung

Im Kontext der direkt in NuT vorhandenen Schleifen-Algorithmen ist mit `nut_matAllocateBlockCSR` bereits einer der komplexesten gegeben. Die Crux ist hier, dass sich die eigentlichen Daten-Operationen als zu leichtgewichtig erweisen, um von Thread-Parallelisierung profitieren zu können. Dies gilt auch für alle weiteren Algorithmen dieser Art in NuT. Somit lohnt es sich auch nicht andere Scheduling-Strategien, /SPE 16/, zu untersuchen oder ein Task-basiertes Vorgehen, /PAS 13/, auszutesten. Die zeitliche Signifikanz der von OpenMP beeinflussten Operationen ist zu gering.

Es ergibt sich somit das Fazit, dass im aktuellen Stand des Funktionsumfanges von NuT OpenMP-Direktiven im NuT-Code keinen Mehrwert liefern. Die hiermit verknüpften Ressourcen (CPU-Kerne) sollten von daher NuT direkt in Form von MPI-Prozessen zur Verfügung gestellt werden und/oder ATHLET für eine OpenMP-basierte Parallelisierung. Wie anhand der Performancemessungen in /STE 17b/ zu sehen, geht dies mit einer signifikanten Verbesserung der Laufzeiten einher.

## 2.2 Weitere Optionen zur Präkonditionierung iterativer Verfahren

Aufgrund der (linear-)impliziten Zeitintegration in ATHLET sind pro Zeitschritt mehrere lineare Gleichungssystem zu lösen, welche die Jacobimatrix am aktuellen Zeitpunkt oder eine Approximation hierzu involvieren. Das Numerical Toolkit bietet verschiedene Presets, /STE 20, Absch. 4.4/, die linearen Systeme zu lösen. Entsprechend der Aufgabenstellung dieses Arbeitspunktes wurde die Incomplete LU Decomposition  $ILU(k)$  aus dem Programm Paket Hypre, /FAL 20/, hinsichtlich einer Anwendbarkeit im NuT-Kontext untersucht. Hierbei käme der  $ILU(k)$  die Rolle eines Präkonditionierers für iterative Verfahren zur Lösung obiger linearer Systeme zu.

Das Numerical Toolkit nutzt hauptsächlich PETSc, /BAL 20/, als externe Numerik-Bibliothek, bedient sich aber auch PETSc-interner Schnittstellen zu weiteren eigenständigen Bibliotheken wie MUMPS, /AME 20/, oder METIS, /KAR 20/. Eine analoge Vorgehensweise war für die  $ILU(k)$  aus dem Hypre-Paket angedacht. Es ergibt sich jedoch das

ernüchternde Ergebnis, dass dieser Präkonditionierer weder unter Windows noch unter Linux stabil läuft. Eine Rückfrage bei den Entwicklern von PETSc ergab, dass der Fehler auf der Seite von Hypre zu suchen sei, bisher jedoch von den zuständigen Entwicklern keine positive Resonanz hinsichtlich der Behebung des Missstandes erfolgt sei, s. hierzu auch /BRO 17/. Die Arbeiten zu diesem Punkt wurden hauptsächlich 2018 durchgeführt. Auch in 2020 ist nicht bekannt, wann oder ob ein Bugfix erscheinen wird.

Um dennoch dem Arbeitspunkt im Kern genüge tun zu können, wurde PETScs Implementierung der level-basierten ILU( $k$ ) für dünnbesetzte Matrizen, /SAA 03, Unterabsch. 10.3.3/, in Kombination mit GMRES, /MEI 11/, dem Kanon an Solver-Presets hinzugefügt. Da es sich um eine sequentielle Implementierung in PETSc handelt, steht das Preset nur zur Verfügung, wenn ein einzelner NuT-Prozess genutzt wird. Im Zuge der Implementierungsarbeiten zeigte sich, dass PETScs Standardalgorithmus zur Fill-in-Reduktion (nested dissection) nur unzureichende Ergebnisse lieferte, was oft zur Divergenz des GMRES-Prozesses führte. PETSc bietet jedoch weitere Fill-in-Reduzierer. Der QMD-Algorithmus (*Minimum Degree Algorithm using Quotient Graphs*), /GEO 80/, liefert deutlich bessere Resultate. Zur Performance siehe auch Abschnitt 5.3. Das Preset zur unvollständigen LU-Zerlegung ist über `-solver ilu_k-gmres` aufrufbar, wobei das Level  $k$  ganzzahlig  $\geq 0$  zu wählen ist.

Das Level gibt an, welches Nicht-Null-Muster für eine LU-Zerlegung zur Verfügung steht: ILU(0) darf sich für die Zerlegung nur der Besetzungstruktur der Ausgangsmatrix  $A$  bedienen. Werden die Faktoren  $L$  und  $U$  der ILU(0) multipliziert, ergibt sich eine Matrix, deren Nicht-Null-Muster von der ILU(1) zu berücksichtigen ist. Für die ILU(2) wird dann das Nicht-Null-Muster zu Grunde gelegt, welches die Matrix aufweist, die sich aus dem Produkt der  $L$ - und  $U$ -Faktoren der ILU(1) ergibt. Analog kann diese rekursive Definition für  $k > 2$  weitergeführt werden. Der Vorteil an der level-basierten ILU ist die Möglichkeit, den benötigten Speicherplatz vor der eigentlichen numerischen Zerlegung ermitteln zu können. Nachteilig ist, dass ein level-basiertes Vorgehen bei indefiniten Matrizen nicht zwingend zu einer Verbesserung der Qualität des Präkonditionierers führt, wenn  $k$  erhöht wird. Im Kontext der (linear-)impliziten Zeitintegration in AC<sup>2</sup> kann stets davon ausgegangen werden, dass die betrachteten Matrizen indefinit sind. Alle bisherigen Beispielrechnungen zeigen jedoch konvergentes Verhalten für  $k = 2$ , sofern der QMD-Algorithmus zur Fill-in-Reduzierung genutzt wird. Man beachte, dass jener Algorithmus auf die Ursprungsmatrix angewandt wird. Die hieraus resultierende Matrix gilt dann als Basis für die ILU-Zerlegungen.

Alle Beispiele liefen bisher stabil für das Level  $k = 2$ . Die gute Performance motivierte die

Entscheidung, den QMD-Algorithmus auch für die Standard-LU-Zerlegung in PETSc und somit für die zugehörigen Solver-Presets zu nutzen. Auch hier zeigte sich die Reduzierung an benötigtem Speicher. Das Preset `ilu_k-gmres` kann zu kleinen Verbesserungen in der Geschwindigkeit führen, ist aber nicht von gleicher numerischer Stabilität wie das sequentielle Standard-Preset `lu-gmres`. Ebenso ist vom Nutzer eine zusätzliche Eingabe erforderlich. Somit sollte `ilu_k-gmres` hauptsächlich als Option für Vergleichszwecke angesehen werden.

## 2.3 Umschreiben des Numerical Toolkit-Codes von FORTRAN nach C++

Der NuT-Code wurde von FORTRAN nach C++ umgeschrieben. Dies wirkt sich vorteilhaft auf sämtliche NuT-Softwarekomponenten und deren Anbindung aus.

Ein wesentlicher Nutznießer ist hierbei die Bibliothek `libnut-core`. Diese ist an der Basis der NuT-Softwarearchitektur angesiedelt und greift direkt auf PETSc-Funktionalitäten zu, welche in der dynamisch gelinkten Bibliothek `libpetsc` gekapselt sind. Durch den Wechsel der Programmiersprache ergibt sich zum einen der Vorteil, die native C-Schnittstelle von PETSc nutzen zu können, welche umfangreicher und stabiler als die FORTRAN-Variante ist. Zum anderen konnte mithilfe von C++ die Typsicherheit und Abstraktion von `libnut-core` zum NuT-Worker merklich verbessert werden.

Weiterhin konnte die Wartbarkeit und Flexibilität deutlich erhöht werden. Im Gegensatz zu C und FORTRAN 77 bietet modernes FORTRAN keine standardisierte Binärschnittstelle (ABI). Dies hat den Nachteil, dass die ABI einer kompilierten Bibliothek mit dem eingesetzten Compiler variieren kann. Um dennoch ausreichende Kompatibilität gewährleistet zu können, müssen alle Komponenten mithilfe des gleichen Compilers erstellt werden. Da es sich bei der Zusammenarbeit in der Regel um unabhängige Softwarekomponenten handelt, die je ihren eigenen Entwicklungsprozess besitzen, ist dies eine Nebenbedingung, die besagte Prozesse unnötig eng aneinanderbindet und dadurch unflexibel und mitunter zeitraubend in ihrer Bearbeitung macht. Z. B. kann in der Entwicklung des NuT-Workers nicht direkt auf die Binärform von `libnut-core` zugegriffen werden, ohne vorher sicherzustellen, dass die verwendeten Compiler kompatibel sind. Somit ist es schwierig, eine für alle Entwickler gültige Binärform von `libnut-core` zur Verfügung zu stellen. Mitunter müsste neu kompiliert werden, was Zeit kostet und wodurch sich evtl. Seiteneffekte ob eines anderen Compilers ergeben. Dieser Nachteil stellte sich im Laufe der Projektarbeiten heraus und konnte mit der Umschreibung nach C++ eliminiert werden. Zwar unterstützt auch C++ selbst keine standardisierte ABI, bietet allerdings die

Möglichkeit, öffentliche Schnittstellen abwärtskompatibel in C zu definieren und damit die nötige Kompatibilität zu gewährleisten.

Aufgrund der mit der Umschreibung des Codes einhergehenden Standardisierung und Verbesserung der Schnittstellen, ist die API zwischen Host-Applikation (Anwendungsseite) und NuT-Plugin, sowie NuT-Plugin und NuT-Worker nicht mehr zwangsläufig zueinander kompatibel zu den bisherigen FORTRAN Implementierungen. Um Konformität mit *Semantic Versioning* Punkt 8, /PRE 20/, zu wahren, muss die Major-Version von 1 auf 2 erhöht und damit Minor- und Patch-Version auf 0 zurückgesetzt werden. Die nächste Veröffentlichung erfolgt damit als Version 2.0.0.

Die Architektur bietet mehrere Schnittstellen in unterschiedlichen Schichten an, so dass gängige Anwendungsszenarien bedient werden können.

- NuT-Base: Hier kann sich eine Anwendung nativ über die C++-Schnittstelle an der vorgefertigten abstrakten Algorithmik der NuT-Entitäten und deren Basisroutinen in `libnut-core` bedienen. In diesem Fall ist NuT voll in der Host-Applikation eingebettet. Dieser Ansatz ist aufgrund des direkten Zugriffs auf die NuT-Algorithmen sehr hilfreich, um beispielsweise im Sandkastenprinzip ohne weitere Hürden, numerische Methoden entwickeln und testen zu können. Es wird lediglich innerhalb der gekapselten PETSc-Routinen miteinander kommuniziert. Damit könnte auch leicht eine MPI-freie Konfiguration erstellt werden, wenn ein Anwendungsfall dies erforderlich machen sollte. Außerdem ist dieser Ansatz für Host-Applikationen geeignet, die bereits voll über MPI parallelisiert sind. Hier fügt sich NuT mit der selben Parallelität nativ in die Host-Applikation ein.
- NuT-Host/NuT-Worker: Sobald eine sequenzielle Host-Applikation auf die parallelen Algorithmen von NuT zugreifen soll, erfolgt die Anbindung über die NuT-Host Schnittstelle. Nach außen hin kann diese Host-Applikationen in C, C++ und FORTRAN bedienen, indem sie abstrakte Funktionen zur Verfügung stellt. Intern werden die übergebenen Daten serialisiert und die Ein- und Ausgabe per MPI mit dem NuT-Worker ausgetauscht. Im Wesentlichen handelt sich dabei um eine abgewandelte Form des *Remote Method Invocation* (RMI-)Ansatzes, der sich zusätzlich um den Aspekt der verteilten Parallelität kümmert. Die Anbindung an die Host-Applikation kann per Linker statisch oder dynamisch erfolgen. Ebenfalls stehen entsprechende Schnittstellen zur Verfügung, die das dynamische Laden per Plugin-Mechanismus zur Laufzeit ermöglichen.

Im Einzelnen haben sich folgende Arbeiten und Entscheidungen ergeben:

- Die komplette Repository-Struktur zu NuT wurde neu aufgelegt, um besser die funktionale Aufteilung des Codes zu repräsentieren. Dies wurde gleich in Git umgesetzt, welches den neuen Standard für die Versionierungskontrolle der Codes im AC<sup>2</sup>-Kontext darstellt.
- Aufgrund des Sprachwechsels konnten verschiedene leistungsstarke Programmier-Idiome zum Einsatz gebracht werden. Unter anderem sind dies RAII (*Resource Acquisition Is Initialization*), das *pimpl-idiom, rule of five* sowie das *abstract factory* und *fabric method pattern*. Ebenso wurden durchgehend *smart pointer* statt *raw pointer* eingesetzt und sich an dem *single responsibility principle* orientiert. In der Summe führt dies zu einem robusteren, effizienteren und übersichtlicheren Code. Das Konzept der *dependency injection* ist ebenfalls zum Einsatz gekommen, siehe hierzu auch Abschnitt 4.3.
- Nach wie vor gibt es im Numerical Toolkit eine logische Aufteilung zwischen Kommunikation, abstrakter Algorithmik, Implementation unter Zugriff auf externe Quellen sowie die Kapselung jener Quellen. Auf Datenebene wird dies zweigliedrig umgesetzt. Die Implementation wird nun statisch in den NuT-Worker gelinkt, welcher ebenfalls und wie zuvor Kommunikation und Algorithmik umfasst. PETSc, MUMPS, METIS und (Sca)LAPACK sind in einer dynamischen Bibliothek aggregiert. Der Zugriff erfolgt über die PETSc-eigene C-Schnittstelle. Damit ist neben erhöhter Funktionalität das wesentliche Kriterium der Compiler-Unabhängigkeit der PETSc-Schnittstelle gegeben. Eine solche Unabhängigkeit gilt auch für das neu erstellte NuT-Plugin, welches Kommunikationspartner vom Toolkit auf der Anwendungsseite ist. Näheres hierzu siehe abermals Abschnitt 4.3. Man beachte, dass eine Compilerunabhängigkeit keine Unabhängigkeit von einer bestimmten Laufzeitumgebung nach sich zieht. Die PETSc-Bibliothek wird z. B. für die Nutzung unter MS Windows per Visual Studio in einer Cygwin Umgebung gebaut und binär dem NuT-Repository als *external* zugewiesen. Dies röhrt daher, dass der Erstellungsprozess zur PETSc-Bibliothek komplex und länglich ist. Als Konsequenz ergibt sich, dass NuT mit einer kompatiblen Laufzeitumgebung zu kompilieren ist. Dies stellt kein Problem dar, es ist nur im Hinterkopf zu behalten.
- Um Compiler- und Linkvorgänge plattformübergreifend zu vereinfachen, wurde auf das Programmierwerkzeug CMake, /SCO 20/, zugegriffen. Dieses ist mit einer BSD-artigen Lizenz versehen und somit kompatibel zu den GRS-Richtlinien. CMake bietet eine eigene Skriptsprache, worüber Abhängigkeiten und Zuordnungen in einem gegebenen Kompilier- und Linkprozess definiert werden können. Auf Basis der plattformunabhängigen Skripte kann dann CMake konkrete Kompiliervorschriften für eine gegebene

Plattform (Linux/Windows/MacOS) erstellen. Damit steht ein einheitlicher Standard für den NuT-Kompilier- und Link-Prozess bereit, der zentral gepflegt werden kann und selbstständig automatisch externe Abhängigkeiten (e. g. MMA, siehe Abschnitt 3.5) über entsprechende Repository-Downloads auflöst.

- Zwei Versionschecks zur Nutzung der PETSc-Bibliothek wurden implementiert. Im Toolkit ist ein Zwei-Tupel hinterlegt, welches die PETSc-Version kennzeichnet, die mit dem Toolkit zu benutzen ist. Ein statischer Check zur Kompilier-Zeit ist über ein *static assert* realisiert, welches sich der über Header-Files verfügbaren PETSc-Macros *PETSC\_VERSION\_MAJOR* und *PETSC\_VERSION\_MINOR* bedient. Besteht hier bereits Diskrepanz bricht der Kompiliervorgang mit einer aussagekräftigen Fehlermeldung ab. Da die externen Bibliotheken dynamisch gelinkt werden, wurde ein zweiter Test implementiert, der zur Laufzeit während der Initialisierung aufgerufen wird. Hierbei wird auf die PETSc-eigene Funktion *PetscGetVersionNumber* zurückgegriffen und ein analoger Vergleich durchgeführt.
- Der über PETSc-interne Schnittstellen verfügbare direkte LA-Löser MUMPS beinhaltet eine Schätzung des Speicheraufwandes, der für die numerische LR-Zerlegung der Matrizen benötigt wird. Ist der reservierte Speicher ungenügend, wird iterativ zusätzlicher Speicher angefordert. Dieser Prozess wurde bereits in der FORTRAN-Fassung zur Verfügung gestellt. Die C++-Implementierung liefert eine Verfeinerung dahingehend, dass nun Grenzwertabfragen und Exception-Handling eingearbeitet wurden.
- Weitere Anpassungen wurden auf Host-Seite vorgenommen, die aufgrund ihres Inhaltes AP3 zugeordnet sind und somit in Kapitel 4 erläutert werden.

Abstrahiert gesehen besteht die Schnittstelle zu NuT aus einer Menge von Funktionen mit zugehörigen Ein- und Ausgabeparametern. Diese Parameter können in zwei verschiedene Gruppen aufgeteilt werden. Bei der ersten handelt es sich um rein intrinsische Datentypen wie z. B. skalare Größen oder Vektoren. Diese können leicht zwischen Host-Applikation und NuT ausgetauscht und in C, C++ und FORTRAN direkt interpretiert und erzeugt werden. Die zweite Gruppe beschreibt abstrakte Datentypen, mittels welcher von der Host-Applikation komplexe Objekte in NuT referenziert werden können. Diese müssen von der Host-Applikation jedoch nicht interpretiert, sondern lediglich gespeichert und weitergegeben werden.

Damit die festgelegte Schnittstelle nicht in jeder Komponente in jeweils C, C++ und FORTRAN separat beschrieben werden muss, ist diese abstrakt in einer JSON-Datei spezifiziert. Beim Kompilieren werden daraus automatisch die passenden Implementie-

rungen generiert. Dies reduziert erheblich die Wartungsarbeiten, da das Hinzufügen oder Anpassen von Funktionen nur in einer Datei zentral und abstrakt nach dem Baukastenprinzip stattfindet. Zum anderen kann der Compiler beim Austausch von Daten zwischen unterschiedlichen Programmiersprachen generell keine Typsicherheit überprüfen. Mithilfe des eingesetzten Generatorkonzeptes wird der korrekte Datenaustausch aller Datentypen einer jeden Implementierung einmal ausführlich überprüft und kann dann beliebig oft systematisch und mit garantierter Konsistenz wiederverwendet werden.

### 3 AP2 – Verbesserung und Ausbau der Zusammenarbeit von ATHLET mit dem Numerical Toolkit

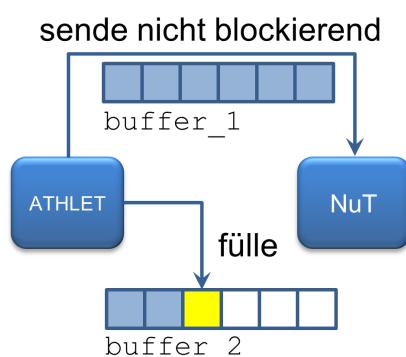
#### 3.1 Aufbau eines Doppelpuffers zur effizienteren Nutzung des ATHLET/NuT-Kommunikationskanals

Ein wesentlicher Datentransfer zwischen ATHLET und NuT tritt genau dann auf, wenn eine neue Jacobimatrix für die linearen Gleichungssysteme zur (linear-)impliziten Zeitintegration berechnet wird. Da dies über finite Differenzen geschieht, müssen die zugehörigen Funktionsauswertungen in ATHLET stattfinden, während die Matrix selbst in NuT gespeichert wird. Übertragen werden dann die nicht-trivialen Elemente der (vektorwertigen) finiten Differenzen. Dies geschieht in der Standardimplementation stets direkt nach der Berechnung jeder einzelnen dieser Differenzen.

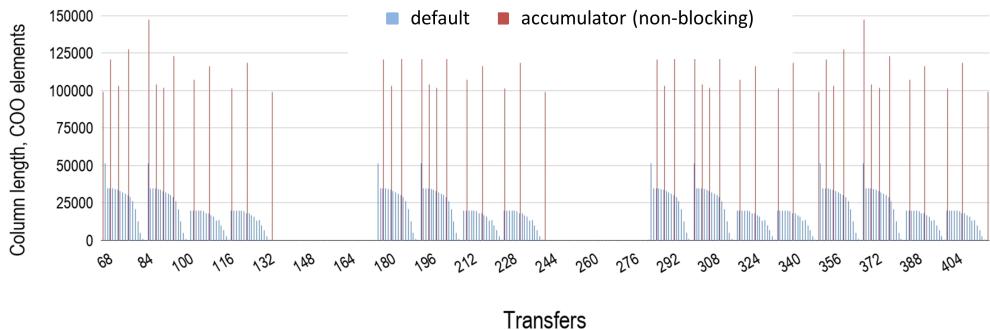
Ziel dieses Arbeitspunktes ist es, die MPI-basierte Kommunikation zwischen den beiden Programmen mittels eines Doppelpuffers effizienter zu gestalten und somit die Gesamt-Performance zu verbessern.

Ein Doppelpuffer akkumuliert in einem Puffer `buffer_1` zu transferierende Daten. Ist ein definiertes Limit erreicht, wird `buffer_1` an NuT gesendet, und in einem zweiten Puffer `buffer_2` werden die nächsten Daten akkumuliert, siehe auch Abbildung 3.1. Dieses Vorgehen wird dann im stetigen Wechsel genutzt, bis alle Daten übertragen sind. Bei entsprechender Puffergröße müssen somit weniger Übertragungen initiiert werden und die Länge der Datennachrichten orientiert sich am Normalisierungswert der Puffergröße, siehe Abbildung 3.2.

Die Arbeiten zu diesem Punkt wurden im Wesentlichen im Rahmen einer Masterarbeit, /DOR 20, Kap. 6/, durchgeführt, welche diesem Bericht angehängt ist, siehe hierzu auch Abschnitt A.1. Eine ausführliche Beschreibung des Vorgehens ist in besagter Arbeit



**Abb. 3.1** Doppelpuffer – Funktionsprinzip



**Abb. 3.2** Ausschnitt des Kommunikationsmusters zu Cube-64, Dimension: 100.657, inklusive Shifts auf horizontaler Achse zur besseren Visualisierung

zu finden. An dieser Stelle sei aber ebenfalls in breve auf die Ergebnisse, wie sie in Tabelle 3.1 zu sehen sind, eingegangen.

**Bemerkung 3.1.** Bei Cube-64 handelt es sich um ein Testproblem, in welchem das Netzwerk einen Würfel ausbildet. Das Problem ist skalierbar und bietet viele Querverbindungen zwischen den einzelnen Systemelementen.

Die hier aufgelisteten Resultate decken die wesentlichen Anwendungsszenarien ab, dass sämtliche Prozesse auf einem Node lokalisiert sind (aber ggf. auf verschiedenen Sockeln). Die dargestellten Zahlen beziehen sich auf eine nicht-blockierende Kommunikation zwischen ATHLET und NuT. Damit kann ein Puffer gesendet werden, während der andere gerade befüllt wird. Offensichtlich ist dies ein Vorteil, der durch einen einzelnen Puffer, wie in der Standardimplementation gegeben, nicht realisiert werden kann.

Obschon die Kommunikationszeit deutlich verbessert werden konnte, fällt der Anteil dieser am Gesamtaufwand zu gering aus, um nennenswerte Verkürzungen der Simulationslaufzeit zu erzielen. Somit kann das Standardvorgehen nach wie vor als effizient angesehen werden.

**Tab. 3.1** Performance des Doppelpuffers zu Cube-64, Dimension: 100.657. Siehe auch /DOR 20, S. 6.1/

	intra-socket	inter-socket
Zeitgewinn [%] (Kommunikation)	13,84	26,26
Zeitgewinn [%] (gesamt)	<0,5	<0,5

**Bemerkung 3.2.** Die durchgeführten Arbeiten am Doppelpuffer wurden anhand der FORTRAN-Implementierung von NuT durchgeführt. Das in Abschnitt 2.3 dargelegte Um-schreiben des NuT-Codes von FORTRAN nach C/C++ fand zu einem späteren Zeitraum in der Projektlaufzeit statt. Da die Ergebnisse keinen merklichen Gewinn verheißen, wurde davon abgesehen, den Doppelpuffer auch in der C/C++-Implementation zu berücksichtigen.

### **3.2 Entwicklung einer effizienten Dimensionierungsstrategie zu den linearen Gleichungssystemen**

Zu diesen Arbeitspunkt wurden Untersuchungen zweier Art vorgenommen.

Zum einen wurde analysiert, inwiefern sich gerade bei großen Problemen, also bei Problemen hoher Dimension, sinnvolle Lösungsalgorithmen zur linearen Algebra definieren lassen. Hierbei wurde insbesondere auf den Aspekt der MPI-Prozessverwaltung eingegangen.

Zum anderen wird auf die Eigenschaft des ATHLET-Differentialgleichungssystems eingegangen, nicht stets jede mögliche Gleichung auch aktiv geschaltet zu haben. Dies ergibt sich auf natürliche Weise aus dem zugrundeliegenden Netzwerk von *control volumes* und *junctions*. Findet z. B. in einem Bereich gar keine Aktivität statt oder liegt nur eine Phase vor, so sind entsprechend Gleichungen abgeschaltet. Durch die Systemdynamik in der Zeitintegration kann sich dies wiederum ändern. Gleichungen können also an- und abgeschaltet werden, und a priori ist das Wechselverhalten nicht im Detail vorhersagbar. Die Dimension des Systems ist somit zeit-variant. Protokolliert wird der Status pro Variable im Array TOP.

#### **3.2.1 Effiziente lineare Algebra**

Die Arbeiten zu diesem Unterpunkt wurden im Rahmen einer Masterarbeit angegangen, /DOR 20/, welche im Anhang dieses Berichts zu finden ist. Eine ausführliche Beschreibung zum behandelten Thema findet sich in Kapitel 5 jener Arbeit. Die dort durchgeführten Untersuchungen bestätigen, dass es sich besonders für sehr große und/oder weit verzweigte Systeme ( $10^5$  oder mehr Unbekannte) anbietet, eine Kombination aus dem direkten Löser MUMPS und dem iterativen Verfahren GMRES zu nutzen. Diese Kombination steht dem Nutzer über das Preset `mumps-gmres` zur Verfügung. MUMPS ist hierbei mit einer großzügigen Fill-in-Strategie belegt, und GMRES übernimmt unterstützend die Rolle eines Iterative Refinements. In diesem Kontext ist es lohnenswert, mehr

als einen Prozess für die numerischen Berechnungen zu reservieren, hierbei aber jene Prozesse möglichst ohne Inter-Socket- und erst recht ohne Inter-Node-Kommunikation auskommen zu lassen. Im Rahmen der bisherigen ATHLET-Modelle sind diese Nebenbedingungen ohne weiteres erfüllbar, da selbst bei großen Problemen mehr als vier dedizierte Numerik-Prozesse keinen nennenswerten Performancegewinn liefern, siehe hierzu auch /STE 17b, Unterabsch. 6.2.3/.

Entsprechend der Hardware-Konfiguration des GRS-internen Linux-Clusters können alle NuT-Prozesse (bis zu einem Maximum von zehn, s. /INT 20/) auf einem der zwei Sockel eines Nodes untergebracht werden, während ATHLET den zweiten Sockel nutzen kann, um mittels OpenMP zwischenzeitliche Prozessausdehnung zur Beschleunigung der internen Berechnungen einzusetzen.

Als weiteres Betriebssystem ist Microsoft® Windows zu berücksichtigen. Hierbei umfasst das AC<sup>2</sup>-Anwendungsszenario praktisch ausschließlich Desktop-Konfigurationen. Nach dem heutzutage üblichen Stand steht hierbei nur ein einzelner Sockel zur Verfügung. Falls ATHLET OpenMP nutzt, ist somit nur darauf zu achten, dass den NuT-Prozessen feste Kerne zugeordnet werden (nebeneinander liegend), um sicherzustellen, dass dynamisch erzeugte OpenMP-Prozesse keine Interferenz mit NuT-Prozessen erzeugen.

### 3.2.2 TOP-basierte Dimensionierungsstrategie

Das Anliegen dieses Unterarbeitspunktes ist es, eine Strategie zu entwickeln, Dimensionsänderungen der System (1.2) vorteilhaft für die lineare Algebra in NuT zu nutzen. Wie eine nähere Untersuchung des Sachverhaltes ergab, kann hierbei auf die titelgebenden Prinzipien des nächsten Abschnittes zurückgegriffen werden.

**Bemerkung 3.3.** Dimensionsänderungen des Systems implizieren unterschiedlich große lineare Systeme. Würde eine Implementierung dies direkt so umsetzen, müsste zusätzlich eine Abbildung definiert werden, die zwischen Lösungsvariablen aus der linearen Algebra und den Systemvariablen übersetzt. Um diesen Schritt zu sparen, wird im Numerical Toolkit der Kniff angewandt stets ein System gleicher Größe zu betrachten, wobei abgeschaltete Variablen dann mit einer Zeile ins Gleichungssystem eingehen, bei der nur das Diagonalelement ungleich 0 gesetzt ist. Der zugehörige Eintrag der rechten Seite steht ebenfalls auf 0. Da die lineare Algebra Inkremente für Systemvariablen berechnet, ergibt sich in einem solchen Fall sofort ein Inkrement von 0. Dies spiegelt genau das erwartete Verhalten, dass die lineare Algebra den Wert zur abgeschalteten Gleichung nicht zu verändern hat. Auch hinsichtlich einer Zerlegung der Matrix oder eines zugehörigen

Seeding-Prozesses<sup>1</sup> ergeben sich keine Nachteile, da derartige Zeilen trivial in solche Algorithmen eingehen. Das skizzierte Vorgehen ist somit legitim und spart zusätzliche Logik.

### 3.2.2.1 TOP-dynamisches und TOP-statisches Prinzip

Im Numerical Toolkit gibt es zwei Modi, mit Aktivierung und Deaktivierung von Gleichungen umzugehen. Als Standardverhalten ist eine TOP-dynamische Behandlung (TOP-D) der Jacobimatrix gesetzt. Das heißt, immer wenn sich ein Wert im TOP-Vektor ändert, wird die Jacobimatrix neu gebaut. Dieses Vorgehen ist speicherplatz-effizient, da keine expliziten Nullen gespeichert werden, birgt aber auch die Mehrarbeit der immer wieder zusätzlich durchzuführenden symbolischen LU-Zerlegung der Systemmatrix (welche gleich der Jacobimatrix inklusive eines Diagonal-Shifts ist).

**Bemerkung 3.4.** Jedes der bis dato angebotenen Solver-Presets involviert die Anwendung eines direkten Läzers für dünnbesetzte Systeme, sei es als ausschließliche Anwendung wie in `mumps` oder im Sinne eines Präkonditionierers wie in `ilu_k-gmres`. Die symbolische LU-Zerlegung ist eine der Phasen im Ablaufprotokoll eines direkten Läzers. Sie muss immer dann ausgeführt werden, wenn sich die Struktur der zu behandelnden Matrix ändert. Für eine Übersicht zu den einzelnen Phasen eines direkten Läzers für dünnbesetzte Systeme siehe /STE 20, Unterabsch. 3.3.1/ im Absatz *MUMPS / METIS*.

Gegensätzlich zum TOP-dynamischen Ansatz steht die TOP-statische Herangehensweise (TOP-S). Dieser Ansatz alloziert Speicher entsprechend der Blockstruktur der Systemmatrix und füllt Blöcke mit expliziten Nullen, sofern entsprechende Gleichungen über TOP abgeschaltet werden. Werden Gleichungen wieder hinzugeschaltet, ist kein neuer Speicher anzulegen und auch keine symbolische Faktorisierung durchzuführen. Nachteil ist, dass generell mehr Speicherplatz benötigt wird und der Fill-in weniger optimiert werden kann.

Durch einen Hysterese-Ansatz sollen die Vorteile beider Ansätze ausgenutzt werden. Dieser wurde auf Konzept-Ebene entwickelt.

**Bemerkung 3.5.** Aufgrund von Änderungen im Arbeitsplan, siehe hierzu auch Kapitel 1, sind die Arbeiten nach der Konzeptebene nicht weitergeführt worden.

---

<sup>1</sup>Zur Begriffsklärung des Seedings siehe Abschnitt 5.2.2.

### 3.2.2.2 Hysterese-Ansatz

Sind sehr viele Gleichungen aktiviert, so bietet sich ein TOP-S-Vorgehen an. Das Abschalten weniger Gleichungen lässt die Struktur nahezu invariant, was wiederum praktisch wenig Vorteile beim Fill-in oder der Speichernutzung nach sich ziehen würde, wenn dynamisch darauf reagiert werden würde. Anders herum verhält es sich bei sehr wenig aktiven Gleichungen: Das Hinzuschalten von Gleichungen geht in der Regel mit wenig Aufwand für eine symbolische Zerlegung oder einen Seeding-Prozess einher, da die Matrix sehr dünnbesetzt ist. Ein solches Szenario favorisiert somit ein TOP-D-Vorgehen. Im Bereich dazwischen, das heißt, wenn moderat viele Gleichungen aktiv sind, setzt das Hysterese-Konzept ein. Dies definiert einen Graubereich, in dem nicht direkt auf das verwandte Konzept geschlossen werden kann. Dies hängt davon ab, ob Gleichungen hinzugefügt oder abgeschaltet wurden. Die Motivation hinter diesem Vorgehen liegt darin, erst dann das Konzept zu wechseln, wenn sich hieraus eine signifikante Verbesserung ergibt. Es ist nicht sinnvoll, eine harte Grenze einzuführen, da dann das An- und Abschalten einer einzelnen Gleichung zu einer ineffizienten alternierenden Anwendung der beiden Konzepte führen könnte. Ein derartiges Hysterese-Vorgehen ist nicht unüblich, wenn es um die Entscheidung geht, zwischen zwei komplementären Konzepten zu wählen, siehe z. B. /GUS 97/ bzgl. der Behandlung von Steifheit.

Eine Implementierung bedarf fester Zahlenwerte, ab wann auf das jeweils andere Konzept gewechselt wird. Die Differenz beider Werte definiert den Hysterese-Bereich. Es bietet sich an, die relative Belegung des TOP-Arrays hierfür heranzuziehen. Sind z. B. 85% der TOP-Werte `false`, so sind viele Gleichungen aktiv und es wird auf TOP-S-Verhalten geschaltet. Sind hingegen nur z. B. 75% mit dem Wert `false` belegt, setzt TOP-D-Verhalten ein. Die Hysterese wäre dann durch das offene Intervall (75%, 85%) definiert. Grundsätzlich sollten derlei Zahlen empirisch bestimmt werden und nicht dem Nutzer zur freien Auswahl stehen. Es bedarf viel Einsicht in das numerische System, um hierzu sinnvolle Einschätzungen geben zu können.

**Bemerkung 3.6.** Es sei angemerkt, dass TOP-D und TOP-S bedingt durch die Gleitkomma-Arithmetik des Rechners unterschiedliche numerische Ergebnisse erzeugen können. Explizit gesetzte Nullen im TOP-S-Vorgehen zeichnen hierfür verantwortlich, da diese im Sinne der Struktur nicht-triviale Werte sind. Folglich können die symbolischen Zerlegungen und damit die numerischen Zerlegungen anders ausfallen. Ein Hysterese-Ansatz würde also sehr wahrscheinlich andere Ergebnisse liefern, als wenn nur eines der beiden Konzepte genutzt würde.

### 3.3 Nutzung von NuT zur Behandlung der linearen Algebra während der Startrechnung

Zur Berechnung von Anfangswerten  $y_0$  für das zu lösende System von Differentialgleichungen (1.1) während der transienten Phase, wird in ATHLET die s.g. Startrechnung durchgeführt. Diese basiert auf dem Ansatz, approximativ einen stationären Zustand für das Gesamtsystem zu ermitteln. Es hat also zu gelten  $f(y_0) \approx 0$  mit  $f$  aus (1.1). Die zugehörige Logik bedient sich algebraischer Beziehungen zwischen den Lösungsvariablen sowie iterativer Prozesse zum Auflösen implizit definierter Größen, siehe /LER 19, Absch. 8.1/. Hierbei treten lineare Gleichungssystem als Subproblem auf.

Zur Handhabung der linearen Systeme bedienen sich die Release-Versionen von ATHLET proprietärer Numerik. Zwar sind die relevanten Matrizen dünnbesetzt – die Systemvariablen interagieren direkt nur mit wenigen anderen Variablen –, jedoch wird dies nicht bei der Speicherung ausgenutzt. Es wird eine volle Matrix inklusive expliziter Nullen gespeichert. Anschließend wird *in situ* eine Gauß-Elimination mit partieller Pivotisierung durchgeführt, siehe ATHLET-Routine DIMS. Ziel dieses Arbeitspunktes ist es, das Numerical Toolkit einzusetzen, um die linearen Systeme mittels effizienter und skalierbarer Algorithmen anzugehen sowie die Organisation der Anwendungs-Numerik weiter zu vereinheitlichen.

Offensichtlich ist die Standard-Vorgehensweise nicht effizient, wenn größere Systeme betrachtet werden. Entsprechend gab es bereits in der Vergangenheit Bemühungen, einen dedizierten Löser für dünnbesetzte Systeme zum Einsatz zu bringen. Konkret handelt es sich um die Bibliothek LIS (Library of Iterative Solvers for linear systems), /NIS 20/, welche über den Plugin-Mechanismus von ATHLET angesprochen wird. Derartige Modifikationen wurden jedoch zu keinem Zeitpunkt offiziell in den Code-Kanon aufgenommen. Im Gegensatz zu PETSc ist LIS mit deutlich weniger Funktionalität ausgestattet und bietet auch nicht ein breites Spektrum an internen Schnittstellen zu weiteren numerischen Bibliotheken. Es lohnt sich also auch nicht, LIS in NuT zu integrieren. Alles, was benötigt wird, ist bereits vorhanden.

Die wesentliche Logik zur Numerik der Startrechnung ist in der ATHLET-Routine DENTH zu finden. Die Matrixdaten liegen in einer Koordinatenform vor (Zeile, Spalte, Wert). Somit muss keine neue Schnittstellenfunktion für NuT geschrieben werden, da genau dieses Format bereits in der Konstruktion der Jacobimatrix während der transienten Phase genutzt wird, um die Daten an NuT zu senden. Im Zuge der Arbeiten wurden in DENTH lediglich drei zusätzliche Arbeitsarrays eingeführt, um erst alle Werte und Positionsinformationen zu sammeln und dann in einem einzigen Aufruf an NuT zu senden.

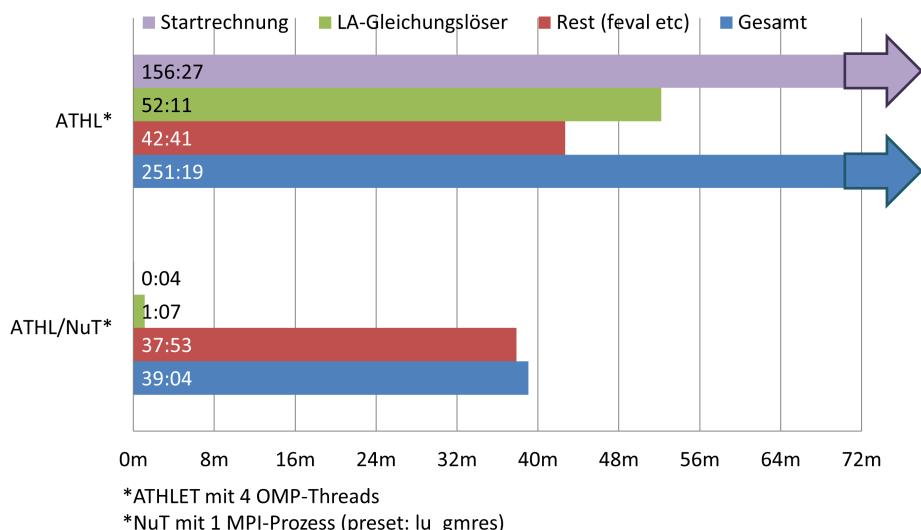
Die Zusammenarbeit mit NuT wird über eine neue Entität geregelt, welche in der Datei nut.F90 über den Aufruf

```
call linalg_new(linalg_start_entity, "athlet-nut", &
                "athlet-linalg-start", shared=0)
```

in NuT erstellt wird. Hierbei bezeichnet athlet-linalg-start die Entitäten-ID-String und athlet-nut den genutzten MPI-Kommunikator. ATHLET kann dann über das Handle linalg\_start\_entity mit der NuT-seitigen Entität interagieren. Diese neue Entität arbeitet auf dem gleichen Kommunikator und nutzt das gleiche Solver-Preset wie die Haupt-Entität zur linearen Algebra im FEBE-Kontext. Bei Bedarf kann zu einem späteren Zeitpunkt leicht angepasst werden.

## Performance

Gerade bei großen und komplexen Systemen ist es recht wahrscheinlich, signifikante Verbesserungen durch die NuT-Modifikationen zu erreichen, siehe Abbildung 3.3. Teils mag der Vorteil auch eher moderat ausfallen und in der Gesamtperformance weniger Gewicht haben. Vorteil ist aber stets, dass vermittels NuT die relevante Numerik basierend auf modernen Algorithmen zentral koordiniert und gepflegt werden kann.



**Abb. 3.3** Performance-Vergleich zum Datensatz SPX (Modellierung des Natrium gekühlten Brutreaktors Superphénix), Dimension: 106.757

### **3.4 Weiterentwicklung und Ausbau der Differentialgleichungsnumerik im Numerical Toolkit**

Wie bereits in Kapitel 1 beschrieben steht im Kern der Berechnungen von ATHLET das Lösen des Anfangswertproblems (1.1), i. e.

$$y' = f(t, y), \quad y(t_0) = y_0. \quad (3.1)$$

Die Funktion  $f$  beschreibt hierbei die Systemdynamik und wird durch die ATHLET-Routine AFK repräsentiert, welche wiederum Unterroutinen zur Abbildung der Modell-Logik aufruft. Ziel der Arbeiten zu diesem Arbeitspunkt ist es, eine geeignete Methodenlogik ins Numerical Toolkit auszulagern und in Verbindung mit einer ATHLET-seitigen Steuerlogik die Zeitintegration zu (3.1) durchzuführen. Die geleisteten Arbeiten sind theoretischer und konzeptueller Natur. Konkrete Implementierungsarbeiten stehen noch aus.

Die Zeitintegration in ATHLET geht mit speziellen Anforderungen einher. Nachfolgend werden diese beschrieben. Die konzeptuellen Ausarbeitungen im Rahmen einer neuen Methodenlogik werden in Abschnitt 3.4.2.3 bis Abschnitt 3.4.2.6 vorgestellt. Zu Vergleichszwecken wird auch auf das aktuelle Vorgehen durch die ATHLET-eigene Integrationsroutine FEBE eingegangen.

Spezielle softwaretechnische Belange und erarbeitete Konzeptideen hierzu werden in Abschnitt 3.4.3 erörtert.

#### **3.4.1 ATHLET-spezifische Anforderungen an die Zeitintegration**

Für die numerische Behandlung von (3.1) sind verschiedene Eigenschaften des Systems für eine Steuer- und Methodenlogik zu beachten.

- a) Aufgrund der Modell-Logik kann es in  $f$  zu Unstetigkeiten kommen. Hierzu existiert eine gesonderte Steuerlogik, welche über das Flag HXX charakterisiert wird.
- b) Eine  $f$ -Auswertung während der lokalen Zeitintegration kann relativ teuer im Vergleich zum Lösen eines linearen Gleichungssystems sein (Zerlegung der involvierten Matrix bereits vorausgesetzt).
- c) Aufgrund von nicht relevanten aber möglicherweise auftretenden Mikroschwingungen im System bedarf es keiner Methoden hoher Ordnung. Ebenso ist ein großes Stabilitätsgebiet des Zeitintegrators von Vorteil, da hiermit ein dämpfendes Verhalten einhergeht und dies wiederum einen glättenden Effekt auf die berechnete Lösung hat.

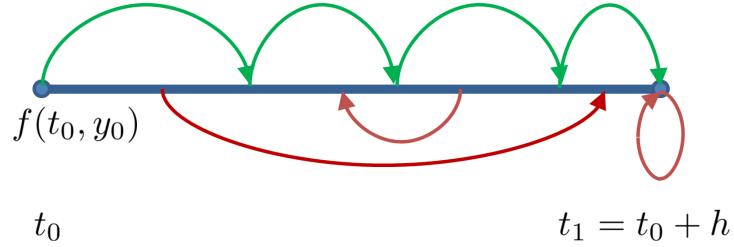
Zur genauen Definition des Stabilitätsgebietes siehe z. B. /STE 17b, Unterabsch. 3.2.2/ oder die Standardliteratur /HAI 96; BUT 16/.

- d) Teilergebnisse aus der  $f$ -Auswertung werden für ein Zusatzsystem bereitgestellt, welches zur Massenkorrektur genutzt werden kann, siehe hierzu /AUS 19, Gl.2-10 und 2-11/.
- e) Für das Ergebnis einer  $f$ -Auswertungen im lokalen Zeitintegrationsintervall kann es relevant sein, zu welcher Subzeitschrittweite die ggf. vorherige  $f$ -Auswertung stattgefunden hat. Dies liegt an einer AFK-internen losen Kopplung zum HECU-Modul (Wärmetransport), da hierdurch ein Memory-Effekt erzeugt wird. Faustregel ist, sich mit einem Auswertungspfad stets in Richtung des positiven Zeitverlaufs durch das Intervall zu bewegen. Hierbei dürfen die Subschritte nicht zu groß werden, da sonst der HECU-Einfluss nicht hinreichend Berücksichtigung findet. Als Grenzwert kann hier eine relative Subschrittweite von  $1/2$  gesehen werden, da es sich auch im ATHLET-eigenen Zeitintegrator FEBE um die maximale relative Subschrittweite handelt. Neben der Anforderung, nicht zu große Subschritte zu gehen, ist darauf zu achten, nicht am gleichen Zeitpunkt mehrmals hintereinander auszuwerten. Rückwärtsschritte werden ebenfalls nicht unterstützt. Stattdessen ist wieder vom Anfang des lokalen Intervales loszugehen. Siehe auch Abbildung 3.4 für eine Veranschaulichung der Verhältnisse.
- f) Ist für den lokalen Zeitschritt eine Approximation  $y_1$  berechnet, wird abschließend  $f(t_1, y_1)$  ausgewertet, um sicherzustellen, dass das System valide bleibt, bevor der Schritt beendet wird. Dies sollte möglichst kostengünstig möglich sein.

**Bemerkung 3.7.** Aufgrund von den in Punkt e) beschriebenen Eigenschaften von  $f$  handelt es sich im klassischen Sinne nicht um eine mathematische Funktion. Ein konsistentes Verhalten von  $f$  lässt sich jedoch erzwingen, wenn man für Auswertungen an einer festen Stelle stets den gleichen Pfad zugrunde legt. Angestrebte Funktionsauswertungen weiter rechts im lokalen Zeitintervall sind somit auf Hilfsfunktionsauswertungen angewiesen. Im Sinne der Effizienz sollten diese vom Zeitintegrator bereits für andere Berechnungen sinnvoll genutzt worden sein.

### 3.4.1.1 FEBE-Integrationsansatz

Zur Zeitintegration wird in ATHLET die Routine FEBE (Forward Euler / Backward Euler) genutzt, in der ein Extrapolationsansatz auf Basis eines Zusammenspiels aus explizitem und linear-impliziten Euler-Verfahren implementiert ist. In der Praxis ist de facto nur



**Abb. 3.4** Beispielhafter valider  $f$ -Auswertungspfad (grün) sowie invalide Pfadabschnitte (rot)

die linear-implizite Handhabung von Relevanz, siehe auch /LER 19, Absch. 1.4/. Die Extrapolation wird bis zu der Ordnung Drei durchgeführt. Die zu lösenden linearen Systeme ergeben sich zu

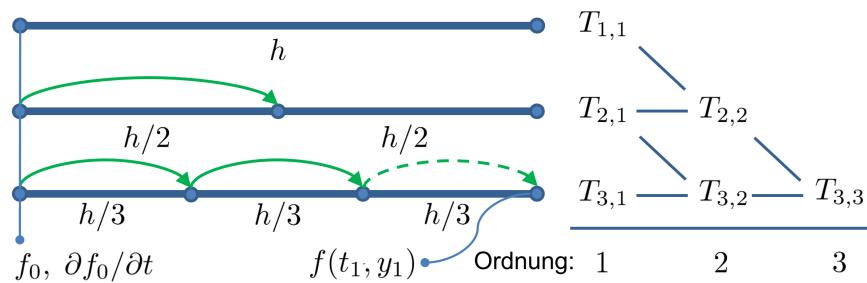
$$(I - h_i J)\delta y_{ij} = h_i f\left(t_0 + j \cdot h_i, y_0 + \sum_{\ell=0}^{j-1} \delta y_{i\ell}\right) + h_i^2 \frac{\partial f_0}{\partial t} \quad (3.2)$$

$$j = 0, \dots, i-1, \quad h_i = h/i, \quad i = 1, 2, 3.$$

Hierbei ist  $J$  eine Approximation für  $\partial f_0 / \partial y$ . Die Startwerte zur Extrapolation, Abbildung 3.5, sind dann gegeben über

$$T_{i,1} := y_0 + \sum_{j=1}^i \delta y_{ij}, \quad i = 1, 2, 3. \quad (3.3)$$

Es zeigt sich, dass der FEBE-Extrapolationsansatz gut auf obige Anforderungen und Eigenschaften abgestimmt ist. Das Verfahren ist einfach nachzuvollziehen, bedarf weniger Funktionsauswertungen und im gegebenen Kontext sind insbesondere das dämpfende Verhalten des linear-impliziten Euler-Verfahrens sowie die abzulaufenden Pfade zur Funktionsauswertung von Vorteil. Zum letzten Punkt siehe auch Abbildung 3.5.



**Abb. 3.5**  $f$ -Auswertungspfade sowie Extrapolationsschema zum FEBE-Integrator

Neben den Vorteilen des Verfahrens gibt es aber auch Nachteile:

- FEBE nutzt das Prinzip der lokalen Extrapolation: Es wird der lokale Fehler von  $T_{3,2}$  per Differenz mit  $T_{3,3}$  geschätzt, jedoch mit  $T_{3,3}$ , der Approximation höherer Ordnung, die Integration fortgeführt, d. h.  $y_1 = T_{3,3}$ .<sup>2</sup> Diese Approximation ist nicht A-stabil, was zu Problemen mit oszillatorischer Steifheit führen kann. Siehe /STE 17b, Absch. 3.2/ für Details.
- Für jedes der drei zu berechnenden  $T_{i,1}$  ist eine neue Matrix für die zugehörigen linearen Systeme (3.2) zu zerlegen, da der Diagonal-Shift von  $i$  abhängig ist. Dies kann mitunter sehr zeitaufwendig sein.
- In /STE 17b, Absch. 3.2/ wird gezeigt, dass der FEBE-Extrapolationsansatz als Rosenbrock-Wanner-Methode (ROW) interpretiert werden kann. Die asymptotische Konsistenzordnung der Approximationen  $T_{i,j}$  hängt also nicht von  $J$  ab. Dies gilt jedoch nicht für die Fehlerschätzung basierend auf der Differenz  $T_{3,3} - T_{3,2}$ . Dies folgt sofort aus /HAI 96, Gl. IV.7.32/, da nicht alle Bedingungen zur Negierung des Einflusses von  $J$  für  $T_{3,2}$  erfüllt sind.
- Zur Bestimmung der Approximation  $y_1$  steht an  $t_0 + h$  (voller Schritt) nur eine  $f$ -Information niedriger Ordnung zur Verfügung. Dies muss nicht zwingend nachteilig sein. Die Implizitheit des Verfahrens wird dadurch aber geschwächt und dies kann sich negativ auf die Stabilität auswirken.
- Der Extrapolationsansatz ist in der gegebenen Form nicht steif-akkurat und kann darüber hinaus unter Ordnungsreduktion leiden. Dies ist ausführlich in /STE 17b, Unterabsch. 3.2.2/ untersucht worden. Entsprechend unattraktiv zeigt sich das Abschneiden bei den Tests in /STE 17b, Unterabsch. 3.5.2/ zu klassischen steifen Problemen.

### 3.4.2 Der FiterRK-Ansatz

Insbesondere die in Bemerkung 3.7 diskutierte Pfadbedingung zu  $f$ -Auswertungen, erschwert es, vollimplizite Verfahren zum Einsatz zu bringen. Um dennoch den Nachteilen des FEBE-Extrapolationsansatzes entgegenwirken zu können, wurde in /STE 17b, Absch. 3.3/ das Konzept der FiterRK-Methode entwickelt und mit den Arbeiten zu diesem Arbeitspunkt verfeinert.

Grundsätzlich zeichnen sich FiterRK-Methoden dadurch aus, dass eine a priori bekannte und feste Anzahl an Newton-Iterationsschritte pro Stufe ausgeführt wird. Somit gehören auch ROW-Methoden, und damit der FEBE-Extrapolationsansatz, zu dieser Klasse, da

---

<sup>2</sup>Es ist Zufall, dass der Begriff *Extrapolation* in einer zweiten Semantik auftaucht. Hierbei handelt es sich um die übliche Bezeichnung für das beschriebene Vorgehen, siehe z. B. /HAI 96/.

genau ein Newtoninkrement pro Stufe berechnet wird. Die FEBE-Extrapolation basiert auf einem diagonal-impliziten Runge–Kutta-Verfahren (DIRK), welches per definitionem keine hohe Stufenordnung aufweisen kann und zusätzlich nicht steif-akkurat ist. Genau diese Nachteile erbt FEBE.

Im Sinne einer Verbesserung des Methodenansatzes wurden als Basis für die Untersuchung und Entwicklung von FiterRK-Methoden  $s$ -stufige steif-akkurate Runge–Kutta Methoden betrachtet, i. e.

$$Y_i = y_0 + \sum_{j=1}^s a_{ij} k_j \quad \text{mit} \quad k_i = h f(t_0 + c_i h, Y_i) \quad \text{und} \quad \begin{array}{l} c_1, \dots, c_s \text{ verschieden,} \\ c_s = 1, \quad y_1 = Y_s. \end{array} \quad (3.4)$$

welche zusätzlich von hoher Stufenordnung sind und ein Ein-Punkt-Spektrum vorweisen.

### 3.4.2.1 Stufenordnung und Ein-Punkt-Spektrum

Seien

$$\mathcal{A} := (a_{ij})_{i,j=1,\dots,s}, \quad C := \text{diag}(c_1, \dots, c_s) \text{ sowie } e := (1, \dots, 1)^T \in \mathbb{R}^s.$$

Gelte für natürliches  $p > 0$  die Bedingung

$$\mathcal{C}(p) : \quad \mathcal{A}C^{q-1}e = \frac{1}{q}C^q e, \quad q = 1, \dots, p. \quad (3.5)$$

Dann ergibt sich

$$\mathcal{C}(p) \implies \begin{cases} y(t_0 + c_i h) - Y_i \in \mathcal{O}(h^{p+1}), & i = 1, \dots, s, \quad \text{siehe /BUT 64/} \\ \mathcal{A} \text{ enthält einen Block } A_m \text{ der Größe } m \text{ mit } m \geq p \end{cases}$$

$\mathcal{C}(p)$  wird aus offensichtlichen Gründen Stufenordnungsbedingung (zur Ordnung  $p$ ) genannt. Das erfüllen dieser Bedingung schützt eine Methode der Ordnung  $p$  vor Ordnungsverlust bei der Integration steifer Systeme, siehe /HUN 03, Kap. II.2/. Zusätzlich ergibt sich der Vorteil, dass es gerechtfertigt ist, bereits die Stufenwerte  $Y_i$  für Aussagen über das System an den Punkten  $t_0 + c_i h$  heranzuziehen.

Der Preis für diese vorteilhaften Eigenschaften liegt in einer höheren Methoden-Komplexität, da die die Methode einen Block von  $m \geq p$  gegenseitig abhängigen Stufen enthalten muss. Sei  $n$  die Dimension des Problems. Um nicht einen Newton-Prozess der Dimension  $n \cdot m$  angehen zu müssen, wird  $A_m$  einer Ähnlichkeitstransformation zu einer Dreiecksmatrix unterzogen. Entsprechend ergeben sich dann  $m$  verzahnte Newton-Prozesse der Dimension  $n$ , deren Iterationsschritte sukzessive ineinandergreifen, siehe /HAI 96, Kap. IV.8/. Besagte Dreiecksmatrix enthält die Eigenwerte von  $A_m$  auf der Diagonalen.

Diese gehen wiederum direkt als Faktor  $\gamma$  in den Diagonal-Shift für die Jacobimatrix-approximation  $J$  in den zugehörigen linearen Gleichungssystemen (1.2) ein. Es ist somit erstrebenswert, ein Ein-Punkt-Spektrum für  $\mathcal{A}$  und damit  $A_m$  sicherzustellen.

$$(J - (h\gamma)^{-1}I)\Delta Y_i = b_i, \quad (3.6)$$

Wird  $m = p$  gewählt, ist dies gleichbedeutend mit der Anforderung an die  $c_i$  skalierte Nullstellen von Laguerre-Polynomen entsprechenden Grades zu sein, siehe /HAI 96, Lemma IV.8.1/. Mehr Flexibilität ergibt sich für  $m = p + 1$  wie in /STE 17a/ ausgeführt.

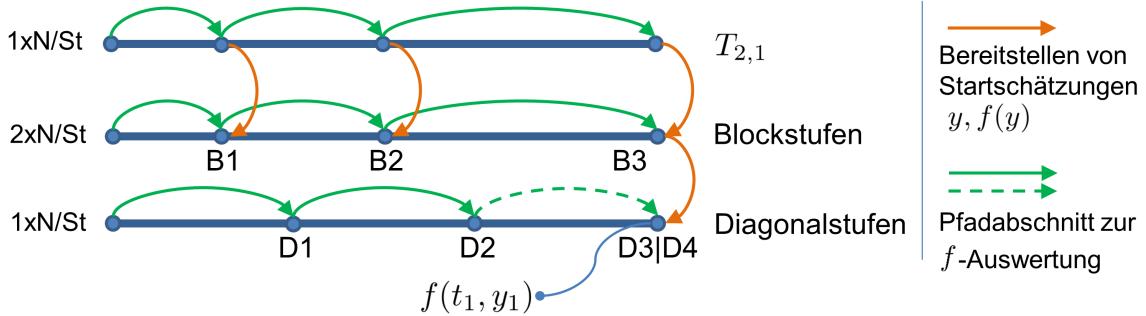
### 3.4.2.2 FiterRK durch Newton

Per Newton-Prozess werden die  $k_i$  der Basismethode (3.4) iterativ bestimmt. Die Idee ist, die Approximationseigenschaften des Newton-Prozesses auszunutzen, um nach einer festen Anzahl an Schritten die Iteration abbrechen zu können und dennoch im Sinne der Konsistenzordnung der Methode hinreichend gute Approximationen  $k_i^\ell$  zur Verfügung zu haben. Legitimiert wird dieses Vorgehen durch /STE 17b, Satz 3.26/.

### 3.4.2.3 FiterRK32-Prototyp

Die Kombination aus Auswertungspfaden zu  $f$  und dem Anspruch  $f(t_1, y_1)$  günstig berechnen zu können, ist eine sehr spezielle Anforderung des ATHLET-Kontextes an die Entwicklung einer neuen Methode. Es wurde ein A-stabiler FiterRK-Prototyp der (Stufen-)Ordnung  $p = 2$  erstellt, der neben einem  $y_1$  der Ordnung  $p = 2$  zum vollen Schritt auch ein  $\tilde{y}_1$  der Ordnung  $p + 1 = 3$  zur lokalen Fehlerschätzung und Extrapolation berechnet. Die  $f$ -Auswertungspfade sind in Abbildung 3.6 dargestellt. Hier sind auch s.g. *Stage-Mappings* aufgezeichnet, welche günstig Startapproximation für Stufenwerte und zugehörige Funktionsauswertungen liefern. Dies spart  $f$ -Auswertungen. Die Anzahl der Blockstufen ist hier zu  $m = p + 1 = 3$  gesetzt. Ein Ein-Punkt-Spektrum wird über den Einsatz der Techniken aus /STE 17a/ sichergestellt. Der Eigenwert ist zu  $\gamma = 2$  gewählt, wodurch ein dämpfendes Verhalten forciert wird. Zusätzlich zu den Blockstufen sind vier Diagonalstufen definiert. Ein  $k_i$  zu einer Diagonalstufe hängt nur von sich selbst und vorherigen  $k_j$  ab. Damit können Diagonalstufen sukzessive berechnet werden. Wie bereits in /BUT 90/ gezeigt, liefern Diagonalstufen mehr Freiheitsgrade, um z. B. die Stabilität oder Fehlerkonstante  $C_{err}$  des Verfahrens verbessern zu können. Um auch für  $\mathcal{A}$  ein Ein-Punkt-Spektrum zu sichern gilt  $a_{ii} = \gamma$  für alle Diagonalstufen. Somit wird genau nur eine Matrix für die linearen Systeme (3.6) benötigt.

Der FiterRK-Charakter zeigt sich linksseitig in Abbildung 3.6, wo die Anzahl an Newtonschritten pro Stufe (N/St) angegeben sind. Um eine gute Startnäherung für die Blockgrößen zu haben, wird  $T_{2,1}$  wie für die FEBE-Extrapolation berechnet. Aufgrund von  $\gamma = 2$  wird auch hier die gleiche Matrix für die linearen Systeme (3.2) genutzt wie für die Hauptmethode.



**Abb. 3.6**  $f$ -Auswertungspfade zum FiterRK32-Prototyp sowie Bereitstellung von Startwerten, B1-B3 beziehen sich auf Blockstufen, während D1-D4 Diagonalstufen bezeichnen

### 3.4.2.4 Fehlerschätzung über Approximationsgüte der Stufen

Im Kontext der Beispielmethode der Stufenordnung  $p = 2$  bleibend kann der lokale Fehler über die Differenz von  $y_1$  und  $\tilde{y}_1$  gemessen werden:

$$y(t_0 + h; y_0) - y_1 = \tilde{y}_1 - y_1 + \mathcal{O}(h^4).$$

Dies gilt sowohl für die implizite Basismethode (3.4) als auch für die FiterRK-Adaption. Eine Idee aus /CHE 98/ aufgreifend kann der Fehler alternativ unter Ausnutzung der Stufenordnung geschätzt werden. Nach Definition der Fehlerkonstanten  $C_{err}$  gilt für den lokalen Fehler die Darstellung

$$y(t_0 + h; y_0) - y_1 = C_{err} \cdot h^3 y_0^{(3)} + \mathcal{O}(h^4), \quad (3.7)$$

wobei  $y_0^{(3)}$  die dritte Ableitung von  $y$  nach  $t$  an  $t_0$  bezeichnet. Diese Ableitung lässt sich für drei beliebige  $k_i$  aus (3.4) aufgrund von

$$\begin{pmatrix} k_{i_1} \\ k_{i_2} \\ k_{i_3} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & c_{i_1} & c_{i_1}^2/2 \\ 1 & c_{i_2} & c_{i_2}^2/2 \\ 1 & c_{i_3} & c_{i_3}^2/2 \end{pmatrix}}_{=:W} \begin{pmatrix} h y_0^{(1)} \\ h^2 y_0^{(2)} \\ h^3 y_0^{(3)} \end{pmatrix} + \mathcal{O}(h^4). \quad (3.8)$$

schätzen. Hierbei sind  $y_0^{(1)}$  und  $y_0^{(2)}$  analog zu  $y_0^{(3)}$  definiert. Seien  $c_{i_1}$ ,  $c_{i_2}$ ,  $c_{i_3}$  paarweise verschieden. Dann ist  $W$  regulär, da es sich um eine gewichtete Vandermonde-Matrix handelt. Folglich ergibt sich

$$y(t_0 + h; y_0) - y_1 = C_{err} \cdot (0, 0, 1) W^{-1} \begin{pmatrix} k_{i_1} \\ k_{i_2} \\ k_{i_3} \end{pmatrix} + \mathcal{O}(h^4). \quad (3.9)$$

Die in Abbildung 3.6 angegebene Anzahl an Newtonschritten pro Stufe der FiterRK-Adaption ist hinreichend gewählt, so dass die Beziehung (3.8) ebenso für die approximativen  $k_i^{l_i}$  gilt und damit auch (3.9). Weder für die Fehlerschätzung über (3.7) noch über (3.9) hat die Jacobimatrixapproximation  $J$  in der Asymptotik einen Einfluss.

### 3.4.2.5 Handhabung des Zusatzsystems zur Massenkorrektur

Das in Abschnitt 3.4.1 erläuterte Zusatzsystem zur Massenkorrektur wird in ATHLET über den gleichen Extrapolationsansatz wie das Hauptsystem berechnet, jedoch wird dabei auf die Teil- $f$ -Auswertungen aus dem linear-impliziten Vorgehen zurückgegriffen. Damit ist die Zeitintegration im Zusatzsystem explizit und die Euler-Schritte können als Anwendung einer Quadraturformel interpretiert werden. Anschließende Extrapolation sorgt für die entsprechende Ordnung.

Im FiterRK-Kontext bedarf die Berücksichtigung des Zusatzsystems ein wenig mehr Überlegung, da das Hauptsystem ggf. mehr als einen Newton-Iterationsschritt pro Stufe ausführt und die Abfolge der Stufen nicht wie beim FEBE-Ansatz als Anwendung eines Grundverfahrens für verschiedene Subschrittweiten interpretiert werden kann. Auf der anderen Seite kann das Vorgehen derart gesehen werden, dass die berechneten  $k_i$  lediglich den Input für die Linearkombination zur Berechnung des  $y_1$  liefern. Da  $y_1$  ebenfalls Stufewert ist, siehe (3.4), folgt aus (3.5), dass die zugehörigen Koeffizienten  $a_{sj}$  Gewichte einer Quadraturformel der Ordnung  $p = 2$  sind. Somit muss nur gewartet werden, bis die finalen  $f$ -Auswertungen pro Stufe durchgeführt wurden, um anschließend die durch die  $a_{sj}$  charakterisierte Quadraturformel mit dem Faktor  $h$  auf die Teil- $f$ -Auswertungen anzuwenden. Man beachte, dass hierbei die Teil- $f$ -Auswertungen zur Anwendung kommen und nicht die  $k_i$  selbst. Das ist aber nicht nachteilig, da bereits für die finalen  $f$ -Auswertungen hinreichend gute Approximationen zum Einsatz kommen.

### 3.4.2.6 Implizite Null

Als Verbesserung zum Grundkonzept des FiterRK-Ansatzes wurde das Konzept der impliziten Null entwickelt. Die Untersuchungen hierzu sind noch nicht abgeschlossen. Die Idee ist es, statt einer klassischen expliziten Nullstufe, i. e.,  $c_0 = 0$  sowie  $k_0 = f(t_0, y_0)$  eine oder vielmehr mehrere implizite Nullstufen zu nutzen. Eine explizite Nullstufe ist einfach zu generieren und bietet der Methode mehr Freiheitsgrade, macht es auf der anderen Seite aber schwieriger die Methode  $A$ -stabil zu halten.

Eine implizite Nullstufe wird wie eine übliche Diagonalstufe der Stufenordnung  $p$  definiert, setzt hierbei aber die relative Schrittweite  $c_N = 0$  an. Als Startapproximation für den Stufenwert  $Y_N$  bietet sich sofort  $Y_N^0 = y_0$  an, was zu der Funktionsauswertung  $f(t_0, y_0)$  führt. Diese ist im ATHLET-Kontext bereits vorhanden, da sie Teil der Validitätsprüfung zu  $f$  ist. Aufgrund der guten Startnäherung bedarf es nur einer Newtoniteration, um die implizite Nullstufe von hinreichender Qualität vorliegen zu haben. Die Berechnung ist also sehr günstig.

Die bisherigen Untersuchungen ergaben, dass eine hervorragende Anwendung für eine implizite Nullstufe die Bestimmung einer Startnäherung für eine Diagonalstufe ist. Gute Startnäherungen sind wichtig im FiterRK-Kontext, um die Anzahl an Newton-Iterationsschritten, und damit die Anzahl an  $f$ -Auswertungen, niedrig zu halten. Es kann gezeigt werden, dass mithilfe der impliziten Nullstufe und im gegebenem Kontext von Methoden hoher Stufenordnung stets eine Startnäherung konstruiert werden kann, welche das gleiche lineare Stabilitätsverhalten wie die Stufe selbst zeigt. Dies erleichtert die Konstruktion einer Methode erheblich, da sich auf die Stabilitäts-eigenschaften der einzelnen Stufen konzentriert werden kann. Anschließend wird pro Diagonalstufe die dedizierte Nullstufe ermittelt und hierüber die Startnäherung.

Als weitere mögliche Anwendungsgebiete ergeben sich folgende Themen. Die Untersuchungen hierzu laufen jedoch noch.

- Abschätzung zur Güte der Jacobimatrixapproximation  $J$ . Dies könnte als weiterer Aspekt in der Steuerlogik zum Update von  $J$  Anwendung finden.
- Dingfest machen von Störungen in höheren Ableitungen der  $k_i^\ell$ . Dies könnte für eine adaptive Ordnungskontrolle genutzt werden. Für ein Verfahren der (Stufen-)Ordnung  $p$  bräuchte es den Zugriff auf  $y_0^{(p+2)}$ . Das ist bereits im Kontext impliziter Methoden nicht-trivial, siehe /BUT 98/. Durch den FiterRK-Aspekt kommt eine weitere Komplexitätsstufe hinzu, die es zu beachten gilt.

### 3.4.2.7 Aufwand und Nutzen

Die Abbildungen 3.5 und 3.6 lassen es bereits vermuten, dass ein Schritt des FiterRK-Prototyps mehr  $f$ -Auswertungen und zu lösende lineare Gleichungssysteme benötigt als der FEBE-Extrapolationsansatz. Auf der anderen Seite zeigt die Methode seitens der Theorie deutlich besseres Ordnungs- und Stabilitätsverhalten. Des Weiteren bedarf es lediglich der Handhabung einer einzelnen Matrix für die linearen Systeme (3.6) im Gegensatz zu den drei Matrizen in (3.2). Es hängt vom konkreten Problem (3.1) ab, wie teuer Matrixzerlegungen im Vergleich zu  $f$ -Auswertungen sind. A priori ist dies schwierig abzuschätzen. Es bedarf einer konkreten Implementation, um dies auszutesten. Dies sollte in einer etwaigen Weiterführung des Themas der nächste Schritt sein.

Seitens der Stabilität ist deutlich besseres Verhalten zu erwarten. Auch kann noch an Details gefeilt werden, weswegen die konstruierte Methode als Prototyp gilt. Z. B. könnte die Berechnung von  $\tilde{y}_1$  eingespart und der Fehler über (3.9) geschätzt werden. Ebenso gilt es zu untersuchen, ob durch die impliziten Nullstufen Newtoniterationen oder Diagonalstufen eingespart werden können. Der FiterRK-Ansatz bietet auf alle Fälle viel Potential.

### 3.4.3 Anforderung an die technische Umsetzung

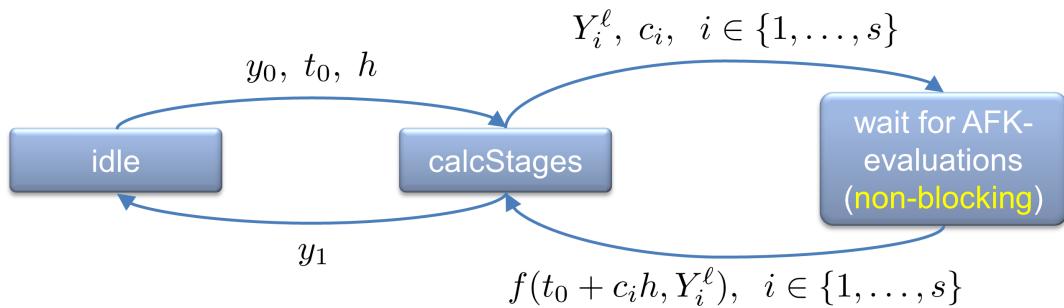
Um eine Methoden-Logik zu ROW- und FiterRK-Methoden im Numerical Toolkit zu ermöglichen, bedarf es zusätzlicher Schnittstellenfunktionen, die auf den Austausch von Information zu  $f$ -Auswertungen abzielen. Nach Übergabe der Startwerte  $y_0$ ,  $t_0$ ,  $h$  generiert eine Methode approximative Stufenwerte  $Y_i^\ell$ , für die  $f$ -Auswertungen benötigt werden. Da dies einen Aufruf der ATHLET-seitigen AFK-Routine impliziert, muss die entsprechende Information aus dem Numerical Toolkit an ATHLET gesendet werden. Nach der AFK-Auswertung wird das Ergebnis zurückgesandt. Es ergibt sich also ein deutlich komplexeres Kommunikationsverhalten im Vergleich zur NuT-Funktionalität der linearen Algebra, welche bei gespeichertem  $J$  einen Diagonal-Shift und eine rechte Seite empfängt und dann ohne weitere Zwischenkommunikation die numerische Lösung liefert.

Da das Numerical Toolkit mehreren Host-Applikationen zur Verfügung stehen soll, siehe auch Abschnitt 4.3, ist es nicht sinnvoll, das Warten auf die AFK-Auswertungen blockierend zu gestalten. Dies macht es notwendig, Status und zugehörige Übergänge innerhalb des Numerical Toolkits bezüglich laufender Tätigkeiten zur Differentialgleichungsnumerik zu definieren. Ein sinnvolles Mittel hierfür ist das Konzept der s.g. *State-Machine*. Grundlegende zu berücksichtigende Zustände und Übergänge finden sich in Abbildung 3.7. Eine

*State-Machine* sollte ausgearbeitet sein, bevor es an die Implementation der eigentlichen Methoden-Logik geht, damit diese direkt darauf abgestimmt werden kann.

Neben der State-Machine und der Methoden-Logik ist auf folgende Besonderheiten einzugehen, welche sich aus obiger Diskussion ergeben haben:

- Logik zum *Stage-Mapping* – Abschnitt 3.4.2.3,
- Handhabung des Zusatzsystems – Abschnitt 3.4.2.5,
- Einbinden impliziter Nullstufen – Abschnitt 3.4.2.6.



**Abb. 3.7** Grundlegende Zustände und Übergänge in NuT, welche durch eine *State-Machine* abgebildet werden müssen

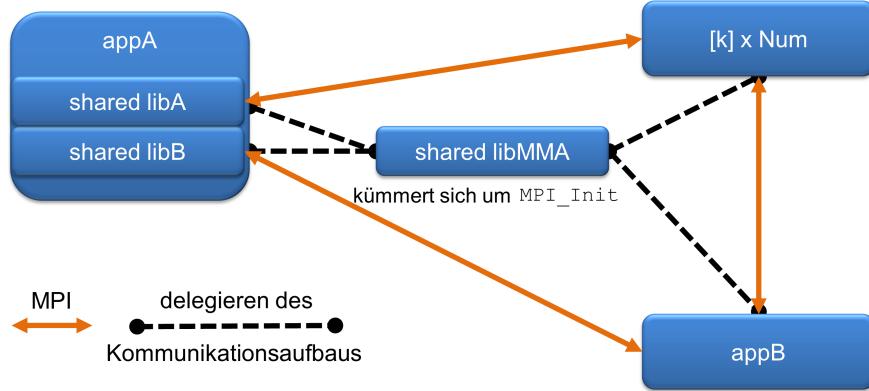
### 3.5 Einbringen der Kommunikations-Bibliothek MMA

MMA steht für *MPI for Multiple Applications*; hierbei handelt es sich um eine Kommunikations-Bibliothek, welche dem Nutzer ein einfach zu bedienendes Interface liefert, um transparent MPI-Prozesse verschiedenen Kommunikatoren zuzuordnen und diese zu nutzen. MMA ist eine Open-Source-Bibliothek und ist über GitLab verfügbar, /JAC 20/.

Aufgrund der Arbeiten zu diesem Punkt steht MMA nun im AC<sup>2</sup>-Kontext zur Verfügung. Die Motivation, MMA einzubinden, liegt darin, einen gekapselten und anwendungsunabhängigen Mechanismus an der Hand zu haben, welcher die Initialisierung der Kommunikation zwischen heterogenen MPI-Prozessen übernehmen kann, siehe Abbildung 3.8. Im Rahmen von AC<sup>2</sup> ist dies kein triviales Thema, da z. B. auf geteilte Ressourcen wie die NuT-Prozesse zugegriffen werden können soll. Siehe hierzu auch Abschnitt 4.3.

Neben der eigentlichen Funktionalität bietet MMA weitere allgemeine Vorteile:

- leichtgewichtig und als dynamische Bibliothek einbindbar,
- geschrieben in C (ISO-C-Standard) mit zusätzlicher Schnittstelle für FORTRAN,



**Abb. 3.8** Delegieren der Kommunikationsinitialisierung an MMA

- CMake-Unterstützung (Version 3.8 oder höher) für Builds unter Linux und MS Windows,
- 2-Klausel-BSD-Lizenz, /INI 20/, und damit kompatibel zum GRS-Lizenzmodell.

Speziell im Kontext von ATHLET ist darüber hinaus folgender Vorzug zu verzeichnen:

#### Prozess-interne dedizierte Kommunikatoren

ATHLET nutzt die Open-Source-Bibliothek *Fortran Development Extensions (libfde)*, siehe /SCH 20/, um zusätzliche Funktionalität zur Laufzeit über ein Plugin-Mechanismus zur Verfügung stellen zu können. Hierüber wird z. B. bei Bedarf die NuT-Schnittstelle oder ATHLET-CD geladen. ATHLET als Hauptapplikation sowie etwaige Plugins laufen im gleichen Prozess-Adressraum. Ebenso gilt dies für MMA, sobald es geladen wird. Der Vorteil hiervon ist, dass sich ATHLET wie auch die Plugins in der gleichen Liste für Kommunikatoren bei MMA registrieren, da diese prozess-global ist. Somit besteht grundsätzlich die Möglichkeit, dass z. B. ein Plugin wie ATHLET-CD einen dedizierten Kommunikator mit NuT für interne Berechnungen nutzen könnte. Die Kommunikation erfolgte direkt; die übergeordnete ATHLET-Instanz wäre nicht involviert.

**Bemerkung 3.8.** Zwar ist die aktuelle Einbindung von MMA noch nicht darauf ausgelegt, dass Plugins wie ATHLET-CD Kommunikatoren anmelden können, dies ließe sich aber leicht durch den Einsatz des Hook-Mechanismus von libfde bewerkstelligen, welcher bereits durch ATHLET-CD genutzt wird.

#### 3.5.1 Einbinden von MMA

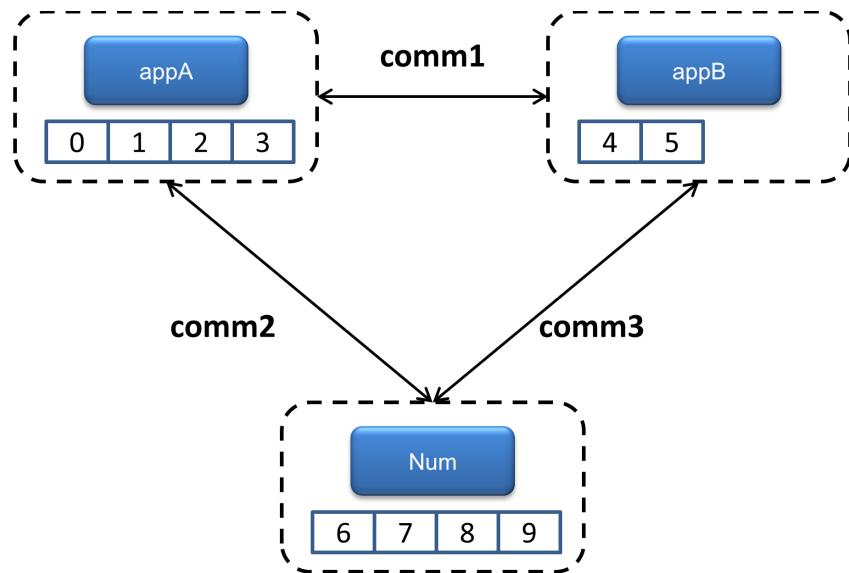
In der aktuellen Implementation wird die MMA-Bibliothek über das NuT-Plugin eingebunden, da hier der Anwendungsbezug für MMA liegt. Hierdurch ergibt sich auch der Vorteil, dass ATHLET keine direkte Abhängigkeit von MMA besitzt. Alternativ könnte ein dediziertes Plugin erstellt werden.

Das Schnittstellenmodul zu MMA ist eine für den gegebenen Kontext angepasste Version von `mma.fmod`, welches die generische FORTRAN-Schnittstelle zu MMA darstellt und Teil des Code-Kanons von MMA ist. Eingebunden wird die Adaption vom NuT-Wrapper. Das Aushandeln von Kommunikatoren wird dann im Zuge der Initialisierung des NuT-Plugins durchgeführt.

### 3.5.2 Funktionsweise von MMA

MMA ist darauf ausgelegt, Gruppen heterogener Prozesse paarweise über Kommunikatoren miteinander zu verknüpfen, siehe Abbildung 3.9. Hierbei wird in zwei Schritten vorgegangen:

- Registrierung der gewünschten Kommunikatoren in einer prozess-globalen Liste
- Aushandeln der Kommunikatoren, hierbei gibt es zwei Subschritte
  - Veröffentlichung der Liste
  - Kommunikatoren beitreten

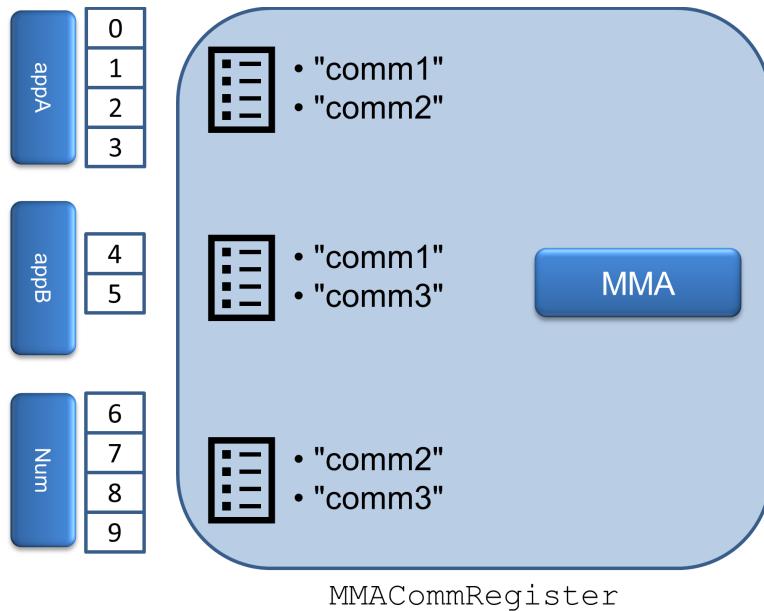


**Abb. 3.9** Paarweise Verknüpfung heterogener Prozessgruppen.

#### 3.5.2.1 Registrierung der Prozesse

Hierfür stellt MMA die dedizierte Routine `MMACommRegister` zur Verfügung, welche pro Prozess aufgerufen wird, siehe Abbildung 3.10. Kommunikatoren werden über Strings identifiziert. Dies vereinfacht die Koordination zwischen den verschiedenen Applikationen.

Sollen zwei Gruppen heterogener Prozesse miteinander kommunizieren, muss in beiden Gruppen der gleiche String-Identifikator registriert werden.



**Abb. 3.10** Registrierung der gewünschten Kommunikatoren pro Prozessgruppe in MMA mittels String-Listen

### 3.5.2.2 Aushandeln der Kommunikatoren

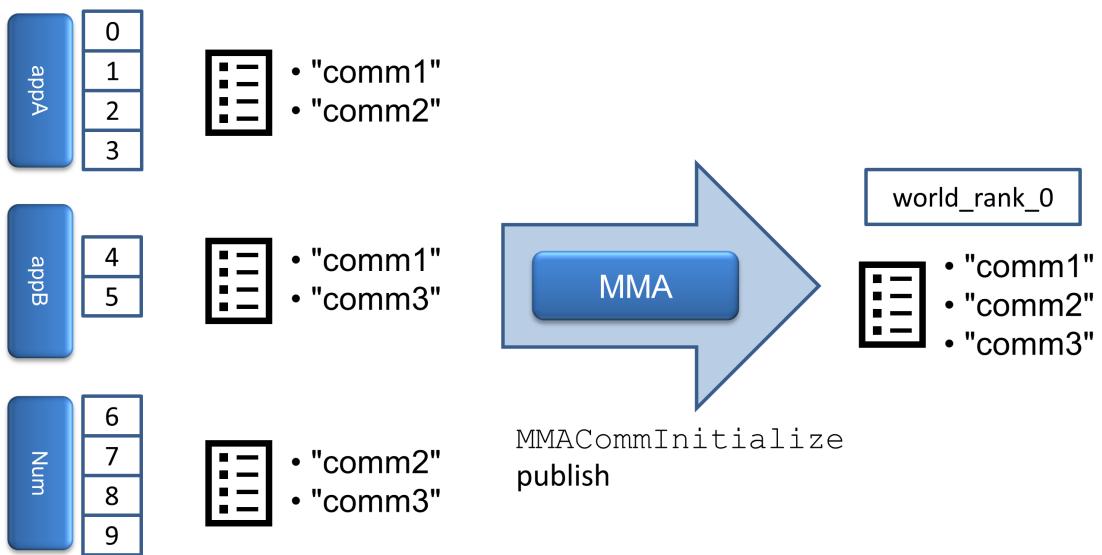
Für diesen Schritt wird ausgenutzt, dass sich sämtliche Prozesse standardmäßig im Kommunikator MPI\_COMM\_WORLD befinden. Die Ranks in diesem Kommunikator werden an dieser Stelle als *World-Ranks* bezeichnet. *World-Rank 0* übernimmt die Koordination des Aushandelns.

#### publish – Veröffentlichung der Listen

Sämtliche Prozesse ungleich *World-Rank 0* schicken ihre Liste an diesen. Dort wird eine globale Liste erstellt, welche jeden Kommunikatornamen genau einmal enthält, siehe Abbildung 3.11.

#### join – Kommunikatoren beitreten

Ausgehend von *World-Rank 0* werden nach dem Aushandeln die Kommunikatoren der globalen Liste sukzessive abgegangen. Jeder Prozess gleicht mit seiner lokalen Liste ab und nimmt bei Übereinstimmung nicht-trivial am MPI\_COMM\_SPLIT teil, siehe Abbildung 3.12. Zusätzlich werden Subkommunikatoren gebildet, um Prozesse zum gleichen Executable zu bündeln. Hierfür wird eine eindeutige `exe_id` genutzt, welche sich aus dem Pfad zum zugehörigen Executable generiert. Die `exe_id` wird zusätzlich zusammen mit dem

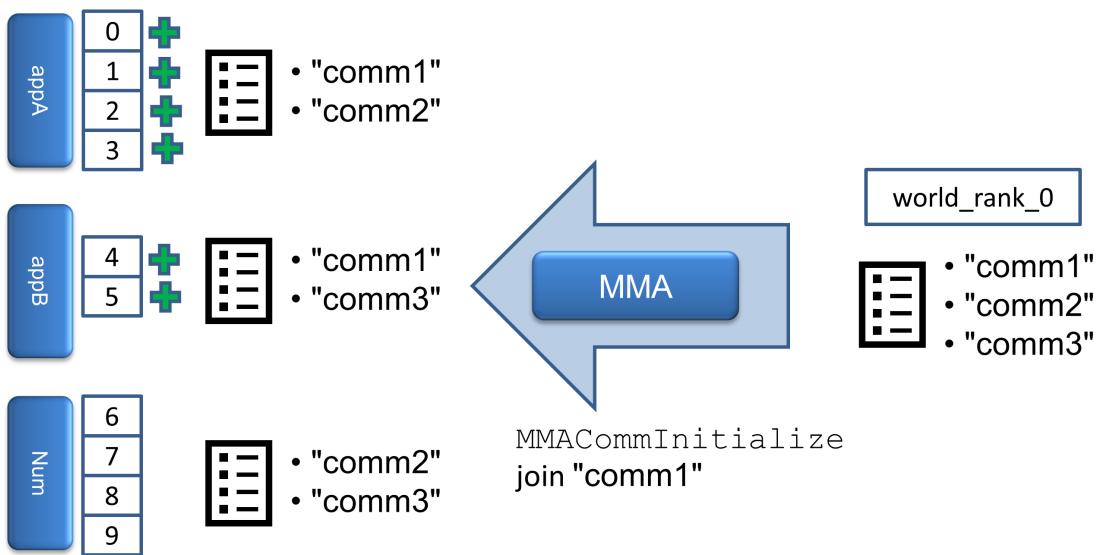


**Abb. 3.11** Veröffentlichung der individuellen Kommunikatorlisten an MMA und dortiges Generieren einer globalen Liste auf *World-Rank 0*

*World-Rank* eingesetzt, um die Prozesse sowohl im gemeinsamen Paar-Kommunikator als auch in den Subkommunikatoren in eine Reihenfolge zu bringen, so dass die Ordnung, welche durch die *World-Ranks* gegeben wird, erhalten bleibt, siehe Abbildung 3.13.

### Reservierung von Ressourcen und Zugriff auf diese

MMA speichert die relevanten (Sub-)Kommunikatorinformationen in einer Datenstruktur, die jedem Prozess für jedwede beliebige weitere MPI-Kommunikation (*Bcast / Send / Recv* etc.) zur Verfügung steht. Hinsichtlich des Zugriffs auf gemeinsame Ressourcen unterstützt MMA das Vermeiden von Verklemmungen ( gegenseitiges Blocken von benötigten Ressourcen). Innerhalb des AC<sup>2</sup>-Kontextes ist dies für evtl. gemeinsam genutzte NuT-Prozesse von Relevanz, vgl. Abschnitt 4.3. Hierbei ist entscheidend, dass MMA die Prozesse in den Subkommunikatoren in oben beschriebener Art und Weise an den *World-Ranks* ausrichtet, siehe Abbildung 3.14. Wird dann eine Reservierung von Ressourcen stets über eine Kommunikation zum *otherRank0* initiiert und in aufsteigender Ordnung der *Ranks* fortgesetzt, kann es nicht zu Verklemmungen kommen. Es kann somit nie der Fall eintreten, dass eine Prozess-Ressource reserviert wird und dann ob einer Verklemmung nicht genutzt werden kann. Modifikationen am NuT-Code sorgen zusätzlich dafür, dass ein ewiges Warten auf Ressourcen ausgeschlossen werden kann, siehe wieder Abschnitt 4.3.



**Abb. 3.12** Teilnahme an einem gewünschten und von MMA über *World-Rank 0* propagierten Kommunikator

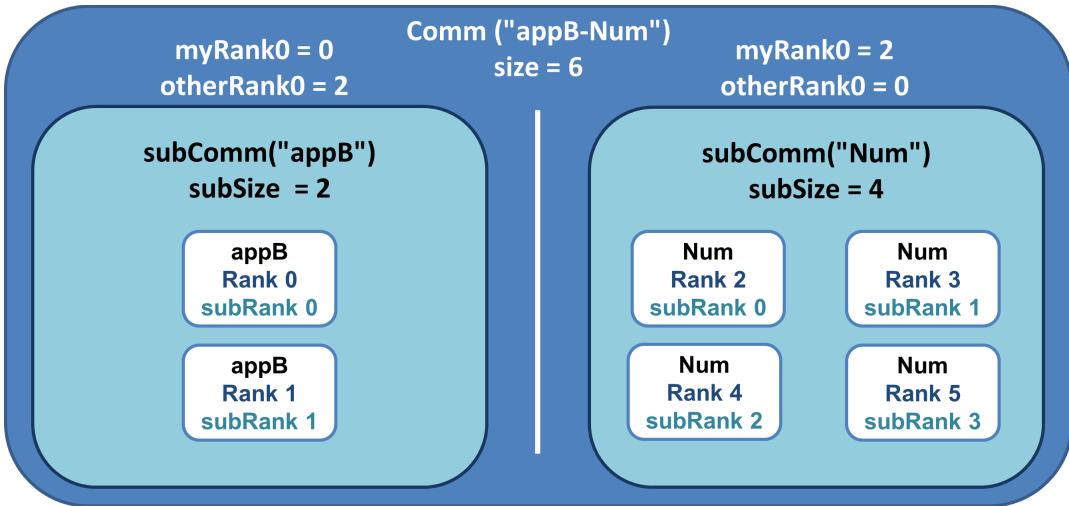
### 3.6 Speichern und Laden der Jacobimatrix im Restart-Szenario mittels NuT

ATHLET bietet dem Nutzer über die Input-Datei, siehe /LER 19, Kap. 5/, die Möglichkeit, während eines Simulationslaufes Restart-Punkte zu setzen, von welchen aus anschließend in einem neuen Lauf die Simulation wieder aufgenommen werden kann. Seitens der Mathematik heißt dies, dass sich  $t_0$  in (1.1) nach vorn verschiebt und ein entsprechend angepasstes  $y_0$  vorliegt. Ein möglicher Anwendungsfall ist, die Rechenzeit zu verkürzen, wenn sich ein zu untersuchendes Phänomen nicht in der zeitlichen Nähe des ursprünglichen Anfangswertes  $y_0$  befindet, welcher approximativ einen stationären Zustand beschreibt, siehe Abschnitt 3.3. Des Weiteren können z. B. Variantenrechnungen mitten im Ablauf ermöglicht werden.

#### Rahmenbedingungen

Obwohl die analytische Beschreibung des Vorganges recht geradlinig ist, geht dieser auf numerisch-softwaretechnischer Ebene mit einer gewissen Komplexität einher. Dies betrifft im Wesentlichen die Definition des Zustandes, welche Informationen aus dem vorherigen Rechenabschnitt übernommen und welche neu initialisiert werden. Dies gilt insbesondere für die Informationen zur Jacobimatrix aus (1.2).

Das Standardverhalten von ATHLET bei einem Restart liegt darin, stets eine neue Jacobimatrix zu erstellen. Keine Matrix muss gespeichert werden und nahezu alle Statusinfor-



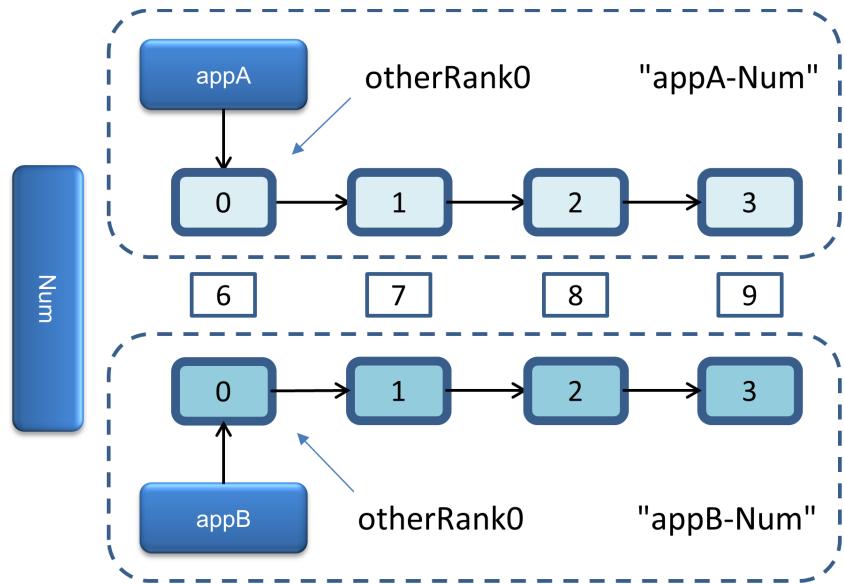
**Abb. 3.13** Via MMA definierte Kommunikatorstruktur inklusive Subkommunikatoren

mationen zur Jacobimatrix können neu initialisiert werden. Dies hat auch Auswirkung auf den ursprünglichen Verlauf, da dort ebenfalls die Berechnung einer neuen Jacobimatrix initiiert wird, wenn ein Restart-Punkt erreicht wird. Damit wird gewährleistet, dass sich der ursprüngliche Verlauf (mit aktiven Restart-Punkten) und ein neuer Simulationslauf ausgehend von einem Restart-Punkt hinreichend decken. Nachteil des Vorgehens ist, dass eine Jacobimatrix möglicherweise aus dem einzigen Grund, einen Restart-Punkt erreicht zu haben, gebaut wird. Die eigentliche Steuerlogik zur Differentialgleichungsnumerik würde an gegebener Stelle ggf. davon absehen. So ist das Potential für unnötige Mehrarbeit gegeben.

Ebenfalls als nachteilig kann es angesehen werden, dass sich eine Abweichung von dem Simulationslauf ergeben kann, welcher keinerlei Restart-Punkte berücksichtigt und somit Jacobimatrizen ausschließlich auf Basis der Steuerlogik aufdatiert. Zwar ist aufgrund von Gleitkomma-Arithmetik und der Verwendung approximativer Algorithmen ein berechneter Wert stets nur eine mögliche Option von vielen. Jedoch bietet sich obige Art des Simulationslaufes als Referenz für Vergleichsrechnungen an, da der Einfluss willkürlich gesetzter Restart-Punkte entfällt.

### Ziel und Vorgehen

Die Idee dieses Arbeitspunktes ist es, im Restart-Szenario die Jacobimatrix zu sichern und bei einem Restart entsprechend zu laden. Die Zeit zur Berechnung einer neuen Jacobimatrix wird eingespart. Somit wird sich direkt an der Referenz orientiert, und die Steuerlogik zur Differentialgleichungsnumerik wird nicht durch eine Restart-Logik überdeckt. Hierbei wird ausschließlich der Fall berücksichtigt, dass NuT aktiviert ist und somit



**Abb. 3.14** Verklemmungsfreier Ressourcenzugriff dank der MMA-induzierten Sortierung der Prozesse

für die Datenhaltung zur Jacobimatrix verantwortlich zeichnet. Dies ist sehr vorteilhaft, da PETSc bereits entsprechende I/O-Operationen für Matrizen anbietet.

Neben der Einbindung der entsprechenden I/O-Routinen seitens NuTs bedurfte es zur Realisierung der beschriebenen Funktionalität einer Erweiterung des Kommunikationsmodells zwischen ATHLET und NuT sowie eines Eingriffs in bestehenden ATHLET-Code. Diffizil ist bezüglich letzterem, welche zusätzlichen Daten gespeichert werden und zu welchem Zeitpunkt im Restart-Szenario diese wieder zur Verfügung stehen müssen.

Die Arbeiten innerhalb NuTs und zum Kommunikationsmodell waren sehr gut auf Basis der Arbeiten aus Abschnitt 2.3 und Abschnitt 4.2 durchführbar. Es stehen nun die NuT-bezogenen Routinen `linalg_jac_save` sowie `linalg_jac_load` für ATHLET (genau genommen für eine beliebige Host-Applikation) zur Verfügung. Die Jacobimatrixinformation wird in einer gesonderten Datei mit Endung `.nut.jac` abgelegt. Es findet ein PETSc-spezifisches Binärformat Anwendung. Vermittels entsprechender PETSc-Skripte kann eine derart gespeicherte Matrix leicht in Matlab oder Python eingelesen werden, um z. B. eine Spektral-Analyse durchzuführen oder die Besetzungssstruktur zu diskutieren. Zwar existierten begrenzte Möglichkeiten in diese Richtungen in älteren ATHLET-Versionen. Jedoch wurde auf proprietären Datenstrukturen gearbeitet, und es konnte nicht flexibel auf oben genannte dedizierte und funktionsmächtige Software zurückgegriffen werden. Das aktuelle Vorgehen geht somit mit einem echten Mehrwert einher und setzt einen neuen Standard hinsichtlich der Analysemöglichkeiten zur Jacobimatrix. Des Weiteren steht nun

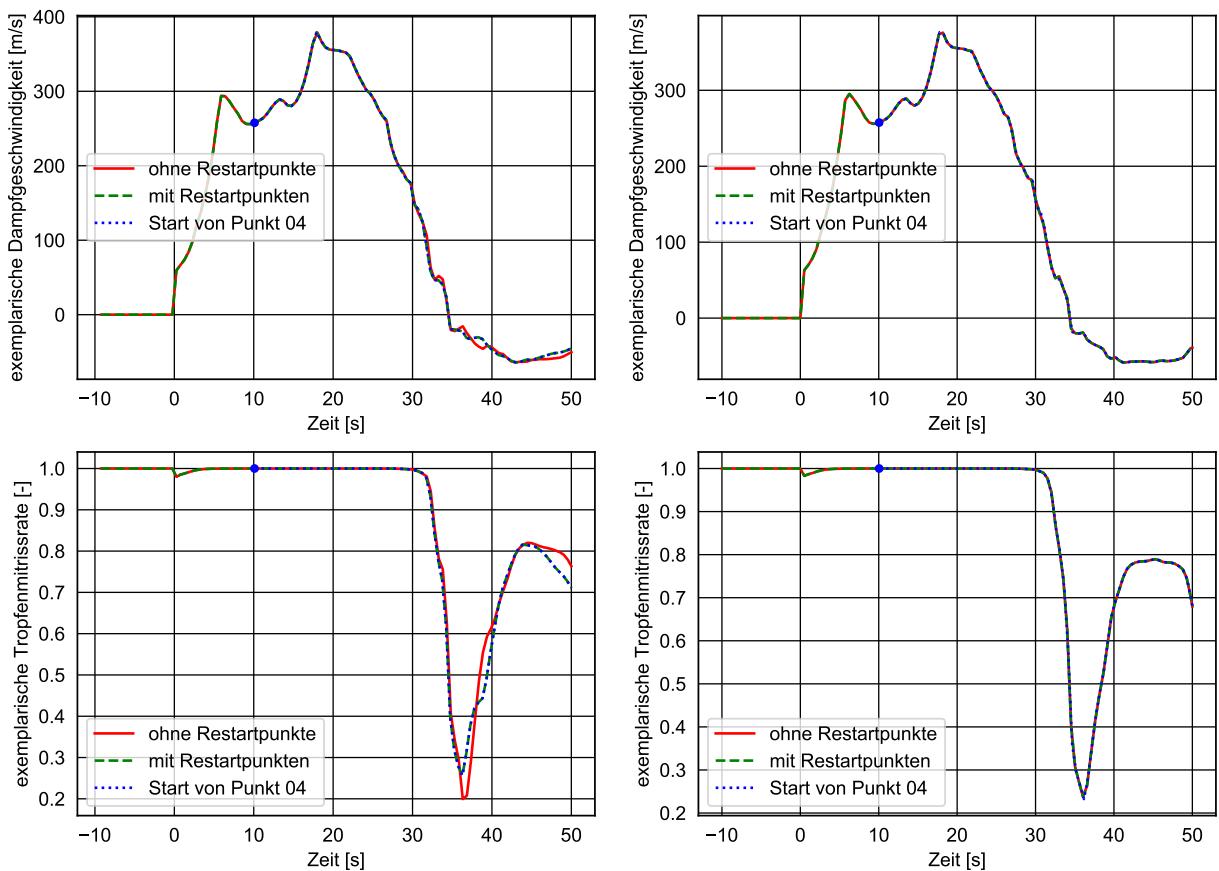
auch die Option offen, Jacobimatrixinformationen mit Dritten zu teilen und im Sinne von Eigenschaften und passenden Algorithmen zu diskutieren. Dies kann komplett jenseits des Anwendungshintergrundes passieren. Keine AC<sup>2</sup>-Installation ist hierfür notwendig.

Als Ergebnis einer längeren Untersuchung des relevanten ATHLET-Codes ergaben sich auch hier entsprechende Modifikationen. Im Wesentlichen betreffen diese die Routinen FREST und FEBE. Zusätzliche Zustandsinformation werden über FREST in das bekannte ATHLET-Datenformat zum Restart gespeichert. Sowohl formatiertes als auch unformatiertes Schreiben der Daten wird unterstützt. Hinsichtlich des Auslesens von Restartdaten werden die neu hinzugefügten Zustandsinformationen nicht direkt in die jeweiligen ATHLET-Variablen geladen, sondern zunächst neu definierte Größen mit dem Zusatz-Suffix \_REST zugeordnet, z. B. IQ\_REST statt IQ. Siehe FREST für eine komplette Liste dieser Variablen. Im LMSET-Block von FEBE, welcher genau einmal pro Lauf/Restart ausgeführt wird, finden dann die Zuordnungen an die tatsächlichen Variablen statt, z. B. IQ = IQ\_REST.

Die geleistete Implementierung geht mit folgenden Verhaltensweisen einher:

- Wird NuT für einen Simulationslauf aktiviert, werden automatisch die erweiterten Restart-Mechanismen berücksichtigt.
- An einem Restart-Punkt wird die Jacobimatrix stets gespeichert, selbst wenn von der Steuerlogik im nächsten Schritt eine neue Matrix gefordert wird.
- Es stellt kein Problem dar, einen Restart ohne NuT durchzuführen, wenn vorher NuT in Benutzung war, da die gespeicherten Daten eine Obermenge zu den Daten bilden, die für einen Restart von ATHLET ohne NuT benötigt werden.
- Ein Restart mit aktiviertem NuT ist nicht möglich, wenn vorher keine Jacobimatrix abgelegt wurde. Bei Bedarf kann die Logik aber durch zusätzliche Anpassungen erweitert werden.

Bisherige Tests haben ein zufriedenstellendes Ergebnis geliefert, siehe Abbildung 3.15. Es ist aber nicht auszuschließen, dass für andere Datensätze weitere Informationen gesichert werden müssen. Dies bleibt abzuwarten. Sind solche ggf. notwendigen Informationen aber erst einmal ausgemacht, ist es ein Leichtes, diese in den Restart-Prozess aufzunehmen, da das bisherige Vorgehen in gleicher Form weitergeführt werden kann.



**Abb. 3.15** Vergleich von Simulationsläufen im Restart-Szenario anhand exemplarischer Systemgrößen zum Datensatz PWR3Dsym. Linkseitig ein ATHLET-Lauf ohne NuT, rechtsseitig mit Einbindung von NuT.

## **4 AP3 – Aktualisierung der DGL-Numerik in der AC<sup>2</sup>-Komponente COCOSYS**

Ein wesentlicher Teil der Rechenzeit bei einer (linear-)impliziten Zeitintegration wird durch die zugehörige lineare Algebra (1.2) beansprucht. Jacobimatrizen müssen erstellt und lineare Gleichungssysteme gelöst werden. Geleistete Arbeiten zu diesem Arbeitspaket beziehen sich auf Anpassungen am AC<sup>2</sup>-Code. Der Vorsatz hierbei ist, NuT für die lineare Algebra im Thermohydraulik-Modul von COCOSYS zum Einsatz bringen zu können. Neben einem erhofften Performance-Gewinn liegt eine wesentliche Zusatz-Motivation in der Nutzung von NuT darin, die Zentralisierung und Vereinheitlichung der AC<sup>2</sup>-Numerik weiter voranzutreiben. Denn es ist ein übergeordnetes Ziel, den Code pflegeleichter, besser erweiterbar und zukunftssicherer zu machen.

Grundsätzlich lassen sich die Aufgaben zu diesem Paket in drei Kategorien einteilen:

- funktionale Anpassungen des relevanten COCOSYS-Codes,
- kommunikationsbezogene Modifikationen auf Seiten von COCOSYS, um mit NuT Daten austauschen zu können,
- Bereitstellung einer allgemeinen Kommunikationsstrategie im AC<sup>2</sup>-Kontext.

Die geleisteten Arbeiten decken bereits den ersten und den letzten Punkt ab, siehe hierzu Abschnitt 4.1 und Abschnitt 4.2 beziehungsweise Abschnitt 4.3. Aufgrund der bestehenden Softwarestruktur und Compiliervorgänge von COCOSYS traten unerwarteten Schwierigkeiten während der Arbeiten auf, COCOSYS zugänglich zu einer allgemeinen AC<sup>2</sup>-Kommunikationsstrategie zu machen. Siehe hierzu auch die Anmerkungen in Abschnitt 4.2.3. Damit ist der zweite Punkt nicht abschließend erfüllt, und es konnten noch keine Beispielsimulationen für das Tandem COCOSYS/NuT durchgeführt werden.

### **4.1 Analyse des Ist-Zustandes der Differentialgleichungsnumerik im Thermohydraulik-Modul von COCOSYS**

Um sich ein Bild der allgemeinen Zeitintegrationsprozesse in COCOSYS zu machen und zu eruieren, ob diese hinreichend unabhängig voneinander sind, wurde zunächst eine umfangreiche Untersuchung über mehrere Module hinweg initiiert. In allen diesen Modulen wird eine Variante von FEBE (Forward Euler / Backward Euler) genutzt. Hierbei handelt es sich um ein Extrapolationsverfahren basierend auf dem expliziten und linear-impliziten Euler-Verfahren, siehe z. B. /HOF 81/ oder auch Abschnitt 3.4.

## **AFP**

Dieses Modul beinhaltet gleich zwei Integratoren: Nach Kommentarangaben für das FIPHOST-Modell beziehungsweise für den Aerosol-Code AMEROS/MGA. Hinsichtlich des FIPHOST-Modells zeigt sich der zugehörige FEBE-Code FHFEBE als *expliziter* Löser. Keine Jacobimatrixinformationen werden generiert oder verarbeitet. Für den Aerosol-Code ist die Logik der Zeitintegration in der Routine AFEBE implementiert. Hier wird auf linear-implizite Zeitintegration zurückgegriffen. In entsprechenden weiteren Routinen wird die Jacobimatrix per finite Differenzen bestimmt sowie die Systemmatrix (diagonal-geshiftete Jacobimatrix) in Hessenberggestalt gebracht und anschließend LU-faktorisiert. Dies ist ein klassisches Vorgehen wie z. B. in /HAI 96, Kap. IV.8/ beschrieben. Nachteil ist, dass eine Dünnbesetztheit der Jacobimatrix hierbei nicht berücksichtigt wird und durch die Hessenberg-Transformation in der Regel komplett zerstört wird. Es wird an jener Stelle also mit vollen Systemen gerechnet. Es existieren keine Strukturinformationen zur Dünnbesetztheit.

## **NEWAFP**

Im Gegensatz zu den FORTRAN-77-Routinen in AFP werden in diesem Modul bereits FORTRAN90/95-Techniken eingesetzt, um eine numerische Zeitintegration zu verwirklichen. Sämtliche Logik hierzu ist im FORTRAN-Modul `febe_module` gekapselt, was die Navigation und Klärung der Abhängigkeiten deutlich vereinfacht. Auch hier wird sich der gleichen klassischen Techniken für die lineare Algebra bedient wie in AFEBE mit den gleichen Nachteilen hinsichtlich der Ausnutzung von Dünnbesetztheit.

## **TH**

Der FEBE-Routine im Thermohydraulik-Modul von COCOSYS ist deutlich die Verwandtschaft mit der gleichnamigen ATHLET-Routine anzusehen, obschon an dieser Stelle wieder FORTRAN-77-Strukturen vorherrschen und weniger Funktionalität geboten wird. Z. B. gibt es keine partiellen Updates für die Jacobimatrix. Dies deutet darauf hin, dass der Fork (unabhängige Weiterentwicklung von einer gemeinsamen Basis ausgehend) vor mehreren Dekaden stattgefunden hat. Auf der anderen Seite wird bereits Struktur-Logik über das Flag IFTRIX in der untergeordneten FTRIX-Routine berücksichtigt. Die Option IFTRIX=3 wird unterstützt und beschreibt wie in ATHLET ein Vorgehen, in der linearen Algebra die Dünnbesetztheit der Jacobimatrix sowohl in Speicherstrukturen als auch dem Lösungsprozess auszunutzen. Neben IFTRIX sind auch Arrays wie TOP (verfolgt nach, ob Gleichungen an- oder abgeschaltet sind) oder TIP (protokolliert, ob Gleichungen explizit oder implizit berechnet werden) vorhanden. Hinsichtlich der Speicherstrukturen zur Jacobimatrix ergeben sich zum Bekannten aus ATHLET jedoch im Detail einige

Unterschiede. Hierauf wird im nächsten Abschnitt eingegangen. In der Summe sehen die Voraussetzungen seitens der Code-Logik aber gut aus, NuT zum Einsatz bringen zu können.

### **Ausbleiben von Seiteneffekten**

Im Rahmen der Untersuchung ist kein Anzeichen dafür entdeckt worden, dass sich die einzelnen FEBE-Derivate gegenseitig beeinflussen oder gemeinsame Datenstrukturen nutzen. Somit wird davon ausgegangen, dass die Strukturen in TH-FEBE ohne Seiteneffekte modifiziert werden können.

## **4.2 Zugang von COCOSYS zur AC<sup>2</sup>-einheitlichen Handhabung der DGL-Numerik über die NuT-Architektur**

Entsprechend der Erläuterungen in der Präambel zu diesem Kapitel und auf Basis der Analyse im vorherigen Abschnitt wurden Anpassungen am bestehenden COCOSYS-Code zum Thermohydraulik-Modul durchgeführt. Ziel dieser Anpassungen ist es, seitens der Anwendung, d. h. COCOSYS, diesen auf funktionaler Ebene kompatibel zu NuT zu machen. Arbeiten, welche generellen Zugriff auf NuT-Objekte respektive auf Referenzen hierzu erlauben, sind nicht Gegenstand dieses Arbeitspunktes. Siehe hierzu auch Abschnitt 4.2.3 sowie Abschnitt 4.3.

### **4.2.1 Funktionsbasierte Modifikationen am Code des Thermohydraulik-Moduls**

Die erbrachten Modifikationen werden pro COCOSYS-Subroutine beschrieben. Die ATHLET-Pendants hierzu sind in eckigen Klammern angegeben.

#### **FTRIX, [FTRIX]**

Diese Routine übernimmt die allgemeine Koordination der Funktionalitäten zur linearen Algebra und ruft entsprechende UnterROUTinen auf. Wesentliche zu koordinierende und delegierende Aufgaben sind hierbei

- Abbildung der Netzwerkstruktur auf Matrixstrukturen und Vektorvariablen (FTAB1 sowie FORG [letzteres wird zusätzlich auch in FTAB1 aufgerufen])
- Seeding zur Jacobimatrix<sup>1</sup> (FJAORD sowie FTSEQB)

---

<sup>1</sup>Siehe Abschnitt 5.2.2 zur Erläuterung des Seeding-Prozesses.

- Permutation der Spalten und Zeilen der Jacobimatrix zur Fill-in-Reduktion (FDEFMI über FTAB1-Aufruf sowie FPERM)
- Befüllen der Jacobimatrix mit Werten über finite Differenzen (FMABLT und FPACSB)
- Block-LU-Zerlegung der Systemmatrix aus (1.2) sowie Lösen der linearen Systeme (FSBMS0)

Das Numerical Toolkit ist auf die Strukturinformationen zur Jacobimatrix angewiesen, welche durch FTAB1 und FORG generiert werden, siehe die nächsten beiden Paragraphen. Initial, oder wenn sich die Struktur ändert, wird in FTRIX ein Flag newTAB gesetzt, um die anwendungsseitige NuT-Logik auf die Verarbeitung einer neuen Matrix vorzubereiten, siehe auch Abschnitt 4.2.2.

Der Seeding-Prozess wird von NuT selbst durchgeführt. Insofern ist nur die Information, weiterzugeben, dass eine neue Matrix vorliegt. Dies passiert ebenfalls über das Flag newTAB. Alle weiteren Punkte in obiger Liste werden direkt von NuT übernommen. Es muss also lediglich dafür gesorgt werden, dass obige dedizierte Routinen nicht aufgerufen werden und folgende Ersetzungen stattfinden

```
FMABLT → FMANUT,  
FSBMS0 → linalg_jac_solve.
```

Konkret werden die beschriebenen Modifikationen über if-else-Anweisungen realisiert, um die originäre Funktionalität nach wie vor parat zu haben. Die Routine FMANUT kümmert sich um die Handhabung der Jacobimatrix, siehe auch Abschnitt 4.2.2. Mittels linalg\_jac\_solve werden die auftretenden linearen Systeme gelöst. Als Parameter ergeben sich die jeweilige rechte Seite sowie der Diagonalshift zur Jacobimatrix, vgl. (1.2). NuT nutzt dann die zuvor über FMANUT gespeicherte Jacobimatrix, um (1.2) zu lösen. Das Ergebnis wird als out-Parameter zurück geschrieben.

### **FTAB1, [FTAB1]**

Die Modifikationen in dieser Subroutine beschränken sich auf das Auslassen, die Routine FDEFMI auszuführen, wenn NuT aktiv ist. Mittels FDEFMI wird eine Variante des Minimum-Deficiency-Algorithmus von Tinney und Walker, /TIN 67/, ausgeführt, um den Fill-in in den LU-Faktoren der Systemmatrix zu verkleinern, wenn zur Lösung von (1.2) die Systemmatrix Gauß-zerlegt wird. Innerhalb des Numerical Toolkits wird diese Fill-in-Reduktion direkt von den einzelnen Lösern in PETSc oder MUMPS übernommen. Details hierzu finden sich in /STE 20, Unterabsch. 3.3.1/. Der Minimum-Deficiency-Algorithmus findet auch im ATHLET-Code Anwendung. Dieser liefert gute Ergebnisse, ist jedoch gerade für größere

Systeme sehr langsam. Im Vergleich hierzu ist der Zeitaufwand für die Fill-in-Reduktion in den NuT-Algorithmen vernachlässigbar, siehe /STE 17b, Unterabsch. 6.2.3/.

### **FORG, [FTABNB]**

Strukturell existiert die Jacobimatrix auf zwei Ebenen, zum einen auf Blockebene, zum anderen auf Elementebene. Die Blockstruktur korrespondiert direkt mit dem zugrundeliegenden Netzwerk aus *control volumes* und *junctions*. Ein Block an der (Block-)Position  $(i, j)$  ist genau dann gesetzt, wenn das  $i$ -te Netzwerkelement mit dem  $j$ -ten in direkter Wechselwirkung steht. Da diese Relation symmetrisch ist, gilt dies auch für die Blockstruktur der Jacobimatrix. Die Elementebene korrespondiert zu den Gleichungen, welche die Netzwerkelemente im Sinne des thermohydraulischen Verhaltens beschreiben.

In dieser Routine werden die beschriebenen Strukturinformationen der Jacobimatrix auf Basis des Netzwerkes ermittelt und in Arrays gespeichert, namentlich in ITAB2 und ITAB5.

ITAB2 verknüpft die Block- mit der Elementebene, indem die Anzahl der Gleichungen pro Netzwerkelement gespeichert werden. Diese ist somit auch gleich der Anzahl von Elementzeilen pro Blockzeile in der Jacobimatrix. Des Weiteren werden die Indexbereiche im permutierten Lösungsvektor vermerkt, die jeweils zu einem Netzwerkelement korrespondieren. Damit entspricht das hier definierte ITAB2 genau dem Gegenstück, welches im ATHLET-Code in FTABNB genutzt wird.

Im Array ITAB5 finden sich die Spaltenindizes der nicht-trivialen Blöcke pro Blockzeile wieder. Trotz gleicher Benennung ist die Struktur des Arrays *nicht* die gleiche wie die des in FTABNB verwandten Arrays. In COCOSYS wird pro Blockzeile für jede Spalte ein Wert gespeichert. Zunächst werden die Spaltenindizes der nicht-trivialen Blöcke gesichert, dann wird mit Nullen aufgefüllt. In ATHLET wird bereits eine komprimierte Speicherform genutzt, welche die explizit gespeicherten Nullen herausnimmt. Dies macht es nötig, ein weiteres Array einzuführen, welches die Indizes in ITAB5 markiert, an denen eine neue Blockzeile beginnt. Diese Information ist in ATHLET im Array ITAB7 zu finden. Trotz eines weiteren Arrays liegt der Vorteil auf der Hand. Bezeichne  $n$  die Dimension der Blockmatrix, so speichert COCOSYS  $\mathcal{O}(n^2)$  Werte, wohingegen ATHLET nur  $\mathcal{O}(n + nnz)$  Speicherplätze benötigt. Hierbei bezeichnet  $nnz$  die Anzahl an Nicht-Null-Elementen der Blockmatrix.

NuT ist darauf ausgelegt, mit Strukturinformationen zu arbeiten wie sie durch ITAB2, ITAB5 (ATHLET) und ITAB7 gegeben sind. Letztere zwei Arrays entsprechen auch im Wesentlichen einem standardisierten CSR-Format (Compressed Sparse Row) für Matrizen. Es wurden somit zwei neue Variablen in den FORG-Code eingeführt, ITAB5c sowie

ITAB7, und analog zum ATHLET-Vorgehen in FORG mit den passenden Werten befüllt. Die Einbindung geschieht über use-Assoziation:

```
use nut_jac, only: init_ITAB7, ITAB7, init_ITAB5compressed,  
*                      ITAB5c => ITAB5
```

Hierbei sind zusätzlich zwei Funktion zum Initialisieren der beiden Arrays definiert, siehe nächsten Abschnitt. Das Modul `nut_jac` kapselt die anwendungsseitige Funktions-Logik von NuT. Da bereits `ITAB5` in `FORG` existiert, wird der `=>`-Operator genutzt, um `ITAB5` aus `nut_jac` einbringen zu können.

#### 4.2.2      **Funktionsbasierte Modifikationen des anwendungsseitigen NuT-Codes**

Das bereits im ATHLET-Umfeld eingesetzte Modul `nut_jac` musste einige Anpassungen erfahren, um in den gegebenen Kontext zu passen. Wie im vorherigen Abschnitt erläutert war es nötig, die Arrays `ITAB5` sowie `ITAB7` mit der Semantik wie in ATHLET einzuführen. Diese sind als Modulvariablen in `nut_jac` angelegt. Ebenfalls sind hier die Implementationen der Initialisierungsroutinen für diese beiden Arrays zu finden.

Die wesentliche Routine in `nut_jac` ist `FMANUT`, welche die anwendungsseitige NuT-Routine zur Bestimmung der Jacobimatrix darstellt. `FMANUT` wird im ATHLET-Umfeld ohne eine Parameterliste aufgerufen. Sämtliche benötigten Daten sind entweder im Modul vorhanden oder über use-Assoziation einbindbar.

In COCOSYS herrschen FORTRAN-77-Techniken vor: Viele globale Variablen werden über COMMON-Blöcke definiert und stehen über separate Dateien zur Verfügung, die mittels der FORTRAN-Anweisung `INCLUDE` eingebunden werden. Es ist keine Option, für eine Adaption von `nut_jac` auf Legacy-Techniken zurückzugreifen. Stattdessen wurde für die Routine `FMANUT` eine Parameterliste eingeführt, welche die Größen enthält, die nicht über Module zur Verfügung stehen. Wie im vorherigen Abschnitt beschrieben wird `FMANUT` in `FTRIX` aufgerufen. Da es sich bei letztere um eine Routine handelt, die überwiegend FORTRAN 77-Techniken einsetzt, wurden an dieser Stelle zusätzlich benötigte Größen über weitere `INCLUDE`-Anweisungen in `FTRIX` sichtbar gemacht und der Parameterliste von `FMANUT` hinzugefügt.

Im ATHLET-Kontext greift `FMANUT` auf die ATHLET-eigene Funktion `FMADEL` zu, um die Störungen für die definiten Differenzen berechnen zu lassen. Im COCOSYS-Kontext ist diese Funktionalität Teil der Routine `FMABLT`, welche gerade durch `FMANUT` ersetzt wird. Folglich wurde `nut_jac` eine Modul-Routine namens `FMADEL` hinzugefügt, welche die

FMABLT-Variante der Berechnungen der Störungen enthält. Auch hier findet die gleiche Idee Anwendung, die Parameterliste zu erweitern, um Legcay-Techniken zu vermeiden.

Obschon die Grundstrukturen des Codes zur linearen Algebra in COCOSYS und ATHLET sehr ähnlich sind, ist klar ersichtlich, dass die ATHLET-Variante moderner ist und mit einem größeren Funktionsumfang aufwarten kann. Im Wesentlichen umfasst dieser

- partielle Updates der Jacobimatrix,
- effiziente Auswertung der rechten Seite der Differentialgleichung.

Es ist nicht Teil des Arbeitspunktes, diese Erweiterungen auch in COCOSYS verfügbar zu machen. Vielmehr folgen hieraus weitere Anpassungen an FMANUT. Im Wesentlichen bestehen diese darin, entsprechende Codeabschnitte herauszunehmen bzw. auszukommentieren. Dies gilt ebenfalls für Codeabschnitte, die sich um das Laden der Jacobimatrix im Restartfall kümmern, siehe Abschnitt 3.6, oder um Abschnitte, welche Logik bezüglich des Kohärenz-Tests zur Berechnung der Jacobimatrix enthalten, siehe Abschnitt 5.2.

Weitere Anpassungen wurden vorerst nicht vorgenommen. Es ist unwahrscheinlich, dass hinsichtlich der Funktionalität weitere wesentliche Modifikationen vonnöten sein könnten. Sehr wohl stehen aber noch Anpassungen hinsichtlich der Anbindung von NuT aus, siehe hierzu den nächsten Abschnitt.

#### **4.2.3 Ausstehendes Bindeglied**

In Abschnitt 4.3 wird die Notwendigkeit einer neuen und standardisierten Kommunikationsstrategie zwischen den AC<sup>2</sup>-Komponenten beschrieben. Dieser Standard ist bereits in ATHLET und NuT durch die Anbindung der MMA-Bibliothek implementiert. Die Kommunikationsstrategie im COCOSYS-Code wurde ebenfalls schon für die Einbringung von MMA vorbereitet. Parallel werden zu diesem Zeitpunkt Arbeiten zur Vereinheitlichung und Modernisierung des AC<sup>2</sup> Build-Systems durchgeführt, die in ATHLET und NuT bereits abgeschlossen sind, jedoch in COCOSYS noch bearbeitet werden.

Um einen Mehraufwand und doppelte Arbeiten im alten und neuen Build-System zu vermeiden, wird das Einbringen von MMA in COCOSYS erst nach der Umstellung des Build-Systems umgesetzt. Dann kann die Anbindung von MMA kurzfristig und nachhaltig erfolgen.

#### **4.3 Ausarbeitung einer Kommunikations- und Parallelitätsstrategie im Arbeitsverbund der AC<sup>2</sup>-Komponenten COCOSYS und ATHLET mit dem Numerical Toolkit**

AC<sup>2</sup> ist bereits ein komplexer und wachsender Zusammenschluss von Softwarekomponenten, die mehrere Prozesse in Anspruch nehmen können. Je nachdem wie die Komponenten verteilt sind, erfolgt deren Kommunikation entweder innerhalb eines Prozesses oder zwischen zwei Prozessen. Da je nach Anforderung prinzipiell jede Komponente der Möglichkeit bedarf, mit einer beliebigen anderen kommunizieren zu können, wurde eine standardisierte globale Kommunikationsstrategie entwickelt. Zu ihren Anforderungen gehört insbesondere ein hohes Maß an Skalierbarkeit. Damit soll sichergestellt werden, dass sich durch das zukünftige Hinzufügen von neuen Komponenten und zusätzlichen Kommunikationskanälen, die Komplexität und der Wartungsaufwand nicht unverhältnismäßig erhöht.

##### **Skalierbarkeit**

Um gute Skalierbarkeit zu erhalten, sollte jede Komponente dediziert selbst bestimmen, mit wem sie kommunizieren möchte. Das ist mit dem statischen Kommunikationsmodell von MPI nicht ohne Weiteres machbar, da das Erstellen der erforderlichen Kommunikatoren ohne vorhergehende globale Kommunikation nicht möglich ist. So entstehen globale Abhängigkeiten zwischen all jenen Komponenten, die sich in welcher Form auch immer an der Kommunikation beteiligen möchten. Diese unerwünschten Abhängigkeiten wurden mithilfe der standardisierten Erzeugung der Kommunikatoren durch MMA abgekapselt und vereinfacht, siehe Abschnitt 3.5.

##### **Lose gekoppelte Komponenten im gleichen Prozess**

Ebenso ist es möglich, individuelle Kommunikatoren für mehrere lose gekoppelte Komponenten innerhalb eines Prozesses zu erzeugen. Wichtig ist hier, dass jede Komponente seine Kommunikatoren selbst registriert und dieser Vorgang nicht delegiert wird. Als Randbedingung dazu muss die Registrierung zeitlich vor der globalen Aushandlung der Kommunikatoren durch MMA durchgeführt werden. Im Kontext von AC<sup>2</sup> kann dies mithilfe des bereits eingesetzten *Hook*-Konzeptes der libfde-Bibliothek umgesetzt werden. Jede Komponente stellt eine Funktion zur Verfügung, die ihre gewünschten Kommunikatoren bei MMA registriert, und hinterlegt diese beim globalen *Hook*. Die registrierten Funktionen können dann wiederum gesammelt über den *Hook* der Hauptkomponente aufgerufen werden, ohne diese im Vorfeld zu kennen (Konzept *dependency injection*). Somit kann der Hauptprozess gewährleisten, dass alle Komponenten ihre Kommunikatoren registrieren,

bevor MMA initialisiert wird. Des Weiteren werden Abhängigkeiten der Hauptkomponente reduziert. Dies wird beispielsweise für den ATHLET-Prozess benötigt. So kommuniziert eine *Controller*-Komponente mit COCOSYS, ATHLET selbst mit NuT und gegebenenfalls zukünftig die Komponente ATHLET-CD auch direkt mit NuT. Damit ist der ATHLET-Prozess in der Lage für alle Komponenten die Kommunikation zum richtigen Zeitpunkt aufzusetzen, ohne dass unnötig kompliziert zwischen den Komponenten kommuniziert werden muss.

### Berücksichtigung des Numerical Toolkit

Die NuT-Prozesse selbst bestehen nur aus einer Komponente, weshalb dort das *Hook*-Konzept zum Aushandeln der Kommunikatoren nicht nötig ist. Damit die Kommunikation funktioniert, wird im NuT-Worker von jeder AC<sup>2</sup>-Komponente, die NuT prinzipiell in Anspruch nehmen könnte, der Kommunikatorname hinterlegt. Nachdem MMA die Kommunikation aufsetzt, sortiert der NuT-Worker alle Kommunikatoren aus, bei denen sich keine Gegenstelle registriert hat. So können beispielsweise AC<sup>2</sup>-Komponenten wie ATHLET einfach zur Laufzeit entscheiden, ob mit oder ohne NuT gerechnet werden soll. Dies geschieht, indem sie sich optional für den vordefinierten NuT-Kommunikator registrieren oder dies unterlassen. Sollte sich keine Komponente für NuT registrieren, so beendet sich der NuT-Prozess automatisch.

Der NuT-Worker wurde so entwickelt, dass er als Ressource von mehreren Host-Applikationen gleichzeitig genutzt werden kann. Im sequenziellen Betrieb, d. h. ein einziger gestarteter NuT-Worker-Prozess, empfängt er auf allen aktiven Kommunikatoren mögliche Anweisungen der gekoppelten Host-Applikationen. Sobald eine Anweisung empfangen wurde, wird sie unmittelbar abgearbeitet. Zwischenzeitlich können auch weitere Anweisungen der verbleibenden Host-Applikationen empfangen werden, deren Abarbeitung allerdings solange blockiert wird, bis die ursprüngliche Host-Applikation bedient wurde. Falls mehrere Anfragen zur Abarbeitung gleichzeitig zur Verfügung stehen, werden diese im *Round-Robin-Verfahren* nacheinander ausgewählt. Damit wird eine faire Behandlung aller Anfragen durchgesetzt und es kann ausgeschlossen werden, dass bei hoher Auslastung die Arbeiten einzelner Host-Applikationen bevorzugt ausgewählt werden und andere Prozesse dafür ewig warten müssen.

Weitere Vorkehrungen wurden getroffen, damit auch der parallele Betrieb von mehreren NuT-Worker-Prozessen für mehrere Host-Applikationen gleichzeitig möglich ist. Abstrakt betrachtet stellt nun jeder NuT-Worker-Prozess eine Teilressource von NuT dar, die ohne die anderen NuT-Worker-Prozesse nicht arbeitsfähig ist. In dem Fall muss eine Host-Applikation also zuerst alle relevanten NuT-Worker-Prozesse für sich akquirieren,

bevor diese genutzt werden können. Die Akquise der einzelnen NuT-Prozesse erfolgt immer nacheinander in globaler Ordnung der *Ranks*, wie sie von MMA zur Verfügung gestellt werden. Somit kann Verklemmungsfreiheit garantiert werden, also dass zwei Host-Applikationen jeweils einen Teil der Ressourcen akquiriert haben und die Ressourcen der jeweils anderen benötigen, um fortzufahren.

## 5 AP4 – Validierung der Modifikationen

### 5.1 Release und QS zu NuT 1.0

Während der Projektlaufzeit wurde das GRS-Programm-Paket AC<sup>2</sup> in der Fassung 2019 zum ersten Mal dem Endnutzer zur Verfügung gestellt. Neben den Simulationsprogrammen ATHLET, ATHLET-CD und COCOSYS ist auch erstmalig das Numerical Toolkit in der Version 1.0 allgemein nutzbar. In 2020 folgte der Patch 2019.1 für AC<sup>2</sup> inklusive der Version 1.0.1 von NuT. Sowohl für Release und Patch wurde die FORTRAN-Variante von NuT genutzt. Zum Release wurde noch an den Neuerungen gearbeitet, die in Abschnitt 2.3 beschrieben sind. Da sich das Interface mit den Umstellungen zu C/C++ geändert hat, wird die neue Implementation NuT 2.0 erst mit dem nächsten Major Release von NuT allgemein zur Verfügung stehen, vgl. hierzu auch /PRE 20/.

Die Arbeiten zum Release und Patch wurden je von einem umfänglichen Validierungsprozess begleitet. Diese richteten sich nach dem Teilkernprozesses TKP 03-05 „Softwareentwicklung“ /HEC 20/ der GRS aus. Es wurde eine ausführliche Nutzeranleitung /STE 20/ für NuT erstellt und mit den Tätigkeiten zum Patch verfeinert. U.a. werden folgende wichtige Themengebiete in der Anleitung behandelt:

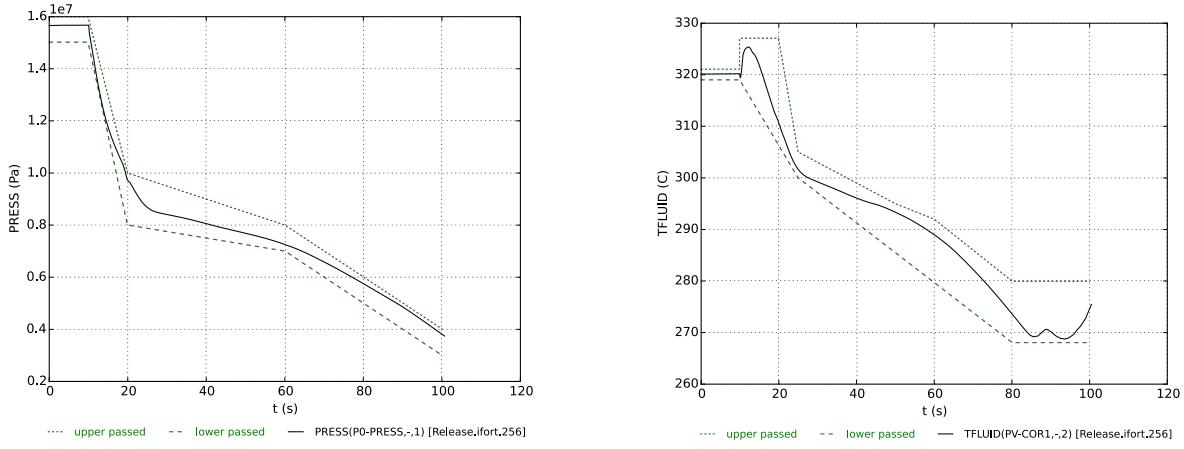
- Anwendungsgebiet von NuT und Vorteile für den Nutzer,
- Auflistung der angebotenen Solver-Presets zur linearen Algebra nebst Erläuterungen,
- Aktivierung von NuT per AC<sup>2</sup>-GUI und per Kommandozeile,
- Ausführliche Beschreibung des NuT-bezogenen Outputs.

Zusätzlich wurden die Vorgaben aus dem TKP 03-05 genutzt, um einen detaillierten QS-Plan aufzustellen, welcher mit den Arbeiten zum Patch um eine entsprechende Sektion erweitert wurde. Hierbei ist sowohl ein Einzeltest zu NuT sowie ein Tandem-Test mit ATHLET definiert und für Release sowie Patch durchgeführt und archiviert worden.

#### NuT-Einzeltest

Dieser Test bindet die Schnittstelle zu NuT direkt über use-Assoziation ein und geht wichtige NuT-Funktion zur linearen Algebra sukzessive anhand eines Beispielsystems ab. Hierzu zählen

- Anlegen der Speicherstruktur zur Jacobimatrix,
- Befüllung mit Werten,
- Lösen des Systems und Ausgabe der berechneten Lösung.



(a) Druck (P0-PRESS,-,1)

(b) Temperatur, Fluid (PV-COR1,-,2)

**Abb. 5.1** Sample 1 – Kontrollvariablen zur Validierung, Solver-Preset `mumps-gmres`

Die rechte Seite des linearen Systems ist derart gewählt, dass die Lösung zum Eins-Vektor verifiziert werden kann.

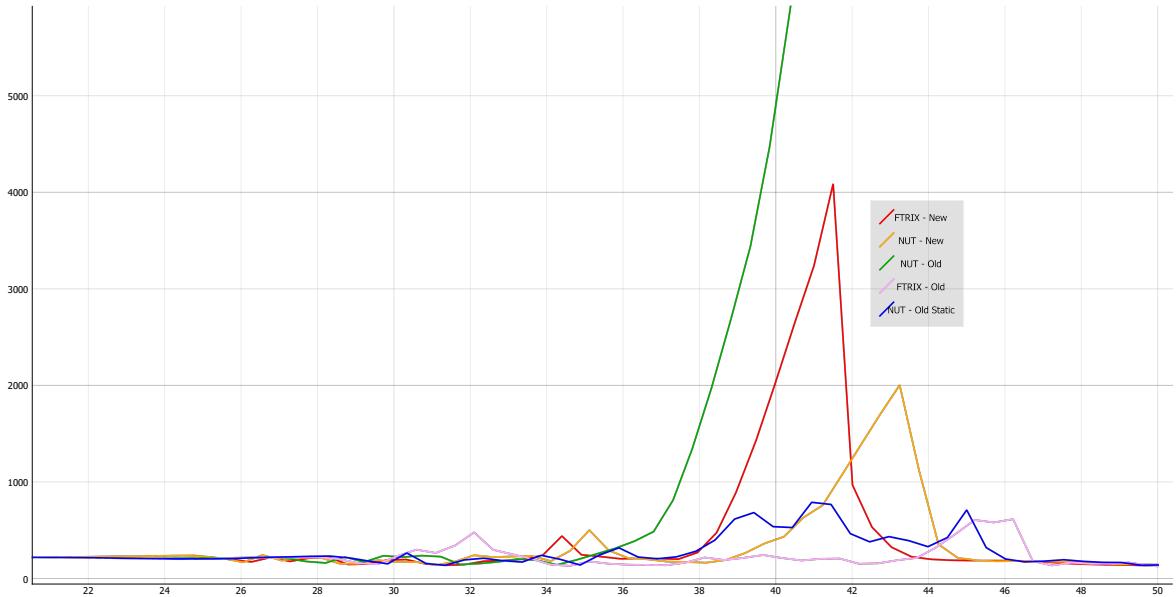
### Tandem-Test mit ATHLET

In diesem Test wird NuT in Konjunktion mit ATHLET für den Datensatz `sample1.in` ausgeführt, welcher zum Kanon von Beispieldatensätzen zu ATHLET gehört. Als Solver-Preset dient `mumps-gmres`. Hiermit kann auch gleich überprüft werden, ob die PETSc-interne MUMPS-Schnittstelle korrekt angesprochen wird. Die Kontrollvariablen zur Validierung sind in Abbildung 5.1 zu sehen.

### 5.2 Kohärenz-Test zur Berechnung der Jacobimatrix

Im Zuge der Umstellung des NuT-Codes von FORTRAN auf C/C++, siehe Abschnitt 2.3, wurde ein Kohärenz-Test zur Jacobimatrixstruktur entwickelt. Dieser vergleicht die von ATHLET generierten Strukturinformationen zur Jacobimatrix mit den tatsächlichen Systemantworten aus der AFK-Routine, welche im ATHLET-Kontext die rechte Seite des Differentialgleichungssystems zur Zeitintegration stellt. Die Ergebnisse des Tests werden sowohl während der Rechnung über Standard-Output ausgegeben als auch in ATHLETs Output-Datei gespeichert.

Die Erstellung eines solchen Tests war nicht konkrete in ursprünglichen Arbeitsplan vorgesehen. Da dieser Test ein wertvolles Instrument bei der allgemeinen Verifikation und Validierung von NuT ist, wurde er mit Priorität als Teil der allgemeinen Validierungsaktivitäten implementiert.



**Abb. 5.2** Verlauf der Dampftemperatur in CV 283 für verschiedene Löser, vor (old) und nach (new) Behebung der Diskrepanz. Das Attribut *Static* bezieht sich auf eine TOP-statische vorgehensweise, siehe hierzu Unterabschnitt 3.2.2.1.

### 5.2.1 Motivation

Die Motivation zur Entwicklung des Tests ergab sich aus dem instabilen Verhalten einer Temperaturvariablen während eines Simulationslaufes zum Datensatz PWR-3D-Sym.in<sup>1</sup>, siehe grüne Linie in Abbildung 5.2. Hintergrund des Laufes war das Austesten einer modifizierten Datenstruktur zur Jacobimatrix in NuT. Anschließende ausführliche Untersuchungen ergaben, dass Diskrepanzen zwischen Strukturinformationen zur Jacobimatrix und dem Verhalten von AFK vorhanden waren. Mithilfe des Tests wurden diese lokalisiert sowie sichtbar gemacht und der Bug in ATHLET anschließend behoben. Entsprechend stabil ergab sich daraufhin wieder das Simulationsverhalten, siehe abermals Abbildung 5.2.

### 5.2.2 Mögliche Auswirkungen einer Kohärenz-Diskrepanz

Um die Auswirkungen einer Diskrepanz zu verdeutlichen, bedarf es zunächst einiger Hintergrundinformationen.

#### Bestimmung der Jacobimatrix per Seeding

Approximationen der Jacobimatrix werden in ATHLET via finite Differenzen bestimmt. Sei ATHLETs AFK-Routine im Sinne einer mathematischen Funktion mit  $f$  bezeichnet. Dann werden an einer Stelle  $y$  mit passenden Richtungen  $s_i$  und skalaren Störungen  $\varepsilon_i$  für die

<sup>1</sup>Dieser wird mit ATHLET als Beispieldatensatz ausgeliefert.

Berechnung der Jacobimatrix  $J = \partial f(y)/\partial y$  die Differenzen

$$\frac{f(y + \varepsilon_i s_i) - f(y)}{\varepsilon_i} = J s_i + \mathcal{O}(\varepsilon_i) \quad (5.1)$$

gebildet. Da Funktionsauswertungen nicht-triviale Rechenzeit kosten, ist es wünschenswert, mit möglichst wenigen Richtungen  $s_i$  sämtliche Jacobimatrixinformationen zu erfassen. Das zugehörige Graphenfärbungsproblem ist NP-vollständig, weswegen auf Heuristiken zurückgegriffen wird. Eine solche liefert eine Matrix  $S \in \{0, 1\}^{n \times c}$  mit paarweise orthogonalen Spalten, so dass gilt

$$JS =: B \in \mathbb{R}^{n \times c} \quad \text{mit} \quad (J)_{kl} \neq 0 \quad \Rightarrow \quad \exists!(u, v) : (B)_{uv} = (J)_{kl}.$$

Obiger Prozess zur Bestimmung von  $S$  wird als *Seeding* bezeichnet. Die Matrix  $S$  heißt *Seed* oder *Seeding-Matrix*. Die Spalten  $s_i$  von  $S$  sind die *Seeding-Vektoren* und werden zur Berechnung der Differenzen in (5.1) herangezogen. Da  $J$  dünnbesetzt ist, kann davon ausgegangen werden, dass tatsächlich  $c \ll n$  gilt und somit viel Rechenzeit gespart werden kann im Vergleich zum trivialen Seed  $S_{triv} = I$ .

**Bemerkung 5.1.** Hinsichtlich der Störung  $\varepsilon_i$  sind die Verhältnisse in (5.1) vereinfacht dargestellt. Dies ist für die Argumentation an dieser Stelle jedoch nicht von Belang. Der Vollständigkeit halber sei auf /STE 17b, S. 42f/ verwiesen. Dort wird erläutert, warum es valide und sinnvoll ist, nicht nur mit einer skalaren Störung  $\varepsilon_i$  für die Richtung  $s_i$  zu arbeiten, sondern stattdessen  $\varepsilon_i$  ebenfalls vektoriell anzusetzen und als Gesamtstörung dann  $E_i s_i$  mit regulärem  $E_i := \text{diag}(\varepsilon_i)$  zu nutzen. Die finite Differenz ergibt sich dann zu  $E_i^{-1} (f(y + E_i s_i) - f(y))$ .

### Ein illustratives Beispiel zur Auswirkung

Sowohl ATHLET als auch NuT greifen auf die Heuristik von Courtis, Powell, and Reid (CPR) als Basisalgorithmus zu, wobei ATHLET eine Variante des smallest orderings nutzt, wohingegen NuT ein incidence degree ordering bemüht, siehe /COL 83/ für Details. Entscheidend ist, dass ausschließlich auf die Strukturinformation von  $J$  zur Bestimmung von  $S$  zugegriffen wird. Diese speist sich wiederum aus dem Aufbau des Netzwerks aus *control volumes* und *junctions*. Die AFK-Routine ist für das Seeding irrelevant, tritt aber in Erscheinung, wenn die Differenzen (5.1) zum Befüllen von  $J$  an den durch die Strukturinformation definierten Positionen genutzt werden. Genau dieses Aufeinandertreffen zweier Informationsstränge (Struktur aus Netzwerk versus AFK) birgt Fehlerpotential.

Sei zur Veranschaulichung die Strukturinformation zu  $J$  mittels

$$P = \begin{pmatrix} * & * \\ * & * & * \\ * & * & * \\ * & * \end{pmatrix}$$

gegeben.  $P$  ist symmetrisch gewählt, da  $J$  im ATHLET-Kontext stets struktursymmetrisch ist. Die Seeding-Heuristik liefere einen (optimalen) Seed

$$S = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}.$$

Weiter sei angenommen, dass aus Sicht der AFK-Routine die Struktur von  $J$  wie folgt sei

$$\tilde{P} = \begin{pmatrix} * & * & \# \\ * & * & * \\ * & * & * \\ \# & * & * \end{pmatrix}.$$

Zusatzpositionen im Vergleich zu  $P$  sind hierbei mit  $\#$  gekennzeichnet. Werden nun entsprechend der Seedmatrix die Differenzen (5.1) gebildet, resultiert dies in einer Approximation von  $J$  mit folgender Positionsinformation

$$P_d = \begin{pmatrix} * + \# & * & & \\ * & * & * & \\ * & * & * & \\ * & * & * + \# \end{pmatrix}. \quad (5.2)$$

Offensichtlich birgt dieses Resultat das Potential, problematisch zu sein. Der Einfluss der  $\#$ -Positionen verschiebt sich auf die Diagonale und bleibt an den ursprünglichen Positionen unberücksichtigt. Ersteres kann besonders heikel sein, da die Diagonalelemente einer Matrix über den Spur-Operator direkt mit den Eigenwerten  $\lambda$  der Matrix verknüpft sind.

$$\text{Es gilt } A = (a_{ij}) \in \mathbb{R}^{n \times n} : \quad \sum_{i=1}^n a_{ii} = \text{Spur}(A) = \sum_{i=1}^n \lambda_i.$$

Dies lässt sich z. B. sofort aus /MAC 93, Kor. 8.17/ ableiten. Sind die Werte an den  $\#$ -Positionen signifikant, könnte dies zu einer merklichen Änderung des Spektrums im Vergleich zur tatsächlichen Jacobimatrix führen. Die Stabilität des übergeordneten Differentialgleichungslösers, d. h. FEBE, wäre gefährdet. In der Regel geht dies dann mit einem massiven Einbruch der Zeitschrittweiten einher.

## Weitere Beispiele

Eine Kohärenz-Diskrepanz muss sich nicht auf die Diagonale auswirken oder auf irgend eine Position, die durch die Strukturinformation zu  $J$  gesetzt ist. Siehe

$$P = \begin{pmatrix} * & * \\ * & * \\ & * \end{pmatrix}, S = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{pmatrix}, \tilde{P} = \begin{pmatrix} * & * \\ * & * & \# \\ \# & * \end{pmatrix} \Rightarrow P_d = \begin{pmatrix} * & * \\ * + \# & * \\ & * \end{pmatrix} \quad (5.3)$$

sowie

$$P = \begin{pmatrix} * & * \\ * & * & * \\ * & * \end{pmatrix}, S = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \tilde{P} = \begin{pmatrix} * & * & \# \\ * & * & * \\ \# & * & * \end{pmatrix} \Rightarrow P_d = \begin{pmatrix} * & * \\ * & * & * \\ & * & * \end{pmatrix}. \quad (5.4)$$

## Grad der Auswirkungen

Es ist klar, dass eine Kohärenz-Diskrepanz nur dann möglicherweise Auswirkungen hat, wenn die AFK-Routine Positionen beeinflusst, die nicht durch die Strukturinformation abgedeckt sind (der umgekehrte Fall würde lediglich zu der Speicherung einer expliziten Null an einer entsprechenden Position führen, was keine Verfälschung nach sich zieht). Zusätzlich müssen die Werte an den Positionen  $\#$  signifikant gegenüber der Null sein, Beispiele (5.2)-(5.4), oder gegenüber den Werten, zu denen sie fälschlicherweise hinzugaddiert werden, Beispiel (5.2) und (5.3). Der Grad der Auswirkung wird jedoch noch durch weitere Faktoren beeinflusst.

Ein wichtiger Faktor im gesamten Integrationsprozess sind die Schrittweiten, mit denen die diskrete Zeitintegration vorangetrieben wird. Je kleiner die Schrittweite, desto weniger Einfluss hat die Jacobimatrix auf den (linear-)impliziten Integrationsverlauf. Folglich gilt Gleichtes für eine etwaige Kohärenz-Diskrepanz, da diese ausschließlich dafür sorgt, dass mit einer zusätzlich gestörten Approximation zur Jacobimatrix gearbeitet wird. Wenn die Schrittweiten lokal klein gewählt werden müssen, weil das System aktuell viel Dynamik aufweist, würde dies eine Diskrepanz verschleiern.

Die der linearen Algebra übergeordnete Integrationsmethode FEBE zählt zu den Rosenbrock-Wanner-Methoden, siehe /HAI 96, Kap. IV.7/. Diese linear-impliziten Verfahren haben die Eigenschaft, dass die Jacobimatrixapproximation keinen Einfluss auf die Konsistenzordnung hat. Solche Methoden erlauben also einen gewissen Grad an Verschmierrungen. Generell ist dies eine gute Eigenschaft. Sie rechtfertigt z. B. den Einsatz finiter Differenzen. Auf der anderen Seite wird auch hierdurch das Bemerkern einer Diskrepanz erschwert.

Obige Effekte treten gleichzeitig auf, was es in der Summe schwierig macht, den Grad der Auswirkungen einer Diskrepanz abzuschätzen beziehungsweise überhaupt zu erkennen, dass eine solche vorliegen könnte. Um dies eindeutig zu klären, wurde der Test entwickelt.

### 5.2.3 Funktionsweise und Anwendungsbereich des Kohärenz-Tests

Der Test ersetzt den berechneten Seed durch einen trivialen, setzt also  $S = I$ . Hiermit wird jede Komponente der finiten Differenz zu einem Seed-Vektor auf Nicht-Null-Einträge überprüft. Anschließend werden solche Einträge mit den Strukturinformationen abgeglichen. Bei einer Diskrepanz wird diese durch Netzwerkinformationen ergänzt ausgegeben und in ATHLETs Output-Datei gespeichert, siehe Abbildung 5.3.

Die Nutzung eines trivialen Seeds ist notwendig, um sicherzustellen, das eventuelle Kohärenz-Diskrepanzen entdeckt werden können. Da es sich zeitgleich um den uneffizientesten Seed handelt, nimmt der Test entsprechend Zeit ein.

Der Test steht sowohl für ATHLET/CD als auch für das Tandem ATHLET/CD/NuT zur Verfügung. Die Aktivierung des Tests ist sehr einfach und bedarf keiner Neucompilierung des Codes. Lediglich eine dedizierte Umgebungsvariable ist zu setzen. Simulationsläufe in der zugehörigen Umgebung führen dann automatisch den Test aus.

```
-----  
# check_jacobian: Updating Jacobian data with coherency test activated ...  
Perturbated value:  
    CV      1324  KEY 01 (pressure)  
with ITABs violating AFK-influence on:  
    Junction  1395  KEY 09 (liquid volume flow)  
    Junction  1395  KEY 10 (vapour volume flow)  
...  
# detected mismatches: 32  
-----
```

**Abb. 5.3** Ausschnitt der Ausgabe des Kohärenz-Tests.

### 5.2.4 Empfehlungen zum Einsatz des Tests

Es empfiehlt sich, den Test regelmäßig durchzuführen, wenn Code-Entwicklungen anstehen, welche die Netzwerkstruktur oder die Interaktion der Systemvariablen zum Ziel haben. Im Rahmen von QS-Tätigkeiten zu einem neuen Release/Patch sollte er stets auf einschlägige Datensätze angewandt werden.

Als weiteres Anwendungsgebiet sollte der Test im Kontext der *Continuous Integration* optional zur Verfügung gestellt werden, wenn Merge Requests in den Master-Branch oder einen anderen geschützten Branch anstehen. Hier sollte von Fall zu Fall vom Entwicklungsteam entschieden werden.

### **5.3 Entwicklung dedizierter Programme zum Austesten einzelner Funktionalitäten**

Es wurden drei Testprogramme zur Überprüfung der Schnittstellen sowie wichtiger NuT-Funktionalitäten definiert. Vom Aufbau her gleich decken die Programme die drei zu NuT 2.0 kompatiblen Schnittstellensprachen ab: FORTRAN, C und C++. Zum Einsatz kommt hier das Generator-Programm, welches bereits in Abschnitt 2.3 diskutiert wird, um automatisiert die passenden Modul- bzw. Header-Files zur Verfügung zu stellen. Die Programme agieren gegenüber NuT als Host-Applikation, analog zu den AC<sup>2</sup>-Simulationsprogrammen. Inhaltlich orientieren sie sich am NuT-Einzeltest aus Abschnitt 5.1 mit dem Zusatz, auch den Seeding-Prozess zu überprüfen und den berechneten Seed auszugeben. Alle wesentlichen NuT-Funktionalitäten sind somit in allen drei unterstützten Sprachen abgedeckt.

Die Testprogramme sind in der Compilierkonfiguration zu NuT berücksichtigt und stehen jedem Entwickler nach dem Compilieren und Linken direkt zur Verfügung. Die Ausgabe ist selbsterklärend. Es bedarf lediglich eines klassischen MPI-Aufrufes wie

```
mpiexec -n 1 nut_host_ex0X : -n 1 nut_worker
```

mit  $X \in \{1, 2, 3\}$ .

## 6 Fazit und Ausblick

Das Numerical Toolkit hat sich mit der Release-Version 1.0 als wertvolles Hilfsmittel innerhalb des AC<sup>2</sup>-Kontextes etabliert. Nicht nur Entwickler sondern auch Endnutzer können nun bei der Benutzung von ATHLET und ATHLET-CD (ATHLET/CD) von den Vorteilen skalierbarer und moderner Algorithmen zur numerischen linearen Algebra profitieren.

Parallel zur Fertigstellung der Release-Version wurde im Rahmen des Projektes an neuen Features gearbeitet. Es wurden nicht nur die Strukturen für das Lösen der Gleichungssysteme (1.2) während der transienten Phase verfeinert und bestätigt, Abschnitt 2.2 und Abschnitt 3.2.1, sondern auch Hilfsmittel für die Startrechnung, Abschnitt 3.3, sowie für das Restart-Szenario in ATHLET zur Verfügung gestellt, Abschnitt 3.6.

Mit dem Umschreiben des Toolkits von FORTRAN nach C(++) hat das Numerical Toolkit aus software-technischer Sicht einen großen Sprung nach vorn getan. Entsprechend der weitläufigen Anpassungen, siehe Abschnitt 2.3, hat NuT in seiner Entwicklung den Versionsstand 2.0 erreicht. Für interne Zwecke kann diese Version bereits genutzt werden. Für den Endnutzer wird dies mit dem nächsten Major-Release von AC<sup>2</sup> der Fall sein.

Die Arbeiten am Numerical Toolkit sind nicht nur NuT oder den Interaktionen der Host-Applikation mit dem Toolkit zugute gekommen. Es wurde im Zuge der Validierung der Modifikationen auch ein leistungsstarkes Werkzeug entwickelt, welches die Abbildung der Netzwerkstruktur auf Matrixstrukturen überwacht und mit dem Verhalten der Differentialgleichungs-Dynamik abgleicht, siehe Abschnitt 5.2. Dieser Kohärenz-Test ist unabhängig von NuT im Rahmen der allgemeinen ATHLET/CD-Entwicklung einsetzbar. Mittels des Tests konnten bereits verschiedene Unstimmigkeiten sowohl im Bestands-Code als auch in Neuentwicklungen aufgedeckt und anschließend behoben werden. Somit ergibt sich ein wertvolles Mittel im Rahmen eines allgemeinen Vorgehens zum *Continous Integration* und *Continous Delivery*.

Effizienter und robuster nach innen und außen zeigt sich der NuT-Code flexibel und weitestgehend compilerunabhängig in seinen Schnittstellen. Dies machte sich bereits für obige Funktionalitätserweiterungen bezahlt und lässt auch für zukünftige Arbeiten eine exzellente Ausgangsposition zu. Dies gilt insbesondere für die Arbeiten, welche im Rahmen des Projektes bisher nur auf Konzept-Ebene durchgeführt werden konnten, siehe Abschnitt 3.4, oder die aufgrund von unerwarteten Umständen jenseits der Zuständigkeit dieses Projektes nicht vollständig in die Praxis umgesetzt werden konnten, siehe z. B. Abschnitt 4.2.3.

Im Ausblick zeigt sich, dass mit den Projektarbeiten eine gute Basis für weiterführende Arbeiten geleistet ist. Build-Prozesse wurden verbessert, Schnittstellen standardisiert und Funktionalitäten erweitert. Der NuT-Entwicklungsprozess kann hierbei als Orientierung für andere Entwicklungen im AC<sup>2</sup>-Rahmen genutzt werden. Ein wesentliches Augenmerk für die nahe Zukunft sollte eine konsequente Weiterführung der Vereinheitlichung und Modernisierung von Build-Prozessen sein. Dies ermöglicht nicht nur die praktische Umsetzung bisher ausstehender Arbeiten, sondern fördert ganz allgemein eine gesunde Weiterentwicklung des AC<sup>2</sup>-Programmpaktes, so dass sich bestmöglich auf zukünftige Ansprüche an AC<sup>2</sup> vorbereitet werden kann.

## Literaturverzeichnis

- /AME 20/ P.R. Amestoy u. a. *MUMPS Web page*. <http://mumps.enseeiht.fr>; Zugriff 27. März. 2020 (siehe S. 10).
- /ARN 19/ S. Arndt u. a. *COCOSYS 3.0 – User Manual*. Teil der Dokumentation zur COCOSYS Software. GRS-P-3 / Vol. 1. Mai 2019 (siehe S. 2).
- /AUS 19/ H. Austregesilo u. a. *ATHLET 3.2 – Models and Methods*. Teil der Dokumentation zur ATHLET Software. GRS-P-1 / Vol. 4 Rev. 5. Feb. 2019 (siehe S. 2, 26).
- /BAL 20/ S. Balay u. a. *PETSc Web page*. <https://www.mcs.anl.gov/petsc>; Zugriff 26. März. 2020 (siehe S. 10).
- /BRO 17/ J. Brown. *[petsc-users] Possible to recover ILU(k) from hypre/pilut?* <https://lists.mcs.anl.gov/pipermail/petsc-users/2017-November/033935.html>; Zugriff 04. Oktober, 2018. Nov. 2017 (siehe S. 11).
- /BUT 64/ J. C. Butcher. “Implicit Runge-Kutta Processes”. In: *Math. Comp.* 18.85 (Jan. 1964), S. 50–64 (siehe S. 29).
- /BUT 16/ J. C. Butcher. *Numerical Methods for Ordinary Differential Equations*. 3. Aufl. John Wiley & Sons, Ltd., 2016 (siehe S. 3, 26).
- /BUT 90/ J. C. Butcher und J. R. Cash. “Towards Efficient Runge–Kutta Methods for Stiff Systems”. In: *SIAM J. Numer. Anal.* 27.3 (Juni 1990), S. 753–761 (siehe S. 30).
- /BUT 98/ J. C. Butcher und M. T. Diamantakis. “DESIRE: diagonally extended singly implicit Runge—Kutta effective order methods”. In: *Num. Alg.* 17 (1998), S. 121–145 (siehe S. 33).
- /CHE 98/ D. J. L. Chen. “The Effective Order of Singly-Implicit Methods for Stiff Differential Equations”. Diss. University of Auckland – Department of Mathematics, Juni 1998 (siehe S. 31).

- /COL 83/ T.F. Coleman und J.J. More. "Estimation of Sparse Jacobian Matrices and Graph Coloring Problems". In: *SIAM J. Numer. Anal.* 20 (1983), S. 187–209 (siehe S. 58).
- /DAV 06/ T. A. Davis. *Direct Methods for Sparse Linear Systems*. Fundamentals of Algorithms. SIAM, 2006 (siehe S. 3).
- /DOR 20/ R. Dorozhinskii. "Configuration of a linear solver for linearly implicit time integration and efficient data transfer in parallel thermo-hydraulic computations". Magisterarb. Technische Universität München – Fakultät für Informatik, März 2020 (siehe S. 17 ff., 73).
- /FAL 20/ R. Falgout, R. Li, U. Yang und D. Osei-Kuffuor. *HYPRE*. <https://github.com/hypre-space/hypre>; Zugriff 05. August, 2020. 2020 (siehe S. 10).
- /GEO 80/ A. George und J. W. H. Liu. "A Fast Implementation of the Minimum Degree Algorithm Using Quotient Graphs". In: *ACM Trans. Math. Softw.* 6.3 (1980), S. 337–358 (siehe S. 11).
- /GUS 97/ K. Gustafsson und G. Söderlind. "Control Strategies for the Iterative Solution of Nonlinear Equations in ODE Solvers". In: *SIAM J. Sci. Comput.* 18.1 (Jan. 1997), S. 23–40 (siehe S. 22).
- /HAI 93/ E. Hairer, S. P. Nørsett und G. Wanner. *Solving Ordinary Differential Equations I: Nonstiff Problems*. 2. Aufl. Bd. 8. Springer Series in Computational Mathematics. Springer Berlin Heidelberg, 1993 (siehe S. 3).
- /HAI 96/ E. Hairer und G. Wanner. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*. 2. Aufl. Bd. 14. Springer Series in Computational Mathematics. Springer Berlin Heidelberg, 1996 (siehe S. 3, 26, 28 ff., 46, 60).
- /HEC 20/ K. Heckmann. *Softwareentwicklung (TKP 03-05)*. GRS-Managementhandbuch Kap. 2.2.3.5 / Rev. 1. Mai 2020 (siehe S. 55).
- /HOF 81/ E. Hofer. "An  $A(\alpha)$ -Stable Variable Order ODE-Solver and its Application as Advancement Procedure for Simulations in Thermo- and Fluid-Dynamics".

- In: *Proceedings of the International Topical Meeting on Advances in Mathematical Methods for the Solution of Nuclear Engineering Problems*. München, Apr. 1981 (siehe S. 45).
- /HUN 03/ W. Hundsdorfer und J. G. Verwer. *Numerical Solution of Time-Dependent Advection-Diffusion-Reaction Equations*. Bd. 33. Springer Series in Computational Mathematics. Springer Berlin Heidelberg, 2003 (siehe S. 29).
- /INI 20/ Open Source Initiative. *2-Klausel-BSD-Lizenz*. <https://opensource.org/licenses/BSD-2-Clause>; Zugriff 08. Juli, 2020 (siehe S. 36).
- /INT 20/ Intel Corp. *Intel Xeon Prozessor E5-2680 v2 data sheet*. <https://ark.intel.com/content/www/de/de/ark/products/75277/intel-xeon-processor-e5-2680-v2-25m-cache-2-80-ghz.html>; Zugriff 05. August. 2020 (siehe S. 20).
- /JAC 20/ V. Jacht. *MMA – MPI for Multiple Applications*. <https://gitlab.com/nordfox/mma>; Zugriff 05. August. 2020 (siehe S. 35).
- /KAR 20/ G. Karypis und V. Kumar. *METIS Web page*. <http://glaros.dtc.umn.edu/gkhome/metis/metis/overview>; Zugriff 27. März. 2020 (siehe S. 10).
- /LER 19/ G. Lerchl u. a. *ATHLET 3.2 – User's Manual*. Teil der Dokumentation zur ATHLET Software. GRS-P-1 / Vol. 1 Rev. 8. Feb. 2019 (siehe S. 23, 27, 40).
- /LIP 12/ S.B. Lippman, J. Lajoie und B.E. Moo. *C++ Primer*. 5th. Addison-Wesley Professional, 2012 (siehe S. 3).
- /MAC 93/ W. Mackens und H. Voß, Hrsg. *Mathematik I*. HECO-VERLAG Aachen, 1993 (siehe S. 59).
- /MAR 08/ R.C. Martin. *Clean Code: A Handbook of Agile Software Craftsmanship*. 1. Aufl. USA: Prentice Hall PTR, 2008 (siehe S. 3).

- /MEI 11/ A. Meister. *Numerik linearer Gleichungssysteme*. 4. Aufl. Vieweg+Teubner, 2011 (siehe S. 3, 11).
- /MES 15/ Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard – Version 3.1*. <https://www.mpi-forum.org/docs/mpi-3.1/mpi31-report.pdf>. Juni 2015 (siehe S. 3).
- /NIS 20/ A. Nishida. *LIS – Library of Iterative Solvers for linear systems*. <https://github.com/anishida/lis>; Zugriff 28. Juli, 2020. 2020 (siehe S. 23).
- /OPE 18/ OpenMP Architecture Review Board. *OpenMP Application Programming Interface*. Online veröffentlicht. <https://www.openmp.org/wp-content/uploads/OpenMP-API-Specification-5.0.pdf>; Zugriff 04. August, 2020. Nov. 2018 (siehe S. 5).
- /PAS 13/ R. van der Pas. *OpenMP Tasking Explained*. Online veröffentlicht. <https://www.openmp.org//wp-content/uploads/sc13.tasking.ruud.pdf>; Zugriff 04. August, 2020. Nov. 2013 (siehe S. 10).
- /PRE 20/ T. Preston-Werner. *Semantic Versioning 2.0.0*. <https://semver.org/lang/de/>; Zugriff 04. August. 2020 (siehe S. 13, 55).
- /SAA 03/ Y. Saad. *Iterative Methods for Sparse Linear Systems*. 2. Aufl. SIAM, 2003 (siehe S. 3, 11).
- /SCH 20/ J. Scheuer. *Fortran Development Extensions (libfde)*. <https://github.com/Zorkator/libfde>; Zugriff 05. August. 2020 (siehe S. 36).
- /SCO 20/ C. Scott. *Professional CMake: A Practical Guide*. 7. Aufl. Online, 2020 (siehe S. 14).
- /ŠPE 16/ J. Špeh. *OpenMP: For & Scheduling*. Online veröffentlicht. <http://jakascorner.com/blog/2016/06/omp-for-scheduling.html>; Zugriff 04. August, 2020. Juni 2016 (siehe S. 10).
- /STE 17a/ T. Steinhoff. “Providing a Single Point Spectrum for Runge-Kutta Schemes of High Stage Order Based on Perturbed Collocation”. In: *AIP Conference*

*Proceedings*. Bd. 1863. <https://doi.org/10.1063/1.4992492>. Juli 2017 (siehe S. 30).

- /STE 17b/ T. Steinhoff und V. Jacht. *Ausbau und Modernisierung der numerischen Verfahren in den Systemcodes ATHLET, ATHLET-CD, COCOSYS und ASTEC*. Abschlussbericht, GRS - 469. Juli 2017 (siehe S. 1, 10, 20, 26, 28, 30, 49, 58).
- /STE 20/ T. Steinhoff und V. Jacht. *NuT – User’s Manual*. Teil der Dokumentation zur NuT Software. GRS-P-10 / Vol. 1 Rev. 2. Apr. 2020 (siehe S. 10, 21, 48, 55, 73).
- /TER 08/ V. A. Tereshonok u. a. *Kalinin-3 Coolant Transient Benchmark – Switching-off of One of the Four Operating Main Circulation Pumps at Nominal Reactor Power. Specification- First Edition*. NEA/NSC/DOC(2009)5. 2008 (siehe S. 9).
- /TIN 67/ W. F. Tinney und J. W. Walker. “Direct solutions of sparse network equations by optimally ordered triangular factorization”. In: *Proceedings of the IEEE* 55.11 (1967), S. 1801–1809 (siehe S. 48).

## Abbildungsverzeichnis

2.1	Prinzip dynamisch erstellter Threads per OpenMP .....	6
3.1	Doppelpuffer – Funktionsprinzip .....	17
3.2	Ausschnitt des Kommunikationsmusters zu Cube-64, Dimension: 100.657.	18
3.3	Performance-Vergleich zum Datensatz SPX .....	24
3.4	Beispielhafter valider $f$ -Auswertungspfad sowie invalide Pfadabschnitte	27
3.5	$f$ -Auswertungspfade sowie Extrapolationsschema zum FEBE-Integrator	27
3.6	$f$ -Auswertungspfade zum FiterRK32-Prototyp sowie Bereitstellung von Startwerten.....	31
3.7	Grundlegende Zustände und Übergänge in NuT, welche durch eine <i>State-Machine</i> abgebildet werden müssen.....	35
3.8	Delegieren der Kommunikationsinitialisierung an MMA.....	36
3.9	Paarweise Verknüpfung heterogener Prozessgruppen .....	37
3.10	Registrierung der gewünschten Kommunikatoren pro Prozessgruppe ....	38
3.11	Veröffentlichung der individuellen Kommunikatorlisten.....	39
3.12	Teilnahme an einem gewünschten Kommunikator .....	40
3.13	Kommunikatorstruktur inklusive Subkommunikatoren .....	41
3.14	Verklemmungsfreier Ressourcenzugriff dank MMA-Sortierung .....	42
3.15	Vergleich von Simulationsläufen im Restart-Szenario anhand exemplarischer Systemgrößen .....	44
5.1	Sample 1 – Kontrollvariablen zur Validierung, Solver-Preset <code>mumps-gmres</code>	56
5.2	Verlauf der Dampftemperatur in CV 283 für verschiedene Löser.....	57
5.3	Ausschnitt der Ausgabe des Kohärenz-Tests .....	61

## **Tabellenverzeichnis**

2.1	OpenMP-Performance für K3-2.....	9
3.1	Performance des Doppelpuffers zu Cube-64.....	18

## **Verzeichnis von Codeabschnitten**

Code 2.1	OpenMP-Performance-Messung I: Timer-Einbettung anhand der Beispielroutine nut_matCreateFromBlockCSR .....	7
Code 2.2	OpenMP-Performance-Messung II: Anweisungsdefinitionen anhand der Beispielroutine nut_matAllocateBlockCSR .....	7

# **A Anhang**

## **A.1 Masterarbeit Ravil Dorozhinskii**

Beiliegend findet sich die Abschlussarbeit von Herrn Ravil Dorozhinskii zur Erlangung des akademischen Grades Master of Science in Computational Science and Engineering durch die TU München, /DOR 20/. Die Arbeit wurde im vierten Quartal 2018 begonnen und am 07. März 2019 eingereicht. Der Titel lautet

*Configuration of a linear solver for linearly implicit time integration and efficient data transfer in parallel thermo-hydraulic computations*

Im Rahmen jener Abschlussarbeit wurden sowohl Effizienzuntersuchungen zur linearen Algebra während einer impliziten Zeitintegration thematisiert als auch die Verbesserung der Kommunikation zwischen zwei MPI-Prozessen mittels Puffer-Konzepten. Der Fokus der Arbeit lag hierbei auf dem ersten Thema, gab aber ebenfalls hinreichende und hilfreiche Informationen zum zweiten Themenbereich.

### **A.1.1 Effiziente lineare Algebra**

Die Resultate der Arbeit zeigen, dass die vom Numerical Toolkit angebotenen Solver-Presets, s. NuT-Anleitung /STE 20, Sec. 4.4/, bereits sehr gut verschiedene Anwendungsszenarien abdecken können. Für weitere Informationen siehe Abschnitt 3.2.1. Eine ausführliche Beschreibung findet sich in /DOR 20, Ch. 5/.

### **A.1.2 Implementation eines Doppelpuffers**

Hinsichtlich der Verbesserung der Kommunikation zwischen zwei MPI-Prozessen lag das zu untersuchende Problem in der Interprozesskommunikation zwischen ATHLET und NuT. Zwar konnte durch eine Doppelpufferstrategie die Kommunikationszeit merklich verbessert werden, jedoch fällt diese zu gering verglichen mit der Gesamtaufzeit aus, um den Aufwand zu rechtfertigen. Für weitere Details siehe /DOR 20, Ch. 6/ sowie Abschnitt 3.1.

### **A.1.3 Betreuung**

Die Betreuung wurde sowohl von der TU München als auch der GRS durchgeführt. Details hierzu sind der Arbeit zu entnehmen.



DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Computational Science and Engineering

**Configuration of a linear solver for linearly  
implicit time integration and efficient data  
transfer in parallel thermo-hydraulic  
computations**

**Ravil Dorozhinskii**





DEPARTMENT OF INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Computational Science and Engineering

**Configuration of a linear solver for linearly  
implicit time integration and efficient data  
transfer in parallel thermo-hydraulic  
computations**

**Konfiguration eines Lösen für  
linear-implizite Zeitintegration und  
effizienter Datentransfer in parallelen  
thermohydraulischen Berechnungen**

Author:

Ravil Dorozhinskii

Supervisor:

Prof. Dr. Thomas Huckle

Advisors:

Dr. rer. nat. Tim Steinhoff, Dr. rer. nat. Tobias Neckel

Submission Date:

07 March 2019



I confirm that this master's thesis is my own work and I have documented all sources and material used.

Munich, 07 March 2019

Ravil Dorozhinskii

# Abstract

**Key words.** *numerical time integration, direct sparse linear methods, multifrontal methods, MUMPS, BLAS, parallel performance, distributed-memory computations, multi-threading, MPI, OpenMP, non-blocking communication*

An application of linearly implicit methods for time integration of stiff systems of ODEs results in solving sparse systems of linear equations. An optimal selection and configuration of a parallel linear solver can considerably accelerate the time integration process. A comparison of iterative and direct sparse linear solvers has shown that direct ones are the most suitable for this purpose because of their natural robustness to ill-conditioned linear systems that can occur during numerical time integration. Testing of different direct sparse solvers applied to systems generated by ATHLET software has revealed that MUMPS, an implementation of the multifrontal method, performs better than the others in terms of the overall parallel execution time.

In this study, we have mainly focused on configuring MUMPS with the aim of improving parallel performance of the solver for thermo-hydraulic computations within a single node of GRS compute-cluster. However, the overall approach, proposed in the study, may be considered as a general framework for a selection and adaptation of a linear sparse solver for solving problem-specific systems of linear equations on distributed-memory machines.

Additionally, we have shown that an intelligent application of non-blocking MPI communication in some parts of the existing thermo-hydraulic simulation code, ATHLET, can additionally solve issues of inefficient data transfer preserving the current software design and implementation without drastic changes of the source code.

# Acronyms

**ATHLET** Analysis of THermal-hydraulics of LEaks and Transients.

**BLAS** Basic Linear Algebra Subprograms.

**GEMM** GEneral Matrix-matrix Multiplication (BLAS subroutine).

**GETRF** GEneral TRiangular Factorization (LAPACK subroutine).

**GMRES** General Minimum RESidual method.

**GRS** Gesellschaft für anlagen- und ReaktorSicherheit gGmbH.

**LAPACK** Linear Algebra PACKAGE.

**LRZ** Leibniz-RechenZentrum (Leibniz Supercomputing Centre of the Bavarian Academy of Sciences and Humanities).

**MPI** Message Passing Interface.

**MUMPS** MULTifrontal Massively Parallel sparse direct Solver.

**NUMA** Non-Uniform Memory Access.

**NUT** NUmerical Toolkit.

**ODE** Ordinary Differential Equation.

**OpenMP** Open Multi-Processing library.

**PDE** Partial Differential Equation.

**PETSc** Portable, Extensible Toolkit for Scientific Computation.

**ScaLAPACK** Scalable Linear Algebra PACKAGE.

**TRSM** TRiangular Solver with Multiple right-hand sides (LAPACK subroutine).

# Glossary

**HW1** Hardware installed on GRS cluster.

**HW2** Hardware installed on a LRZ CoolMUC-2 Linux cluster.

# List of Figures

2.1.	One-dimensional finite volume formulation of a thermo-hydraulic problem in ATHLET . . . . .	5
2.2.	A general view of the 6-stage W-method implemented in ATHLET . . . . .	5
2.3.	An example of NUT process groups . . . . .	7
2.4.	ATHLET-NUT software coupling . . . . .	8
4.1.	An example of pinning 5 MPI processes with 2 OpenMP threads per process in case of HW1 hardware . . . . .	15
5.1.	An example of a sparse matrix and its Cholesky factor . . . . .	22
5.2.	An elimination tree of matrix $A$ of the example depicted in Figure 5.1 . . . . .	23
5.3.	Information flow of the multifrontal method . . . . .	25
5.4.	An example of matrix postordering . . . . .	27
5.5.	An example of an efficient data treatment during matrix factorization using stacking . . . . .	28
5.6.	An example of a supernodal elimination tree . . . . .	28
5.7.	Examples of tree-task parallelism . . . . .	30
5.8.	Theoretical speed-up of models 1 and 2 . . . . .	31
5.9.	Comparisons of parallel performance of MUMPS, PaStiX and SuperLU_DIST libraries during 5 point-stencil Poisson matrix (1000000 equations) factorizations . . . . .	38
5.10.	An example of static and dynamic scheduling in MUMPS . . . . .	42
5.11.	An influence of different fill reducing algorithms on parallel factorizations of <i>pwr-3d</i> , <i>cube-5</i> , <i>k3-2</i> and <i>cube-64</i> matrices . . . . .	44
5.12.	An influence of different fill reducing algorithms on parallel factorizations of <i>k3-18</i> and <i>cube-645</i> matrices . . . . .	45
5.13.	Profiling of MUMPS-ParMetis configuration applied to parallel factorizations of relatively small matrices . . . . .	46
5.14.	Comparisons of <i>close</i> and <i>spread</i> pinning strategies applied to parallel factorizations of <i>pwr-3d</i> and <i>cube-64</i> matrices . . . . .	49
5.15.	Comparisons of <i>close</i> and <i>spread</i> pinning strategies applied to parallel factorizations of <i>cube-645</i> and <i>k3-18</i> matrices . . . . .	50

5.16. Comparisons of <i>close</i> and <i>spread</i> pinning strategies applied to parallel factorizations of <i>PFlow_742</i> and <i>CurlCurl_3</i> matrices . . . . .	51
5.17. One dimensional block column distribution of a type 2 node in MUMPS . . . . .	54
5.18. An example of type 2 node factorization implemented in MUMPS . . . . .	54
5.19. Software dependencies between Netlib libraries . . . . .	55
5.20. Comparisons of parallel factorization of GRS matrix set performed on HW1 machine using MUMPS solver linked to different BLAS implementations . . . . .	58
5.21. Comparisons of parallel factorizations of GRS and SuiteSparse matrix sets performed on HW1 machine using MUMPS solver linked to different BLAS implementations . . . . .	59
5.22. Anomalies of parallel executions of MUMPS-OpenBLAS configuration during factorizations of large-sized GRS matrices . . . . .	64
5.23. Thread conflicts of MUMPS-OpenBLAS configuration detected during parallel factorization of matrix <i>k3-18</i> . . . . .	64
5.24. Comparisons of parallel factorizations of small- and middle-sized GRS matrices between applications of the default and optimal MUMPS configurations . . . . .	69
5.25. Comparisons of parallel factorizations of large-sized GRS matrices between applications of the default and optimal MUMPS configurations . .	70
6.1. An example of matrix coloring and compression . . . . .	73
6.2. An example of an efficient Jacobian matrix partitioning . . . . .	75
6.3. A column-length distribution of the example depicted in Figure 6.2 . .	75
6.4. <i>Accumulator</i> concept . . . . .	78
6.5. Technical characteristics of HW1 hardware interconnection . . . . .	79
6.6. An application of the <i>accumulator</i> concept to the example depicted in Figure 6.3 . . . . .	80
6.7. A part of <i>cube-64</i> communication pattern . . . . .	81
6.8. Comparisons of the benchmarks running a recorded part of <i>cube-64</i> communication pattern between two sockets of a node . . . . .	85
6.9. A comparison of <i>BM1</i> benchmark with the original ATHLET-NUT implementation running a recorded part of <i>cube-64</i> communication pattern between two compute-nodes . . . . .	86
A.1. Sparsity patterns of SuiteSparse matrix set . . . . .	92
C.1. An influence of different fill reducing algorithms on parallel factorizations of <i>torso3</i> , <i>consph</i> , <i>CurlCurl_3</i> and <i>x104</i> matrices . . . . .	97

C.2. An influence of different fill reducing algorithms on parallel factorizations of <i>cant</i> and <i>memchip</i> matrices . . . . .	98
D.1. Comparisons of <i>close</i> and <i>spread</i> pinning strategies applied to parallel factorizations of <i>cube-64</i> and <i>torso3</i> matrices . . . . .	100
D.2. Comparisons of <i>close</i> and <i>spread</i> pinning strategies applied to parallel factorizations of <i>consph</i> and <i>memchip</i> matrices . . . . .	101
E.1. Comparisons of parallel factorizations of <i>cant</i> , <i>consph</i> , <i>memchip</i> and <i>x104</i> matrices performed on HW1 machine using MUMPS solver linked to different BLAS implementations . . . . .	103
E.2. Comparisons of parallel factorizations of <i>pkustk10</i> , <i>CurlCurl_3</i> and <i>Geo_1438</i> matrices performed on HW1 machine using MUMPS solver linked to different BLAS implementations . . . . .	104

# List of Tables

1.1.	A list of software developed by GRS . . . . .	2
4.1.	GRS matrix set . . . . .	13
4.2.	SuiteSparse matrix set . . . . .	13
4.3.	Hardware specifications . . . . .	14
5.1.	A list of parallel preconditioning algorithms available in PETSc . . . . .	20
5.2.	Theoretical speed-up of models 1 and 2 . . . . .	31
5.3.	A list of direct sparse linear solvers adapted for distributed-memory computations . . . . .	36
5.4.	Comparisons of parallel performance of <i>cube-5</i> matrix factorizations using MUMPS, PasTiX and SuperLU_DIST solvers with their default parameter settings . . . . .	37
5.5.	Comparisons of parallel performance of <i>cube-64</i> matrix factorizations using MUMPS, PasTiX and SuperLU_DIST solvers with their default parameter settings . . . . .	37
5.6.	Comparisons of parallel performance of <i>k3-18</i> matrix factorizations using MUMPS, PasTiX and SuperLU_DIST solvers with their default parameter settings . . . . .	38
5.7.	Assignment of GRS matrices to specific fill-in reducing algorithms . . . . .	47
5.8.	Assignment of SuiteSparse matrices to specific fill-in reducing algorithms . . . . .	47
5.9.	Comparisons of MUMPS parallel performance at the saturation points in case of factorization of GRS matrix set . . . . .	52
5.10.	Comparisons of MUMPS parallel performance at the saturation points in case of factorization of SuiteSparse matrix set . . . . .	53
5.11.	Commercial and open source BLAS libraries . . . . .	56
5.12.	Comparisons of different MUMPS-BLAS configurations applied to GRS matrix set . . . . .	60
5.13.	Comparisons of different MUMPS-BLAS configurations applied to SuiteSparse matrix set . . . . .	60
5.14.	Comparisons of different hybrid MPI/OpenMP modes used for parallel factorization of GRS matrix set on HW1 . . . . .	65

---

*List of Tables*

---

5.15.	Comparisons of different hybrid MPI/OpenMP modes used for parallel factorization of GRS matrix set on HW2 . . . . .	65
5.16.	Comparisons of different hybrid MPI/OpenMP modes used for parallel factorization of SuiteSparse matrix set on HW1 . . . . .	66
5.17.	Comparisons of different hybrid MPI/OpenMP modes used for parallel factorization of SuiteSparse matrix set on HW2 . . . . .	66
6.1.	Time reduction of data transfers with respect to the original implementation in case of execution of <i>cube-64</i> communication pattern . . . . .	84
B.1.	Comparisons of parallel performance of <i>pwr-3d</i> matrix factorizations using MUMPS, PaStiX and SuperLU_DIST libraries with their default parameter settings . . . . .	94
B.2.	Comparisons of parallel performance of <i>k3-2</i> matrix factorizations using MUMPS, PaStiX and SuperLU_DIST libraries with their default parameter settings . . . . .	94
B.3.	Comparisons of parallel performance of <i>cube-645</i> matrix factorizations using MUMPS, PaStiX and SuperLU_DIST libraries with their default parameter settings . . . . .	95

# Listings

5.1.	An example of usage of the PETSc built-in sparse direct linear solver as a sub-preconditioner for the Block Jacobi preconditioning algorithm . . . . .	21
5.2.	Pseudocode of the iterative refinement method . . . . .	33
5.3.	An example of setting <i>spread</i> process pinning using advanced OpenMPI command line options . . . . .	53
6.1.	Pseudocode of the original ATHLET-NUT coupling: ATHLET part . . . . .	76
6.2.	Pseudocode of the original ATHLET-NUT coupling: NUT part . . . . .	77
6.3.	Pseudocode of an auxiliary <i>Accumulator</i> class . . . . .	82
6.4.	Pseudocode of a modified client side of the benchmark . . . . .	83

# Contents

<b>Abstract</b>	iii
<b>Acronyms</b>	iv
<b>Glossary</b>	v
<b>List of Figures</b>	vi
<b>List of Tables</b>	ix
<b>Listings</b>	xi
<b>1. Introduction</b>	1
<b>2. Overview of ATHLET and NUT software</b>	3
2.1. ATHLET . . . . .	3
2.2. NUT . . . . .	6
2.3. ATHLET-NUT coupling . . . . .	6
<b>3. Problem Statement</b>	9
<b>4. Methodology and Experimental Setup</b>	12
<b>5. Configuration of a sparse linear solver</b>	16
5.1. Overview of Sparse Linear Solver Types . . . . .	16
5.1.1. Iterative Methods . . . . .	16
5.1.1.1. Theory Overview . . . . .	16
5.1.1.2. Parallelization Aspects . . . . .	18
5.1.1.3. Preconditioners . . . . .	19
5.1.2. Direct Sparse Methods . . . . .	21
5.1.2.1. Theory Overview . . . . .	21
5.1.2.2. Parallelization Aspects . . . . .	29
5.1.2.3. Threshold Pivoting and Solution Refinement . . . . .	32
5.1.3. Results and Conclusion . . . . .	34

*Contents*

---

5.2.	Selection of a Sparse Direct Linear Solver . . . . .	35
5.3.	Overview of MUMPS Library . . . . .	39
5.4.	Configuration of MUMPS Library . . . . .	43
5.4.1.	Fill Reducing Reorderings . . . . .	43
5.4.2.	MPI Process Pinning . . . . .	48
5.4.3.	Optimized BLAS Implementations . . . . .	53
5.4.4.	Hybrid MPI/OpenMP Computing . . . . .	61
5.5.	Results . . . . .	68
5.6.	Conclusion . . . . .	70
<b>6.</b>	<b>Improvement of ATHLET-NUT Communication</b>	<b>73</b>
6.1.	Jacobian Matrix Compression . . . . .	73
6.2.	Accumulator Concept . . . . .	78
6.3.	Benchmark and Test Data . . . . .	80
6.4.	Results . . . . .	83
6.5.	Conclusion . . . . .	87
<b>Appendices</b>		<b>89</b>
<b>A.</b>	<b>Sparsity Patterns of Matrix Sets</b>	<b>90</b>
<b>B.</b>	<b>Selection of a Sparse Direct Linear Solver</b>	<b>93</b>
<b>C.</b>	<b>Fill Reducing Reorderings</b>	<b>96</b>
<b>D.</b>	<b>MPI Process Pinning</b>	<b>99</b>
<b>E.</b>	<b>Optimized BLAS Libraries</b>	<b>102</b>
<b>Bibliography</b>		<b>105</b>

# 1. Introduction

Nowadays, nuclear energy is one of the main sources of electricity. It comes from splitting atoms in a reactor which, as a result, heats water up to the point where it is converted into pressurized steam. In its turn, the steam rotates turbines which, finally, produce electricity. According to the recent estimations, thermal efficiency of modern nuclear power plants lies in the range of 35-45% which is comparable to conventional fossil fueled power plants [28]. In spite of considerable initial investments, nuclear power plants have low operating costs and long service life which make them particularly cost effective.

In recent years, nuclear power plants have become an attractive means of power generation because of relatively low emission of carbon dioxide. As a result, a level of green house gase emissions to the atmosphere and thus the contribution of nuclear power plants to the global warming are relatively small [45].

Today, nuclear power plants generate almost 30% of the electricity produced in the European Union (EU). There are almost 130 nuclear reactors in operation in 14 EU countries, namely: Belgium, Bulgaria, Czech Republic, Finland, France, Germany, Hungary, Netherlands, Romania, Slovakia, Slovenia, Spain, Sweden, and the United Kingdom [17].

The main problem associated with nuclear power is radioactive waste which is extremely dangerous for people and the environment and has to be carefully looked after for several thousand years after utilization. Any accident in a plant can lead to grave consequences at a scale similar to the Chernobyl disaster. For this reason, nuclear safety is one of the most important topics in this area. It demands a huge amount of testing and analysis to be performed before and during an operation of a nuclear power plant in order to predict any possibility of unwanted outcomes and devise preventive measures against such accidents. The topic has become even more prominent since 2011 Fukushima accident. In response to the disaster, numerous stress tests were conducted to measure the ability of the EU nuclear industry to withstand any kind of natural disaster [17].

---

## 1. Introduction

---

Since 1977, Gesellschaft für Anlagen- und Reaktorsicherheit (GRS) has been the main German scientific research institute in the field of nuclear safety and radioactive waste management [23]. Today, the organization carries out advanced research and analysis in the field of reactor safety, radioactive waste management as well as radiation and environmental protection [23]. Due to the inability to create various nuclear accident test scenarios, which by their very nature could be catastrophic, GRS develops and provides numerous simulation software products to cope with this problem. A short description of the main software packages developed by GRS is provided in Table 1.1.

Name	Description
ATHLET	Thermohydraulic safety analyses for the primary circuit of LWRs
ATHLET-CD	Analyses of accidents with core meltdown and fission product release for LWRs
ATLAS	Analysis simulator for interactive handling and visualisation of several computer codes
COCOSYS	Analyses of severe incidents in the containment of LWRs
DORT/TORT	Solution of time-dependant neutron transport equations for 2D/3D transients analyses
QUABOX/CUBBOX	3-D neutron kinetics core model
SUSA	Uncertainty and sensitivity analyses
TESPA-ROD	Core rod code for design basis accidents
NUT	Container of various numerical tools and algorithms

Table 1.1.: A list of software developed by GRS, [22], where LWR stands for a Light Water Reactor

The main focus of this thesis is dedicated to ATHLET and NUT software packages. The goal of the study is to identify the most compute-intensive parts of the ATHLET-NUT code and possibly accelerate its execution time.

The next chapter continues the introduction and gives a general overview of ATHLET-NUT purpose, design, architecture and coupling. The introduction ends with a clear exposition of the problem statement presented in Chapter 3 where the order of the remaining thesis is described in detail.

## 2. Overview of ATHLET and NUT software

### 2.1. ATHLET

Analysis of THermal-hydraulics of LEaks and Transients (ATHLET) software is developed by GRS for an analysis of the whole spectrum of operational conditions, incidental transients, design-basis accidents and beyond design-basis accidents without core damage anticipated for nuclear energy facilities [21]. The code provides specific models and methods to simulate many types of nuclear power plants, comprising current light water reactors (PWR<sup>1</sup>, BWR<sup>2</sup>, WWER<sup>3</sup>, HPCR<sup>4</sup>), advanced Generation III+ and IV reactors as well as SMRs<sup>5</sup> [21].

Physical processes inside of hydraulic circuits of light water reactors can be naturally described by a two-phase thermo-fluiddynamic model based on equations of conservation of mass, momentum and energy for liquid and vapor phases i.e. Equations 2.1 - 2.7, [7].

1. Liquid mass

$$\frac{\partial((1-\alpha)\rho_l)}{\partial t} + \nabla((1-\alpha)\rho_l\vec{w}_l) = -\psi \quad (2.1)$$

2. Vapor mass

$$\frac{\partial(\alpha\rho_v)}{\partial t} + \nabla(\alpha\rho_v\vec{w}_v) = \psi \quad (2.2)$$

3. Liquid momentum

$$\frac{\partial((1-\alpha)\rho_l\vec{w}_l)}{\partial t} + \nabla((1-\alpha)\rho_l\vec{w}_l\vec{w}_l) + \nabla((1-\alpha)p) = \vec{F}_l \quad (2.3)$$

4. Vapor momentum

$$\frac{\partial(\alpha\rho_v\vec{w}_v)}{\partial t} + \nabla(\alpha\rho_v\vec{w}_v\vec{w}_v) + \nabla(\alpha p) = \vec{F}_v \quad (2.4)$$

---

<sup>1</sup>Pressurized Water Reactor

<sup>2</sup>Boiling Water Reactor

<sup>3</sup>Water-Water Energetic Reactor

<sup>4</sup>High Power Channel-type Reactor

<sup>5</sup>Small Modular Reactor

### 5. Liquid energy

$$\frac{\partial \left[ (1-\alpha)\rho_l(h_l + \frac{1}{2}\vec{w}_l\vec{w}_l - \frac{p}{\rho_l}) \right]}{\partial t} + \nabla \left[ (1-\alpha)\rho_l\vec{w}_l(h_l + \frac{1}{2}\vec{w}_l\vec{w}_l) \right] = -p\frac{\partial(1-\alpha)}{\partial t} + E_l \quad (2.5)$$

### 6. Vapor energy

$$\frac{\partial \left[ \alpha\rho_v(h_v + \frac{1}{2}\vec{w}_v\vec{w}_v - \frac{p}{\rho_v}) \right]}{\partial t} + \nabla \left[ \alpha\rho_v\vec{w}_v(h_v + \frac{1}{2}\vec{w}_v\vec{w}_v) \right] = -p\frac{\partial\alpha}{\partial t} + E_v \quad (2.6)$$

### 7. Volume vapor fraction

$$\alpha = \frac{V_v}{V} \quad (2.7)$$

The notation is as follows:  $p$  - pressure of a liquid-vapor mixture,  $\psi$  - a mass source term,  $\vec{F}$  - an external composite force acting on a CV<sup>6</sup>,  $E$  - an external composite energy source term within a CV, subscripts  $l$  and  $v$  denote liquid and vapor phases, respectively.

Spacial integration of the conservation equations, System 2.1 - 2.7, is performed on basis of the finite volume method using one-dimensional problem formulation, Figure 2.1. Then, the system is transformed to an initial value problem of a system of non-autonomous Ordinary Differential Equations (ODEs) by means of certain additional mathematical transformations, see [7].

$$\frac{dy}{dt} = f(t, y), \quad t_0 \leq t \leq t_F \quad y(t_0) = y_0 \quad (2.8)$$

where  $y \in \mathbb{R}^n$  is a composite vector of variables,  $f$  is a non-linear function such that  $f : \mathbb{R} \times \mathbb{R}^n \supset \Omega \rightarrow \mathbb{R}^n$ .

According to *ATHLET Mod 3.1A – Models and Methods* [7], System 2.8 is stiff and thus must to be solved with an implicit solver. To avoid a high computational cost of a fully implicit method, ATHLET makes use of an extrapolation ansatz based on the linearly implicit Euler method. This may be interpreted as a six stage W-method where the exact Jacobian information is not required. However, fresh information of the Jacobian matrix during a numerical simulation greatly improves stability and robustness of the method. Therefore, several mechanisms have been implemented in ATHLET to closely monitor and, if it is required, to update a Jacobian matrix approximation using the finite difference method.

---

<sup>6</sup>CV - Control Volume

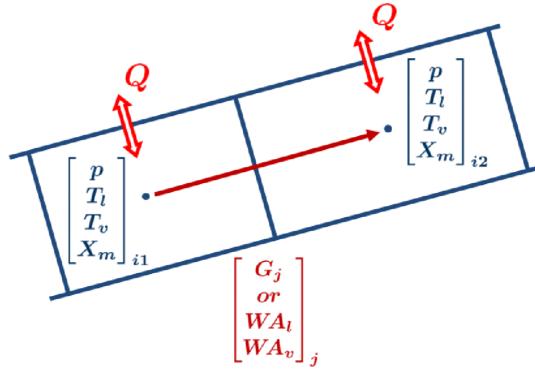


Figure 2.1.: One-dimensional finite volume formulation of a thermo-hydraulic problem in ATHLET, [44], where  $T_l$  - a temperature of liquid inside of a CV;  $T_v$  - a temperature of vapor inside a CV;  $X_m$  - mass quality;  $G$  - mass flow of a liquid-vapor mixture;  $W$  - a velocity component of a liquid-vapor mixture perpendicular to a CV boundary;  $A_l$  - an area occupied by the liquid fraction on a CV boundary,  $A_v$  - an area occupied by the vapor fraction on a CV boundary;  $Q$  - an external heat transfer

In the general case, a step of the W-method method, implemented in ATHLET, can be viewed as a sequence of six stages in the following way. Each stage uses the implicit Euler method and exactly one Newton's iteration to evaluate values of vector  $y$  at the next integration step  $h$  with different accuracy. Then, the obtained values are extrapolated, in the order shown in Figure 2.2, to achieve the third order of numerical

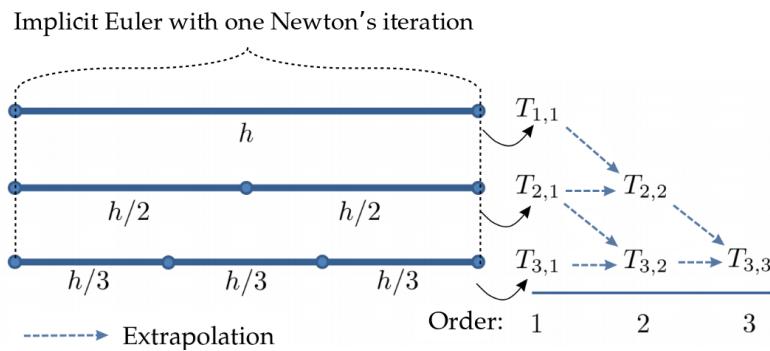


Figure 2.2.: A general view on the 6-stage W-method implemented in ATHLET (based on [44]), where  $T_{1,1} = y_0 + y_{1,0}$ ;  $T_{2,1} = y_0 + y_{2,0} + y_{2,1}$ ;  $T_{3,1} = y_0 + y_{3,0} + y_{3,1} + y_{3,2}$

integration. By and large, the algorithm can be expressed in a compact form of Equation 2.9.

$$(I - h_i J) \delta y_{ij} = h_i f(t_0 + j \cdot h_i, y_0 + \sum_{l=0}^{j-1} \delta y_{il}) + h^2 \frac{\partial f_0}{\partial t} \quad (2.9)$$

where  $J \approx \frac{\partial f}{\partial y}$  - an approximation of Jacobian matrix;  $i = 1, 2, 3$ ;  $j = 0, \dots, i - 1$ ;  $h_i = h/i$ .

## 2.2. NUT

NUmerical Toolkit (NUT) can be viewed as a container of various dense and sparse linear algebra subroutines which can run in parallel on distributed-memory machines. NUT design follows a paradigm of the *Adapter/Wrapper* pattern which provides a uniform common interface for its services to any application that follows a certain communication protocol. Currently, only ATHLET can communicate with NUT. However, more GRS applications are going to adopt the protocol in the future and will be able to operate with NUT as well. The approach, implemented in NUT, helps to achieve re-usability, flexibility and extensibility properties of the code.

Currently, NUT is based heavily on Portable, Extensible Toolkit for Scientific Computation (PETSc). It is one of the most widely used parallel numerical software libraries [47]. It includes a large suite of parallel linear and nonlinear equation solvers as well as a software-infrastructure to handle computations on distributed-memory machines by means of Message Passing Interface (MPI) and specific data structures. Fortunately, through a careful selection of the design pattern, NUT can be easily extended to provide an extra service or an external library access which has not been implemented in PETSc yet.

## 2.3. ATHLET-NUT coupling

Coupling of NUT with GRS tools is based on the client-server architecture where NUT acts as a server and the tools can be viewed as clients. Communication between two parts is done via MPI.

To provide a clear and concise external interface, NUT contains a client module called "NUT Plug-in". It can be considered as a socket, from the client side by analogy with

## 2. Overview of ATHLET and NUT software

---

the Transmission Control Protocol (TCP). The plug-in hides all MPI calls to the sever which considerably improves readability of the code.

In principle, NUT allows multiple clients to work concurrently with the server. To handle communication traffic, NUT splits the default MPI communicator at start-up time of the application into appropriate process groups, as it is shown in Figure 2.3.

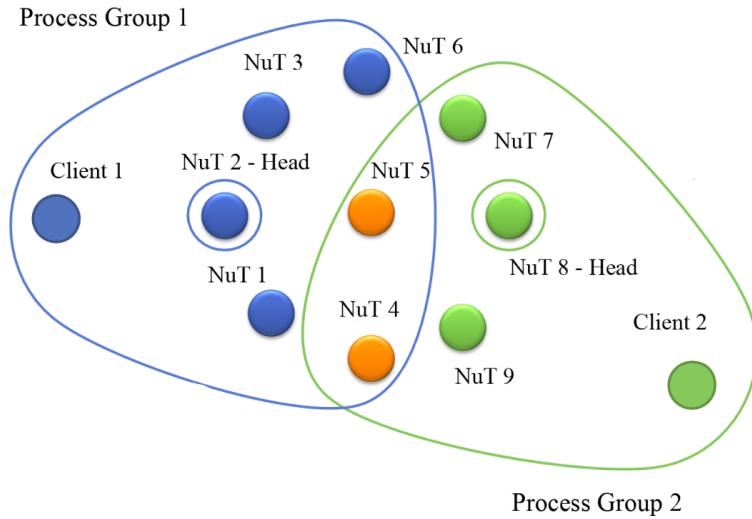


Figure 2.3.: An example of NUT process groups

The design of NUT allows sharing of some NUT-MPI processes among different process groups due to performance reasons i.e. a finite number of processing units on hardware. To resolve possible deadlocks, each process group has its own representative, called the head. Each client has two views on its respective group which is achieved by means of distinct MPI communicators. The first communicator is responsible for client-head communication whereas the second one allows the client to talk to any NUT process within the group.

A general view of client-server communication looks like a 3-way handshake in the following way: a client sends a request to the head which is a signal to reserve all compute-units of the group for an upcoming task. Having possessed the resources and prepared them for a specific service, the head notifies the client about a resource acquisition and the entire process group waits for data. Afterwards, the client sends data either to a specific NUT-process or to the entire group using the second communicator

## 2. Overview of ATHLET and NUT software

---

and waits for a result of the service. According to the current implementation of NUT, the communication between a client and server is synchronous i.e. a client gets blocked while waiting for a result from the server.

A general view on ATHLET-NUT software coupling is given in Figure 2.4 where ATHLET is responsible for marching the numerical time-integration process whereas NUT computes solutions of linear systems derived from Equation 2.9.

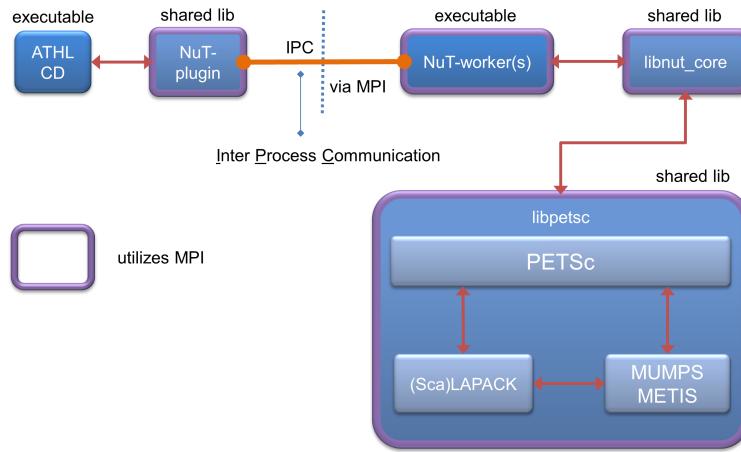


Figure 2.4.: ATHLET-NUT software coupling

Partial and full Jacobian matrix updates derived from finite differences are computed on the client side since only a client has access to the function  $f$  of Equation 2.8. Due to transformations of the underlying system of Partial Differential Equations (PDEs) and specifics of finite volume discretization, the Jacobian matrix is sparse and, therefore, ATHLET uses a matrix compression algorithm described in Section 6.1 to reduce an amount of Jacobian column evaluations. Having computed a matrix column, ATHLET immediately broadcasts it to its entire NUT process group by means of the 3-way handshake mechanism described above. It is worth mentioning that this approach allows to circumvent potential memory limits on the client side and thus store the entire sparse Jacobian matrix in a distributed fashion on the server. In other words, ATHLET never holds the entire Jacobian matrix in its memory; conversely, the matrix is distributed across multiple NUT processes according to the block-row distribution scheme induced by PETSc. In turn, NUT is waiting for the entire Jacobian matrix information from ATHLET and starts solving Systems 2.9 right after receiving the corresponding request from the client.

### 3. Problem Statement

Integration of a system of ODEs by means of W-methods involves solving several systems of linear equations. Equations 2.9 can be rewritten in a form 3.1, after grouping both the right- and left-hand sides in a single matrix and vector, respectively.

$$A_i \delta y_{ij} = b_{ij} \quad (3.1)$$

where  $A_i = (I - h_i J)$  is a  $n \times n$  nonsingular sparse matrix;  $\delta y_{ij}$  and  $b_{ij}$  are  $\mathbb{R}^n$  vectors.

According to the integration scheme, Figure 2.2, and definition of the method, each step of numerical integration requires to solve 6 linear systems with 3 distinct matrices, resulting from the Jacobian matrix by the corresponding shifts of the main diagonal. Therefore, the computational burden of the W-method mainly lies in both solving sparse linear systems and evaluations of non-linear function  $f$ , Equation 2.8. However, because of the time limit of the thesis, the main focus of this study is on solving Systems 3.1 efficiently on GRS computational cluster.

There exist two families of linear sparse solvers, namely: iterative and direct sparse methods. In the general case, execution time of any method, regardless of a solver family, is bounded by  $O(n^2)$  complexity due to matrix sparsity, where  $n$  is the number of equations in a system. However, constants in front of the factor  $n^2$  can vary significantly between the methods which explains differences in execution time. Additionally, it is important to mention that the families follow different approaches for solving sparse linear systems and are greatly different in details. Therefore, they possess different numerical properties. Among all properties, there are some which are particularly important for efficient execution of W-methods, namely:

- robustness of a method to treat, in particular, ill-conditioned systems
- parallel efficiency

These, above mentioned properties, can be treated as non-functional requirements to

### 3. Problem Statement

---

a sparse linear solver for efficient numerical time integration.

Finding solutions of sparse linear systems is a well-known and commonly occurring problem in the field of scientific computing and, therefore, numerous implementations of different kinds of linear solvers exist. However, the NUT project imposes some extra constraints due to the design philosophy adopted by GRS:

- open-source license
- direct interface to PETSc
- technical support and maintainability of a solver/package

In this study, we are primarily concerned with a selection and configuration of a sparse linear solver which can cover all above listed requirements.

This report is organized as follows. Chapter 4 provides information about methodology, data, software and hardware used in this study. Subsections 5.1.1 and 5.1.2 give an overview of the theory and parallelization aspects of iterative and sparse direct methods as well as discussions of some issues related to numerical solution accuracy. Then, in Subsection 5.1.3, we make a conclusion about which type of sparse linear solvers is well suited for numerical time integration governed by W-methods. In Section 5.2, a concrete implementation of a specific method is selected by means of testing. From Section 5.4 onwards, we perform configuration and adaptation of a solver for distributed-memory computations. At the end, Section 5.6 summarizes obtained results and makes a general conclusion with respect to data and compute environment provided by GRS.

An additional topic, considered in this study, is an improvement of ATHLET-NUT communication during Jacobian matrix transfers. As it was described in Section 2.3, ATHLET, the client, transfers a Jacobian matrix in a column-wise fashion. NUT, the server, treats each column transfer as a service and, therefore, each transfer passes through the 3-way handshake described in Section 2.3. Moreover, it is important to mention one more time, due to the current implementation of ATHLET-NUT coupling, client-server communication is blocking. In other words, ATHLET gets blocked till completion of a column transfer.

The main goal of Jacobian matrix compression, described in Section 6.1, is to minimize the number of perturbations of non-linear function  $f$  of Equation 2.8. Additionally it allows to reduce an amount of column transfers as well. Therefore, it improves the

### 3. Problem Statement

---

overall application performance from both computational and communication points of view. However, there are still some aspects to be considered.

Due to specifics of a matrix compression algorithm, described in Section 6.1, column lengths are decreasing between the first and last columns of a compressed Jacobian matrix form which, as a result, leads to unequal MPI message sizes.

In the last part of the study, we introduce a concept called *accumulator* which allows to transfer a compressed Jacobian matrix in equal chunks. This approach potentially solves three important problems at once. First of all, *accumulator* can help to get rid of small MPI messages which improves utilization of network bandwidth. Secondly, it helps to reduce an amount of synchronizations between a client and the server and, therefore, improves operation of NUT as the server. Lastly, it allows to apply non-blocking MPI communication on the client side and thus overlap Jacobian matrix transfers with computations.

In Section 6.1, we briefly describe the Jacobian matrix compression algorithm and the resulting ATHLET-NUT communication problem. In Section 6.2, we present and describe an algorithm which is supposed to resolve the problem. Section 6.3 provides a description of developed benchmarks and test data. Then, we discuss obtained results in Section 6.4. Finally, in Section 6.5, we provide a general conclusion of the performed part of the study and summarize the results.

## 4. Methodology and Experimental Setup

ATHLET is a CFD<sup>1</sup> tool designed for computer-based simulations of transient thermo-hydraulic problems where topology of hydraulic circuits can be changed during a numerical simulation. As a direct consequence, the Jacobian matrix can frequently change with respect to both numerical values and the matrix sparsity structure between time integration steps. Since ATHLET usually generates hundreds of matrices during a simulation, configuration of a linear solver in run-time becomes a time consuming and compute-expensive problem. Moreover, results of such dynamic solver configuration may be difficult to analyze and interpret.

In this study, a static solver configuration approach is used, instead. In other words, a solver is configured with only a small set of matrices, i.e. GRS matrix set, randomly saved during simulations of the most common GRS test scenarios. In the general case, it may lead to inaccurate conclusions, however, it is only one technically feasible approach.

Besides GRS matrix set, a second set was used for verification of testing results. The set, called SuiteSparse matrix set, was generated by downloading a dozen of matrices from SuiteSparse Matrix Collection [11], [12] where we tried to pick out different matrices with respect to both the number of equations  $n$  and matrix density i.e. ratio between the number of non-zero elements  $nnz$  and the number of equations in a system.

The main matrix properties as well as matrix sparsity patterns are shown in Tables 4.1, 4.2 and appendix A.

Approximations of condition numbers, shown in Tables 4.1 and 4.2, were computed using the Rayleigh–Ritz procedure. The reader can become familiar with the procedure in [30]. GMRES solver, configured with 1000 iteration steps before the restart, was applied to un-preconditioned systems to generate a Krylov subspace for each matrix. Then, the resulting Hessenberg matrices were used for approximating eigenspaces and the corresponding eigenvalues. The approximations should be treated as lower bounds since the algorithm overestimates the smallest eigenvalues.

---

<sup>1</sup>Computational Fluid Dynamics

#### 4. Methodology and Experimental Setup

---

Name	<i>n</i>	<i>nnz</i>	<i>nnz / n</i>	Approximate Condition Number	Structure
pwr-3d	6009	32537	5.4147	1.019e+07	SYMM-PTRN
cube-5	9325	117897	12.6431	1.592e+09	SYMM-PTRN
cube-64	100657	1388993	13.7993	7.406e+08	SYMM-PTRN
cube-645	1000045	13906057	13.9054	6.474e+08	SYMM-PTRN
k3-2	130101	787997	6.0568	1.965e+15	SYMM-PTRN
k3-18	1155955	7204723	6.2327	1.947e+12	SYMM-PTRN

Table 4.1.: GRS matrix set, where SYMM - symmetric; NON-SYMM - non-symmetric; SYMM-PTRN- non-symmetric but with symmetric sparsity pattern

Name	<i>n</i>	<i>nnz</i>	<i>nnz / n</i>	Approximate Condition Number	Structure	Problem
cant	62451	4007383	64.1684	5.082e+05	SYMM	-
consph	83334	6010480	72.1251	2.438e+05	SYMM	-
CurlCurl_3	1219574	13544618	11.1060	2.105e+05	SYMM	Model Reduction
Geo_1438	1437960	63156690	43.9210	4.677e+05	SYMM	-
memchip	2707524	13343948	4.9285	1.305e+07	NON_SYMM	Circuit Simulation
PFlow_742	742793	37138461	49.9984	5.553e+06	SYMM	-
pkustk10	80676	4308984	53.4110	5.589e+02	SYMM	Structural
torso3	259156	4429042	7.0903	2.456e+03	NON_SYMM	-
x104	108384	8713602	80.3956	3.124e+05	SYMM	Structural

Table 4.2.: SuiteSparse matrix set, where SYMM - symmetric; NON-SYMM - non-symmetric; SYMM-PTRN- non-symmetric but with symmetric sparsity pattern

Two different hardware were available for this study. The first machine was a compute-cluster installed in GRS (HW1) which was the main target. Additionally, LRZ CoolMUC-2 Linux cluster (HW2) was used every time when some ambiguous results were obtained in order to check whether a problem was hardware, software or algorithmic specific. Table 4.3 shows compute-node specifications of both compute-clusters.

For this study, OpenMPI implementation of the MPI standard was used because of its open-source license and comprehensive documentation. The library has many

---

#### 4. Methodology and Experimental Setup

---

	HW1 (GRS)	HW2 (LRZ Linux)
Architecture	x86_64	x86_64
CPU(s)	20	28
On-line CPU(s) list	0-19	0-27
Thread(s) per core	1	1
Core(s) per socket	10	14
Socket(s)	2	2
NUMA node(s)	2	4
Model	62	63
Model name	E5-2680 v2	E5-2697 v3
Stepping	4	2
CPU MHz	1200.0	2036.707
Virtualization	VT-x	VT-x
L1d cache	32K	32K
L1i cache	32K	32K
L2 cache	256K	256K
L3 cache	25600K	17920K
NUMA node0 CPU(s)	0-9	0-6
NUMA node1 CPU(s)	10-19	7-13
NUMA node2 CPU(s)	-	14-20
NUMA node3 CPU(s)	-	21-27
RAM per node, GB	128	64

Table 4.3.: Hardware specifications

options for processes pinning which was intensively used during the study.

To make process pinning explicit and deterministic, a Python script was developed to automatically generate rank-files based on the number of MPI processes, OpenMP threads per MPI process, the maximum number of processing elements and the number of NUMA domains. The script always leaves appropriate gaps between MPI processes to allow each process to fork the corresponding number of threads within a parallel region.

A rank-file specifies explicit mapping between MPI processes, ranks, and actual processing elements, cores, of a compute-cluster. The script has two modes, namely: *spread* and *close*. Given a certain number of ranks, *spread* mode tries to distribute them as spread as possible across multiple available NUMA domains in a round-robin fashion. In contrast to *spread* strategy, *close* one groups ranks as close as possible to keep the maximum number of ranks within a single NUMA domain. Figure 4.1 shows an example of mapping 5 MPI ranks and 2 OpenMP threads per rank onto a compute

#### 4. Methodology and Experimental Setup

---

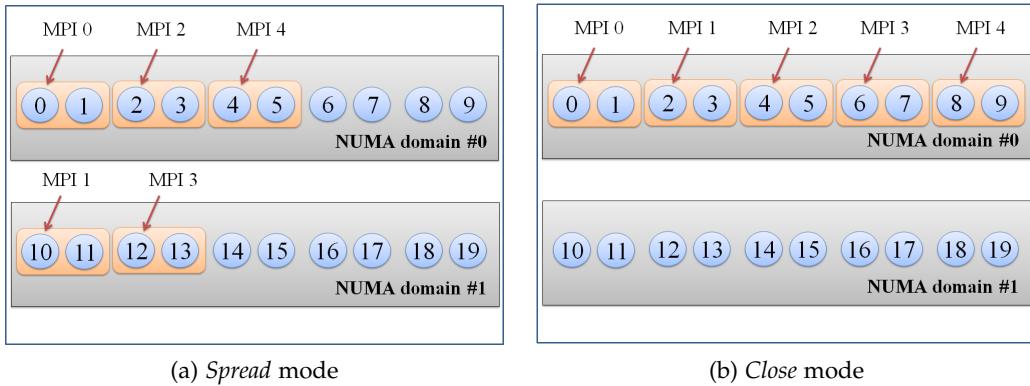


Figure 4.1.: An example of pinning 5 MPI processes with 2 OpenMP threads per process in case of HW1 hardware

node equipped with 20 cores and 2 NUMA domains (HW1).

In this study, PETSc 3.10 and OpenMPI 3.1.1 libraries were chosen and compiled with Intel 18.2 compiler.

# 5. Configuration of a sparse linear solver

## 5.1. Overview of Sparse Linear Solver Types

Next two subsections on iterative and direct sparse methods, respectively, are organized as follows. The first part of each subsection contains a concise theory overview of a linear solver type. The second parts provide discussions of the main sources and means of parallelization of both iterative and direct sparse methods. The last part of each subsection focuses on issues related to numerical solution accuracy of a particular method type.

Then, the section continuous with a discussion and conclusion of which type of linear solvers is well suited for solving systems derived from ATHLET thermo-hydraulic simulations.

### 5.1.1. Iterative Methods

#### 5.1.1.1. Theory Overview

Given an initial guess, an iterative method generates a sequence of approximate solutions by means of a specific rule. Depending on the method and the given problem, there may exist certain conditions such that the aforementioned sequence eventually converges to the exact solution. Iterative methods are preferred for their relatively low computational cost per iteration and storage requirements  $O(nnz)$ . In essence, the methods make use of simple linear algebra kernels at each iteration and thus can handle matrix sparsity efficiently.

The family of iterative methods consists of two distinct classes, namely: stationary and Krylov methods. Nowadays, Krylov methods dominate in the field of scientific computing because of their rather fast convergence in case of solving well conditioned systems or/and a "good" initial guess.

The most well-known methods among the Krylov family are Conjugate Gradient (CG) for symmetric positive definite matrices, MINimal RESidual method (MINRES)

for symmetric indefinite systems and General Minimum RESidual method (GMRES) for non-symmetric systems of linear equations. There also exist different variants of CG such as BiConjugate Gradient Method (BiCG), BiConjugate Gradient STABilized Method (BiCGSTAB), etc.

The key idea is a construction of an approximate solution of a system of linear equations as a linear combination of vectors  $b, Ab, A^2b, A^3b, \dots, A^{n-1}b$  where, without loss of generality, the initial guess  $x_0$  is equal to zero. The combination defines a subspace, also known as Krylov subspace  $\mathcal{K}_n$ . At each iteration, the subspace is expanded by adding and evaluating the next vector in the sequence. The methods usually define and expand another subspace  $\mathcal{L}_n$  of the same size as  $\mathcal{K}_n$  such that  $r_n = b - Ax_n \perp \mathcal{L}_n$  which is known as the Petrov-Galerkin condition. A construction of subspace  $\mathcal{L}_n$  is defined by the methods and based on matrix properties.

Some Krylov methods also have interpretations as minimization problems. For example, GMRES aims to minimize the Euclidean norm of the residual  $r_m$  of a solution vector  $x_m$  in the  $m$ th Krylov subspace  $\mathcal{K}_m$ . However, the basis vectors ( $b, Ab, A^2b, A^3b, \dots, A^{m-1}b$ ) of the  $m$ th Krylov subspace are usually close to linearly dependent and, therefore, a solution vector  $x_m$  is constructed in an orthonormal basis  $U_m$  which forms the same subspace  $\mathcal{K}_m$ .

$$r_m = \min_{x_m \in \mathcal{K}_m} \|Ax_m - b\|^2 = \min_{y_m \in \mathbb{R}^m} \|AU_m y_m - b\|^2 \quad (5.1)$$

where a vector  $x_m$  can be written in the basis  $U_m$  as:

$$x_m = U_m y \quad (5.2)$$

The orthogonalization of the basis can be performed in different ways. Saad and Schultz, in [42], proposed to use the Arnoldi algorithm, see [5] for details, for constructing an  $l_2$ -orthogonal basis. As a result, Equation 5.1 can be written as follows:

$$r_m = \min_{y_m \in \mathbb{R}^m} \|U_{m+1} H_{m+1,m} y_m - \|b\| u_1\|^2 = \min_{y_m \in \mathbb{R}^m} \|H_{m+1,m} y_m - \|b\| e_1\|^2 \quad (5.3)$$

where  $H_{m+1,m}$  is a  $(m+1) \times m$  Hessenberg matrix with non-zero entries defined by the Arnoldi algorithm.

An application of the Givens rotation algorithm results in computing the corresponding QR decomposition of matrix  $H_{m+1,m}$ . After substitution of the matrix with the corresponding QR decomposition and some mathematical transformations, Equation 5.3 can be written as follows:

$$r_m = \min_{y_m \in \mathbb{R}^m} \|Q^T R y_m - \|b\| e_1\|^2 = \min_{y_m \in \mathbb{R}^m} \left\| \begin{pmatrix} R_m \\ 0 \end{pmatrix} y_m - \begin{pmatrix} \tilde{b}_m \\ \tilde{b}_{n-m} \end{pmatrix} \right\|^2 \quad (5.4)$$

Now, the minimization problem 5.4 can be solved as:

$$R_m y_m = \tilde{b}_m \quad (5.5)$$

Given the decomposition of a vector in the orthonormal basis  $U_m$ , Equation 5.2, a solution vector  $x_m$  can be written as:

$$x_m = U_m y_m \quad (5.6)$$

A solution significantly improves with growth of subspace  $\mathcal{K}_m$ . This, as a direct consequence, leads to a considerable increase of a computational cost and storage space. Therefore, the algorithm usually runs till 20 - 50 column vector evaluations of the corresponding Krylov subspace and restarts using a computed approximate solution as an initial guess for the next iteration.

### 5.1.1.2. Parallelization Aspects

In the general case, iterative methods usually make use of dot and matrix-vector products for solving systems of linear equations. Applications of these linear algebra kernels allow to efficiently handle sparsity of linear systems and thus reduce computational complexity of the methods. Additionally, it allows to distribute vectors and matrices across multiple compute-units and solve systems of equations in parallel, efficiently exploiting data-based parallelism. Hence, the main drop of parallel performance mainly comes from process-communication overheads.

However, some methods can have sequential parts that may affect on parallel performance as well. For instance, a triangular solve operation, Equation 5.5, of the GMRES method is usually computed in a single processor because of its small size which depends on the number of iterations before the restart. Hence, if the underlying system

of equations is relatively small then such sequential operations can become a bottleneck in solving the corresponding system.

### 5.1.1.3. Preconditioners

The most important criterion of Krylov methods is convergence. In case of a "bad" initial guess  $x_0$ , the convergence of iterative methods strongly depends on an involved matrix and, in particular, on its condition number. For instance, Equation 5.7 shows dependence of an error reduction in the solution on the corresponding matrix condition number in case of the CG method. It can be clearly observed that a big condition number may result in a very slow error reduction and, therefore, in a slow convergence rate.

$$\|e^i\|_A \leq 2\left(\frac{\sqrt{k}-1}{\sqrt{k}+1}\right)^i \|e^0\|_A \quad (5.7)$$

where  $k = \frac{\lambda_{\max}}{\lambda_{\min}}$  - a matrix condition number

In practice, a linear transformation, known as preconditioning, is applied to the original system of equations in order to reduce its condition number. As a result, the solution process of the modified system can be accelerated. The transformation can be applied in different ways. For instance, Equations 5.8 and 5.9 show applications of a preconditioning matrix  $P$  from left and right sides, respectively.

$$PAx = Pb \quad (5.8)$$

$$AP(P^{-1}x) = b \quad (5.9)$$

In the general case, a good preconditioning algorithm should result in *a low computational cost, low storage space and a low condition number of the transformed system*. Additionally, computations of large linear systems require an algorithm to be adapted for parallel executions as well.

There exist numerous techniques to compute a preconditioner  $P$  for given a matrix  $A$  e.g. (point) Jacobi, Block-Jacobi, incomplete LU decomposition (ILU), multilevel ILU (ILU(p)), threshold ILU (ILUT), incomplete Cholesky factorization (IC), sparse approximate inverse (SPAI), multigrid as a preconditioner, etc. Experience has shown that some techniques can work particularly well for matrices derived from a certain

## 5. Configuration of a sparse linear solver

---

Package name	Origin	Method	Tuning parameters	Comments
block Jacobi	PETSc	block Jacobi	-pc_bjacobi_blocks -sub_pc_type	-
additive Schwarz	PETSc	additive Schwarz	-pc_asm_blocks -pc_asm_overlap -pc_asm_type -pc_asm_local_type -sub_pc_type	-
euclid	hypre	ILU(k)	-nlevel -thresh -filter	deprecated form PETSc
pilut	hypre	ILU(t)	-pc_hypre_pilut_tol -pc_hypre_pilut_maxiter -pc_hypre_pilut_factorrowsize	-
parasail	hypre	SPAI	-pc_hypre_parasails_nlevels -pc_hypre_parasails_thresh -pc_hypre_parasails_filter	-
SPAI	Grote, Barnard	SPAI	-pc_spai_epsilon -pc_spai_nbstep -pc_spai_max -pc_spai_max_new -pc_spai_block_size -pc_spai_cache_size	-
BoomerAMG	hypre	algebraic multigrid	-pc_hypre_boomeramg_cycle_type -pc_hypre_boomeramg_max_levels -pc_hypre_boomeramg_max_iter -pc_hypre_boomeramg_tol etc.	39 tuning parameters in total

Table 5.1.: A list of parallel preconditioning algorithms available in PETSc

PDE or a system of PDEs e.g. Poisson, NavierStokes, etc., and discretized in a certain way. However, *sometimes it can take a quite considerable amount of time to tune a particular preconditioning algorithm in order fulfill all above listed requirements.*

As it can be clearly observed from Table 4.1, all matrices contained in GRS matrix set are very ill-conditioned and, as a result, require suitable linear transformations. PETSc provides various preconditioning methods as well as access to some external preconditioning libraries. Table 5.1 contains a list of some widely-used preconditioning algorithms available in PETSc, capable to run in parallel on distributed-memory machines, as well as their short descriptions and tuning parameters.

Detailed descriptions of all tuning parameters listed in Table 5.1 can be found in either PETSc or Hypre users' manuals, [9] and [18], respectively.

It is worth mentioning that both block Jacobi and additive Schwarz algorithms split the original system into smaller blocks where each block is usually computed on a single processor. Therefore, these algorithms require an explicit specification of another preconditioning algorithm or a linear solver, called as sub-preconditioner, for local block computations. As an example, code Listing 5.1 shows an example of usage of the *sequential* PETSc built-in *LU* matrix decomposition subroutine as a sub-preconditioner for the block Jacobi algorithm. As a result, the number of tuning parameters for these two algorithms can grow significantly.

```

1 -pc_bjacobi_blocks 4
2 -sub_pc_type lu
3 -pc_factor_mat_ordering_type rcm
4 -pc_factor_pivot_in_blocks true

```

Listing 5.1: An example of usage of the PETSc built-in sparse direct linear solver as a sub-preconditioner for the Block Jacobi preconditioning algorithm

## 5.1.2. Direct Sparse Methods

### 5.1.2.1. Theory Overview

Direct sparse methods combine the main advantages of direct and iterative methods. In other words, numerical accuracy of the methods is comparable with the standard Gaussian Elimination process while their computational complexity is typically bounded by  $O(n^2)$  [48] due to efficient treatment of non-zero matrix elements. As it is in case of direct dense methods, a solution of a system of equations is computed by means of forward and backward substitutions using *LU* decomposition of the corresponding matrix.

The multifrontal method is probably the most representative example of direct sparse solvers, introduced by Duff and Reid in [16]. The method is, in fact, an improved version of the frontal method [27] and can compute independent fronts in parallel. A front, also called a frontal matrix, can be considered as a small dense matrix resulting from a column elimination of the original system. There also exist left- and right-looking variants of the multifrontal method explained in detail in [41].

In this subsection, the theory of multifrontal method is explained, which helps to understand parallel aspects and strong scaling behavior of direct sparse solvers in case of parallel execution. To keep the overview rather simple, we assume that matrix  $A$  is symmetric positive definite and sparse. Therefore, matrix decomposition can be conveniently written as follows:

$$A = LDL^T \quad \text{with } (D)_{ii} > 0 \quad (5.10)$$

The algorithm starts with symbolic factorization of System 5.10 with the aim of predicting a sparsity pattern of factor  $L$ . Once it is done the corresponding elimination tree can be constructed.

Figure 5.1 shows an illustrative example of a sparse matrix  $A$  and its Cholesky factor  $L$ , taking from [36]. The solid circles represent the original non-zero elements whereas hollow ones define fill-in elements of  $L$ .

$$A = \begin{pmatrix} 1 & a & & & & & \\ 2 & b & \bullet & \bullet & & & \\ 3 & & c & \bullet & & & \\ 4 & & & d & & & \\ 5 & & & & e & \bullet & \\ 6 & & & & & f & \\ 7 & & & & & & g \\ 8 & & & & & & & h \\ 9 & & & & & & & i \end{pmatrix} \quad L = \begin{pmatrix} 1 & a & & & & & \\ 2 & b & & & & & \\ 3 & & c & & & & \\ 4 & & & d & & & \\ 5 & & & & e & & \\ 6 & & & & & f & \\ 7 & & & & & & g \\ 8 & & & & & & & h \\ 9 & & & & & & & i \end{pmatrix}$$

Figure 5.1.: An example of a sparse matrix and its Cholesky factor, [36]

An elimination tree is a crucial part of the method. It can be considered as a structure of  $n$  nodes where node  $p$  is the parent of  $j$  if and only if it satisfies Equation 5.11. It is worth pointing out that Definition 5.11 is not only one possible and one can define a structure of an elimination tree in a different way as well, [36].

$$p = \min(i > j | l_{ij} \neq 0) \quad (5.11)$$

In fact, node  $p$  represents elimination of the corresponding column  $p$  of matrix  $A$  as well as all dependencies of column  $p$  factorization on results of eliminations of its descendants.

Given Definition 5.11 and a sparsity pattern of factor  $L$ , the corresponding elimination tree can be constructed, as it is shown in Figure 5.2.

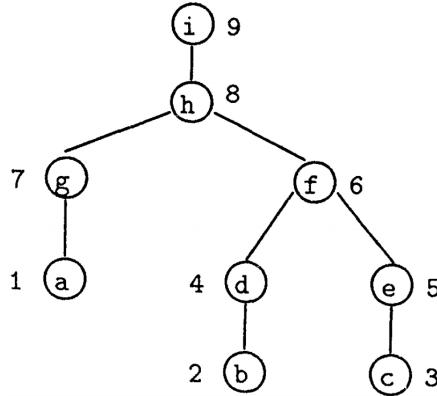


Figure 5.2.: An elimination tree of matrix  $A$  of the example depicted in Figure 5.1, [36]

The fundamental idea of the multifrontal method spins around frontal and update matrices. Frontal matrix  $F_j$  is used to perform Gaussian Elimination for a specific column  $j$  and it is equal to a sum of frame  $Fr_j$  and update  $\hat{U}_j$  matrices, as it can be observed from Equation 5.12

$$F_j = Fr_j + \hat{U}_j = \begin{bmatrix} a_{j,j} & a_{j,i_1} & a_{j,i_2} & \dots & a_{j,i_r} \\ a_{i_1,j} & & & & \\ a_{i_2,j} & & & & \\ \vdots & & & & 0 \\ a_{i_r,j} & & & & \end{bmatrix} + \hat{U}_j \quad (5.12)$$

where  $i_0, i_1, i_2, \dots, i_r$  are row subscripts of non-zeros in  $L_{*j}$  where  $i_0 = j$ ;  $r$  is the number of off-diagonal non-zero elements.

Frame matrix  $Fr_j$  is filled with zeros except the first row and column which contain non-zero elements of the  $j$ th row and column of the original matrix  $A$ . Because of symmetry of matrix  $A$ , the frame matrix is square and symmetric.

In order to describe parts of an elimination tree, notation  $T[j]$  is introduced to represent all descendants of node  $j$  in the tree and node  $j$  itself. In this case, update matrix  $\hat{U}_j$  can be defined as follows:

$$\hat{U}_j = - \sum_{k \in T[j]-j} \begin{bmatrix} l_{j,k} \\ l_{i_1,k} \\ \vdots \\ l_{i_r,k} \end{bmatrix} [l_{j,k} \quad l_{i_1,k} \quad \dots \quad l_{i_r,k}] \quad (5.13)$$

Update matrix  $\hat{U}_j$  is, in fact, can be considered as the second term of the Schur complement i.e. update contributions from already factorized columns of  $A$ .

The subscript  $k$  represents descendant columns of node  $j$ . Hence, only those elements of descendant columns are included and considered which correspond to a non-zero pattern of the  $j$ th column.

Let's consider partial factorization of a 2-by-2 block dense matrix, Equation 5.15, to better understand the essence of update matrix  $\hat{U}_j$ . Let's assume that matrix  $B$  has already been factorized and can be expressed as follows:

$$B = L_B L_B^T \quad (5.14)$$

$$A = \begin{bmatrix} B & V^T \\ V & C \end{bmatrix} = \begin{bmatrix} L_B & 0 \\ VL_B^{-T} & I \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & C - VB^{-1}V^T \end{bmatrix} \begin{bmatrix} L_B^T & L_B^{-1}V^T \\ 0 & I \end{bmatrix} \quad (5.15)$$

The Schur complement, from Equation 5.15, can be viewed as a sum of the original sub-matrix  $C$  and update  $-VB^{-1}V^T$ . The update can be written as a sum of outer products as follows:

$$-VB^{-1}V^T = -(VL_B^{-T})(L_B^{-1}V^T) = - \sum_{k=1}^{j-1} \begin{bmatrix} l_{j,k} \\ \vdots \\ l_{n,k} \end{bmatrix} [l_{j,k} \quad \dots \quad l_{n,k}] \quad (5.16)$$

Firstly, it can be clearly observed that Equations 5.16 and 5.13 are similar. However, Equation 5.13 exploits sparsity of the corresponding rows and columns of factor  $L$  and, therefore, masks unnecessary information. Secondly, frame matrix  $Fr_j$  corresponds to block matrix  $C$  and brings information from the original matrix  $A$ , whereas update matrix  $\hat{U}_j$  adds information about the columns that have already been factorized.

As soon as frontal matrix  $F_j$  is assembled, i.e. the complete update of column  $j$  has been computed, elimination of the first column of matrix  $F_j$  can be started which will result in computing of non-zero entries of factor column  $L_{*j}$ . The process is denoted as

partial factorization of matrix  $F_j$ .

Let's denote  $\hat{F}_j$  as a result of the first column elimination of frontal matrix  $F_j$ . Then, the elimination process of column  $j$  can be expressed as follows:

$$\hat{F}_j = \begin{bmatrix} l_{j,j} & \dots & 0 \\ \vdots & I & \\ l_{i_r,j} & & \end{bmatrix} \begin{bmatrix} 1 & \dots & 0 \\ \vdots & U_j & \\ 0 & & \end{bmatrix} \begin{bmatrix} l_{j,j} & \dots & l_{i_r,j} \\ \vdots & I & \\ 0 & & \end{bmatrix} \quad (5.17)$$

where sub-matrix  $U_j$  represents the full update from all descendants of node  $j$  and node  $j$  itself. Equation 5.18 expresses sub-matrix  $U_j$  as a sum of outer products:

$$U_j = - \sum_{k \in T[j]} \begin{bmatrix} l_{i_1,k} \\ \vdots \\ l_{i_r,k} \end{bmatrix} \begin{bmatrix} l_{i_1,k} & \dots & l_{i_r,k} \end{bmatrix} \quad (5.18)$$

Update column matrix  $U_j$ , also called as a contribution matrix, together with the frontal  $F_j$  and update  $\hat{U}_j$  matrices, forms the key concepts of the multifrontal method. Let's consider an example, depicted in Figure 5.3, to demonstrate the importance of contribution matrices.

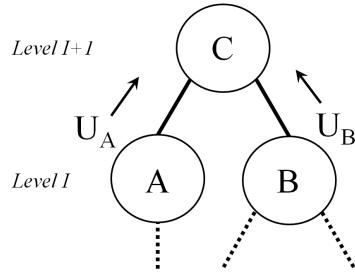


Figure 5.3.: Information flow of the multifrontal method

Let's assume that columns  $A$  and  $B$  have already been factorized and the corresponding contribution matrices  $U_A$  and  $U_B$  have already been computed. According to Equation 5.18, it is known that both  $U_A$  and  $U_B$  matrices contain the full updates of all their descendants including updates of columns  $A$  and  $B$  as well. Therefore, update column matrices  $U_A$  and  $U_B$  have already contained all necessary information to construct update matrix  $\hat{U}_C$ . A detailed proof and careful explanation can be found in [36].

It can happen that only a subset of rows and columns of matrices  $U_A$  and  $U_B$  is needed due to sparsity of column  $C$ . Hence, only relevant elements of the corresponding matrices have to be retrieved to form matrix  $\hat{U}_C$ . For that reason, an additional matrix operation, called *extend-add*, has been introduced in the theory of direct sparse methods.

As an example, taking from [36], let's consider the *extend-add* operation applied to 2-by-2 matrices  $R$  and  $S$  which correspond to indices 5,8 and 5,9 of a matrix  $B$ , respectively.

$$R = \begin{bmatrix} p & q \\ u & v \end{bmatrix}, S = \begin{bmatrix} w & x \\ y & z \end{bmatrix} \quad (5.19)$$

The result of such operation is a 3-by-3 matrix  $K$  which can be written as follows:

$$K = R \text{ } \ddiamond \text{ } S = \begin{bmatrix} p & q & 0 \\ u & v & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} w & 0 & x \\ 0 & 0 & 0 \\ y & 0 & z \end{bmatrix} = \begin{bmatrix} p+w & q & x \\ u & v & 0 \\ y & 0 & z \end{bmatrix} \quad (5.20)$$

Hence, the formation of frontal matrix  $F_j$  can be expressed using the *extend-add* operation and all direct children of node  $j$  as follows:

$$F_j = \begin{bmatrix} a_{j,j} & a_{j,i_1} & a_{j,i_2} & \dots & a_{j,i_r} \\ a_{i_1,j} & & & & \\ a_{i_2,j} & & & & \\ \vdots & & 0 & & \\ a_{i_r,j} & & & & \end{bmatrix} \text{ } \ddiamond \text{ } U_{c_1} \text{ } \ddiamond \text{ } \dots \text{ } \ddiamond \text{ } U_{c_s} \quad (5.21)$$

where  $c_1, c_2, \dots, c_n$  are indices of the direct children of node  $j$ .

At this point, it is worth mentioning that the resulting frontal matrix  $F_j$  forms a small dense block which has to be factorized along the first column. Partial factorization of the block can be efficiently performed by means of the corresponding dense linear algebra kernels.

After partial factorization of matrix  $F_j$ , assembly of contribution matrix  $U_j$  must be completed by adding those elements of  $U_{c_1}, U_{c_2}, \dots, U_{c_s}$  to  $U_j$  that have not been used in factorization of  $F_j$  due to sparsity of column  $j$ . Then, the process continues moving up along the tree. Therefore, complete update matrices are growing in size while the

global elimination process is moving towards the root of the tree.

Manipulations with frontal and contribution matrices play a significant role in performance of the multifrontal method. Sometimes contribution matrices, generated in previous steps, must be stored into a temporary buffer and efficiently retrieved from it later during the global factorization process. This can require to change a column elimination order which can be achieved by some matrix reordering techniques. For instance, *post-ordering*, mentioned by Liu in [36], can be considered as an example of such reordering, in case of symmetric matrices, and can eventually make efficient use of *stack* data structure. Post-ordering is based on topological ordering and thus it is equivalent to the original matrix order. Hence, such reordering results in the same fill-in of the factor [36].

A post-ordered tree implies that each node is ordered before its parent and nodes in each subtree are numbered consecutively. Figure 5.4 shows an example of post-ordering applied to the elimination tree of the matrix shown in Figure 5.1. As a result, consecutive *push* and *pop* operations can be efficiently used during matrix factorization and thus can result in significant simplification of a computer program, see Figure 5.5.

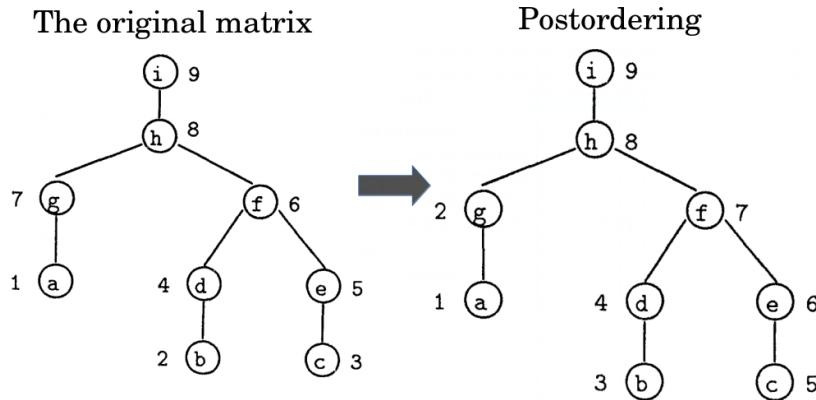


Figure 5.4.: An example of matrix postordering, [36]

In practice, an improved version of the multifrontal method, called the supernodal method, is used. The method tends to shrink an elimination tree by grouping some certain nodes/columns in a single node. As a result, more floating point operations can be performed per memory access by eliminating few columns at once within the same frontal matrix.

A super-node is formed by a set of contiguous columns which have the same off-diagonal sparsity structure. Hence, a super-node has two important properties. Firstly,

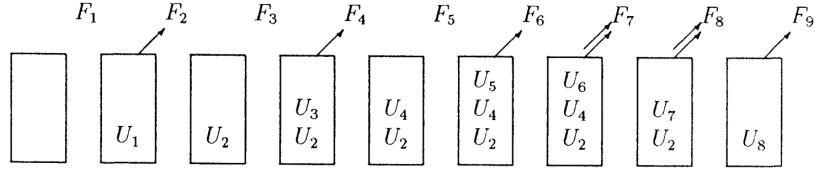


Figure 5.5.: An example of an efficient data treatment during matrix factorization using stacking, [36]

it can be expressed as a set of consecutive column indices, namely:  $\{j, j+1, \dots, j+t\}$  where node  $j+k$  is a parent of  $j+k-1$  in the corresponding elimination tree. Secondly, the size of super-nodal frontal matrix  $\mathcal{F}_j$  is equal to the size of frontal matrix  $F_j$  resulted from the original post-ordered tree. As an example, Figure 5.6 shows a post-ordered matrix  $A$ , its Cholesky factor  $L$  and the resulting super-nodal elimination tree.

$$\bar{A} = \begin{pmatrix} 1 & a & \bullet & & \bullet & \bullet \\ 2 & \bullet & g & & \bullet & \bullet \\ 3 & b & \bullet & \bullet & & \\ 4 & \bullet & d & & \bullet & \bullet \\ 5 & & c & \bullet & \bullet & \\ 6 & & \bullet & e & \bullet & \bullet \\ 7 & & \bullet & & f & \bullet \\ 8 & \bullet & \bullet & \bullet & \bullet & h \\ 9 & \bullet & \bullet & \bullet & \bullet & i \end{pmatrix} \quad \bar{L} = \begin{pmatrix} 1 & a & & & & \\ 2 & \bullet & g & & & \\ 3 & & b & & & \\ 4 & & \bullet & d & & \\ 5 & & & c & & \\ 6 & & & \bullet & e & \\ 7 & & & \bullet & \circ & \bullet & f \\ 8 & \bullet & \bullet & \bullet & \bullet & \circ & h \\ 9 & \bullet & \bullet & \bullet & \bullet & \bullet & \circ & i \end{pmatrix}$$

Figure 5.6.: An example of a supernodal elimination tree, [36]

Equation 5.23 shows an assembly process of super-nodal frontal matrix  $\mathcal{F}_j$ . In contrast to Equation 5.21, frame matrix  $\mathcal{F}_{r_j}$  contains more dense rows and columns. As before, the *extend-add* operation is used to construct the full update block from contribution matrices of the children, namely:  $U_{c_1}, U_{c_2}, \dots, U_{c_s}$ .

$$\mathcal{F}_j = \mathcal{F}_{r_j} \ddagger U_{c_1} \ddagger \dots \ddagger U_{c_s} \quad (5.22)$$

$$\mathcal{F}_j = \begin{bmatrix} a_{j,j} & a_{j,j+1} & \dots & a_{j,j+t} & a_{j,i_1} & \dots & a_{j,i_r} \\ a_{j+1,j} & a_{j+1,j+1} & \dots & a_{j+1,j+t} & a_{j+1,i_1} & \dots & a_{j+1,i_r} \\ \vdots & \vdots & \dots & \vdots & & & \\ a_{j+t,j} & a_{j+t,j+1} & \dots & a_{j+t,j+t} & a_{j+t,i_1} & \dots & a_{j+t,i_r} \\ a_{i_1,j} & a_{i_1,j+1} & \dots & a_{i_1,j+t} & & & \\ \vdots & \vdots & \dots & \vdots & & & \\ a_{i_r,j} & a_{i_r,j+1} & \dots & a_{i_r,j+t} & & & \end{bmatrix} \oplus U_{c_1} \oplus \dots \oplus U_{c_s} \quad (5.23)$$

It is worth mentioning there exist other definitions of super-nodes which allow to amalgamate even more nodes from the original post-ordered tree. For example, Liu pointed out, in [36], that a super-node could be defined without the column contiguity constrain which can result in denser frame matrix  $\mathcal{F}_j$ .

It can be clearly observed that the method consist of three distinct phases, namely: analysis, numerical factorization and solution phases. The analysis phase includes fill reducing matrix reordering, symbolic factorization, post-ordering, amalgamation of nodes, elimination tree construction, etc. During the numerical factorization phase,  $L$  and  $D$ , or  $U$ , factors of the original matrix  $A$  are computed based on sequence of partial factorizations of frontal matrices. Given a matrix decomposition, the solution step computes a solution vector  $x$  by means of backward and forward substitutions.

### 5.1.2.2. Parallelization Aspects

In contrast to iterative methods, parallelization of direct sparse methods mainly comes from task-based parallelism where an elimination tree can be considered as a collection of tasks. In fact, the tree represents data dependencies during column partial factorizations and, therefore, reveals dependent and independent tasks. For example, leaves usually locate in separate branches of an elimination tree and thus represent concurrent tasks that can be executed in parallel. On the other hand, parent nodes represent data dependences from their children and cannot be factorized beforehand. Therefore, sparse direct methods have only limited parallelism which swiftly decreases while computations are moving towards the top of an elimination tree.

Let's consider two simple models, that have been developed for this part of the study, in order to demonstrate potential parallel performance of tree-task parallelism. The models, in fact, are perfectly balanced binary trees with different costs per level. Within

a level, the cost is distributed equally among the nodes. The first model, Figure 5.7a, implies quadratic decrease of a computational cost between nodes of adjacent levels whereas the second one, Figure 5.7b, simulates cubic decay of compute-intensity. The models intend to reflect growth of complete update matrices in size, while moving from bottom to top along an elimination tree, and thus increase of floating point operations.

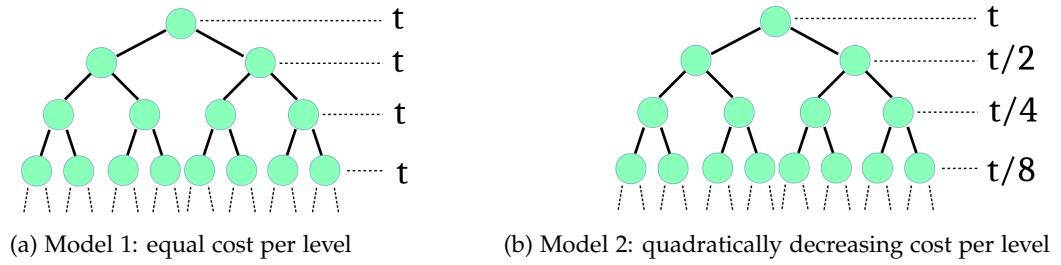


Figure 5.7.: Examples of tree-task parallelism

The models imply parallel computations within a level but sequential execution between them. In other words, to start computations at the next top level, factorizations of all nodes at the current one have to be fully performed. It also means that computations at the next level cannot be started even if there are some available free processors but factorization of the last node at the current level has not been completed yet. Thereby, minimal execution time of both models can be exactly evaluated based on the model descriptions. Essentially, it is equal to a sum of time spent on a single node of each level i.e a sum along the deepest branch. Therefore, it determines asymptotes in the corresponding speed-up graphs.

Figures 5.8a and 5.8b represent strong scaling behavior of both models filled with 65535 nodes i.e. 16 levels. As it can be observed, the models demonstrate a rapid drop of parallel performance, especially in case of the quadratic one, Figure 5.8b . Table 5.2 compares speed-up of two models obtained with 32768 and 20 abstract processors. Number 32768 is equal to the number of leaves at the bottom level and thus implicitly determines the maximum speed-up. It is worth mentioning that two models almost exhaust tree-task parallelism even with 20 processing elements. Further increase of the number of processor can only barely improve the overall parallel performance.

These, rather simple, models reveal the most important fact about tree-task parallelism of sparse direct methods. The performance depends heavily on an elimination tree structure and, in particular, on a distribution of compute-intensity among nodes.

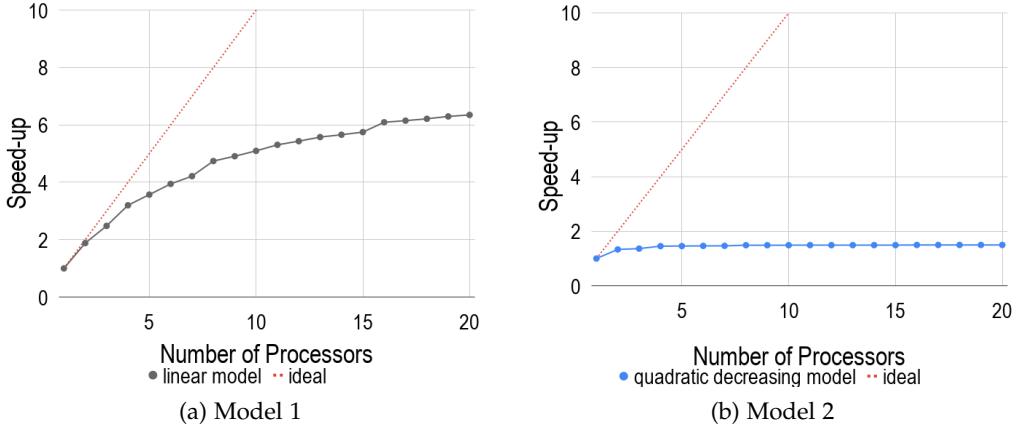


Figure 5.8.: Theoretical speed-up of models 1 and 2

	20 PEs	32768 PEs
Model 1	6.3492	8.0000
Model 2	1.4972	1.5000

Table 5.2.: Theoretical speed-up of models 1 and 2

As it was mentioned above, the intensity is usually centered on the top part of the tree where task-based parallelism is limited due to data dependencies. As an example, Liu observed that factorizations of the last 6 nodes took slightly more than 25% of the total number of floating point operations in case of an application of the multifrontal method to a  $k - by - k$  regular model problem using a nine-point difference operator, [36].

Node-data parallelism is often combined with tree-task one which usually results in an improvement of parallel performance of direct sparse methods. In the general case, frontal matrix  $\mathcal{F}_j$  can be distributed across multiple processors and partially factorized in parallel. However, performance of node-data parallelism depends on both a matrix size and the number of processors assigned to perform factorization. Over-subscription of processing elements to a node can result in slow-down induced by communication overheads. Therefore, data parallelism is only applied to top and middle parts of elimination trees because of fine granularity of bottom levels.

By and large, combined tree-task and node-data parallelism improves performance and strong scaling of sparse direct solvers, however, it cannot change the performance trend induced by tree-task parallelism. Therefore, one can still expect stagnation of

speed-up even with a relatively small number of processors.

A parallel implementation of sparse direct solvers demands to expand the analysis phase by adding two more pre-processing steps, namely: process mapping and load balancing. Both mapping and balancing are usually performed statically during the analysis of an elimination tree.

### 5.1.2.3. Threshold Pivoting and Solution Refinement

Because of an accumulative effect of inexact computer arithmetics due to a floating point representation of real numbers and, as a result, truncations and rounding errors, small numerical values along the main matrix diagonal during the Gaussian Elimination process can result in significant numerical inaccuracy of the process. Therefore, pivoting is a crucial step of Gaussian Elimination. It implies interchanging rows and columns of a matrix in such a way to place distinct and distant values from zero to the main diagonal.

In case of direct dense methods, pivoting is a straightforward operation and can be expressed as multiplication of the original matrix  $A$  by a permutation matrix  $P$ , where each row and column contain a single 1, at the corresponding place, and 0s everywhere. However, treatment of pivoting in direct sparse methods is an issue.

On the one hand, absence of numerical information during the analysis phase makes it impossible to perform pivoting at this step. On the other hand, an application of pivoting during the numerical factorization phase usually distorts all matrix reorderings and, therefore, the elimination tree structure. As a consequence, pivoting can lead to significant fill-in, load unbalance and, as a result, slow-down of numerical factorization. For that reason, threshold pivoting is commonly used, in practice, for direct sparse methods.

Threshold pivoting means that a pivot  $|a_{i,i}|$  is accepted if it satisfies Equation 5.24.

$$|a_{i,i}| \geq \alpha \times \max_{k=i \dots n} |a_{k,i}| \quad (5.24)$$

where  $\alpha \in [0, 1]$  and  $k = i \dots n$  represents row indices of column  $i$  within the fully summed block of a frontal matrix.

Factorization of a column is suspended i.e. delayed, if Equation 5.24 cannot be satisfied within the fully-summed block of a frontal matrix. In this case, the column and the corresponding row are moved to the parent's frontal matrix as a part of its contribution block where the process repeats again. The process is also known as

delayed pivoting and helps to improve numerical accuracy. Higher values of  $\alpha$  lead to more accurate solutions but often generate extra fill-in and lead to load unbalance which, as a result, affect parallel performance. On the other hand, smaller values usually preserve the original elimination tree and, therefore, preserve the load balance computed during the analysis phase by appropriate mapping processors across nodes of the tree. However, in this case, numerical accuracy usually degrades. In practice, values of  $\alpha$  lay in the range between 0.01 and 0.1 [37].

A case when parameter  $\alpha$  is equal to 0 is known as static pivoting which means that no pivoting is being performed during the numerical factorization phase. This allows to better optimize data layout, load balancing, and communication scheduling, [35], before numerical factorization which is supposed to result in better parallel performance.

By and large, solutions computed by direct sparse methods can be numerically inaccurate, in some degree, and may demand to perform solution refinements. As an example, solution accuracy can be improved using the iterative refinement method. Code Listing 5.3 shows a pseduocode of the method where parameter  $\omega$  represents an estimation of the backward error, Equation 5.25, [6]. In practice, the method usually takes 2 or 4 iterations to achieve sufficient numerical accuracy.

```

1 # perform analysis and numerical factorization phases
2 LU = SparseDirectSolver(matrix=A)
3
4 # compute initial solution
5 x = Solve(factorization=LU, rhs=b)
6
7 # compute initial residual
8 r = A * x - b
9
10 while r > ω
11     # find correction
12     d = Solve(factorization=LU, rhs=r)
13
14     # update solution
15     x = x - d
16
17     # update residual
18     r = A * x - b

```

Listing 5.2: Pseudocode of the iterative refinement method

$$\frac{|b - A\hat{x}|_i}{(|b| + |A||\hat{x}|)_i} \quad (5.25)$$

where  $\hat{x}$  is a computed solution;  $|\cdot|$  is the element-wise module operation.

As an alternative to the iterative refinement method, one can use the resulting *LU* decomposition of matrix  $A$  as a preconditioner for an iterative solver, for instance GMRES. Based on experience of ATHLET-NUT users, this approach usually takes 1 up to 3 iterations to achieve desired numerical accuracy even with extreme small values of  $\alpha$ .

Finally, it is worth mentioning that both refinement techniques, mentioned above, exploit only data-based parallelism and, therefore, are scaled well on distributed-memory machines.

### 5.1.3. Results and Conclusion

Nowadays, iterative methods is a common choice for solving sparse systems of linear equations because of their possible fast convergence and high parallel efficiency. However, applications of such methods always demand preconditioning for ill-conditioned systems to make methods converge to numerical accurate solutions. It can be clearly observed from Table 4.1 that numerical integration of thermo-hydraulic simulations in ATHLET entails solving ill-conditioned systems based on estimated condition numbers of matrices form GRS matrix set.

As the first step of the study, we tested various preconditioning algorithms together with their tuning parameters, mentioned in Table 5.1, applied to GRS matrix set. GMRES was chosen as an iterative solver with values of relative and absolute convergence tolerances in the residual norm to be equal to  $1E - 8$  and  $1E - 4$ , respectively. A coarse grid search was used with maximum 3 values for each tuning parameter starting from the default towards more accurate values in order to refine parameter settings of each preconditioning algorithm. Testing results showed that none of them could lead to convergence for the entire set of matrices.

One can assume that a finer grid search can result in finding a suitable preconditioning algorithm with parameter settings that can lead to convergence of GMRES solver for the entire set. However, it is important to point out that the matrices were generated by running the most common GRS thermo-hydraulic test-scenarios and saving them somewhere during the time integration process. Hence, there is no guarantee that the parameters found in such a way can always lead to convergence of GMRES solver in all time steps of any thermo-hydraulic simulation. Therefore, iterative methods may

not satisfy *robustness* criterion stated in Chapter 3 as a non-functional requirement to a sparse linear solver.

Taking into account the above reasoning, we have come to the conclusion that sparse direct methods is the best choice for our problem, in spite of the limited tree-task parallelism described in Subsection 5.1.2.2, because the methods stably result in numerical accurate solutions even in case of ill-conditioned linear systems. Hence, the next objective of the study is to find a suitable sparse direct method and its implementation, and adapt it for HW1 compute-cluster environment in terms of efficient parallel execution.

## 5.2. Selection of a Sparse Direct Linear Solver

*Fair to say, there is no single algorithm or software that is the best for all types of linear systems*

---

— Xiaoye Li, [49]

Nowadays, there exist many different and available sparse direct solvers. Some of them are tuned for specific linear systems whereas others are targeted for the most general cases [49]. Some of them handle tree-task and node-data parallelism in different ways even within the same library depending on sizes of frontal matrices and other criteria [25], [37], [35]. Hence, parallel performance of a direct sparse method depends heavily on its specific implementation. Table 5.3 represents a short summary of almost all available libraries capable to run on distributed-memory machines, at the time of writing, based on works [49] and [8].

It can be clearly observed, from Table 5.3, that only MUMPS, PaStiX and SuperLU\_DIST cover requirements induced by GRS, see Chapter 3, in particular: open-source license and a direct interface to PETSc. Additionally, each development group of these solvers provides technical support for its software package. Therefore, it indirectly means that the solvers are maintainable. It is interesting to notice that all libraries, mentioned above, are implementations of different sparse direct methods, namely: multifrontal (MUMPS), left-looking (PaStiX) and right-looking (SuperLU\_DIST). Moreover, PaStiX and SuperLU\_DIST use only static pivoting [39], [35] whereas MUMPS provides a full implementation of the threshold pivoting strategy [37], described in Subsection 5.1.2.3.

Package	Method	Matrix Types	PETSc Interface	License
Clique	Multifrontal	Symmetric	Not Officially	Open
MF2	Multifrontal	Symmetric pattern	No	-
DSCPACK	Multifrontal	SPD	No	Open
MUMPS	Multifrontal	General	Yes	Open
PaStiX	Left looking	General	Yes	Open
PSPASES	Multifrontal	SPD	No	Open
SPOOLES	Left-looking	Symmetric pattern	No	Open
SuperLU_DIST	Right-looking	General	Yes	Open
symPACK	Left-Right looking	SPD	No	Open
S+	Right-lookin	General	No	-
PARDISO	Multifrontal	General	No	Commercial
WSMP	Multifrontal	General	No	Commercial

Table 5.3.: A list of direct sparse linear solvers adapted for distributed-memory computations, [49], [8], where SPD - Symmetric Positive Definite

To compare the libraries, a couple of flat-MPI tests were performed using GRS matrix set and HW1 compute-cluster. From now onwards, we refer to a flat-MPI test as parallel factorizations of a matrix performed with varying values of the MPI process count from 1 to 20 and conducted on a single compute-cluster node.

PETSc library was compiled and configured with MUMPS (version 5.1.2), PaStiX (version 6.0.0) and SuperLU\_DIST (version 5.4) packages using their default parameter settings. An internal built-in PETSc profiler was used to measure execution time. A time limit of 15 minutes was set up for each test-case to prevent blocking of a cluster compute-node from an unexpected long program execution. Results are summarized in Tables 5.4, 5.5, 5.6 and in appendix B where numerical values are given in seconds.

Some problems were detected during SuperLU\_DIST library testing. First of all, executions of *cube-64* and *k3-2* test-cases exceeded the set time limit. Secondly, it was noticed the library was crashing during processing of *k3-18*, *cube-645* and (partially) *pwr-3d* test-cases. Debugging revealed that a segmentation fault occurred in function *pdgstrf* during the numerical factorization phase. Nonetheless, it is still unclear whether the problem was software or hardware specific. A solution or a reason of such program behavior has not been found at the moment of writing.

MPI	MUMPS	PaStiX	SuperLU
1	7.02E-02	8.72E-02	3.17E+00
2	6.73E-02	7.10E-02	1.43E+00
3	6.36E-02	7.01E-02	1.07E+00
4	6.28E-02	7.11E-02	8.17E-01
5	6.50E-02	7.15E-02	7.51E-01
6	6.72E-02	7.62E-02	6.15E-01
7	6.91E-02	7.69E-02	6.48E-01
8	6.89E-02	8.17E-02	5.41E-01
9	7.50E-02	8.28E-02	5.02E-01
10	7.22E-02	8.52E-02	4.64E-01
11	7.55E-02	8.89E-02	5.82E-01
12	7.61E-02	1.06E-01	4.37E-01
13	7.84E-02	9.72E-02	5.43E-01
14	8.06E-02	1.02E-01	4.22E-01
15	8.20E-02	1.19E-01	3.91E-01
16	8.07E-02	1.19E-01	4.44E-01
17	8.38E-02	1.22E-01	5.19E-01
18	8.40E-02	1.26E-01	3.77E-01
19	8.58E-02	1.33E-01	5.47E-01
20	8.64E-02	1.49E-01	3.39E-01

Table 5.4.: Comparisons of parallel performance of *cube-5* matrix factorizations using MUMPS, PaStiX and SuperLU\_DIST solvers with their default parameter settings

MPI	MUMPS	PaStiX	SuperLU
1	1.36E+00	1.39E+00	time-out
2	1.00E+00	9.82E-01	time-out
3	8.83E-01	1.06E+00	time-out
4	8.17E-01	8.74E-01	time-out
5	7.85E-01	8.50E-01	time-out
6	8.06E-01	8.52E-01	time-out
7	7.71E-01	8.33E-01	time-out
8	7.66E-01	8.33E-01	time-out
9	7.93E-01	8.35E-01	time-out
10	8.07E-01	8.15E-01	time-out
11	7.75E-01	8.15E-01	time-out
12	7.81E-01	8.10E-01	time-out
13	7.85E-01	8.35E-01	time-out
14	7.85E-01	8.18E-01	time-out
15	7.88E-01	8.46E-01	time-out
16	7.81E-01	8.23E-01	time-out
17	6.83E-01	8.49E-01	time-out
18	7.96E-01	8.44E-01	time-out
19	8.04E-01	8.65E-01	time-out
20	6.85E-01	8.87E-01	time-out

Table 5.5.: Comparisons of parallel performance of *cube-64* matrix factorizations using MUMPS, PaStiX and SuperLU\_DIST solvers with their default parameter settings

To complete comparison and evaluate parallel performance SuperLU\_DIST library, an additional test was conducted using a 2D formulation of the Poisson problem with 100000 unknown. According to the results, SuperLU\_DIST managed to complete matrix factorizations within the set time limit without crashing, however, it showed abnormal jagged strong scaling behavior. Moreover, it turned out it was the slowest in comparison to the other solvers. The results are shown in Figure 5.9.

5. Configuration of a sparse linear solver

---

MPI	MUMPS	PaStiX	SuperLU
1	1.55E+02	6.44E+01	crashed
2	6.28E+01	4.84E+01	crashed
3	5.06E+01	5.02E+01	crashed
4	4.17E+01	4.50E+01	crashed
5	2.52E+01	3.98E+01	crashed
6	2.58E+01	4.29E+01	crashed
7	2.65E+01	4.30E+01	crashed
8	2.59E+01	3.73E+01	crashed
9	1.95E+01	4.08E+01	crashed
10	1.91E+01	3.81E+01	crashed
11	1.77E+01	3.81E+01	crashed
12	1.60E+01	3.75E+01	crashed
13	1.42E+01	3.58E+01	crashed
14	1.45E+01	3.59E+01	crashed
15	1.47E+01	3.57E+01	crashed
16	1.41E+01	3.52E+01	crashed
17	1.54E+01	3.45E+01	crashed
18	1.52E+01	3.31E+01	crashed
19	1.52E+01	3.31E+01	crashed
20	1.38E+01	3.16E+01	crashed

Table 5.6.: Comparisons of parallel performance of  $k3\text{-}18$  matrix factorizations using MUMPS, PaStiX and SuperLU\_DIST solvers with their default parameter settings

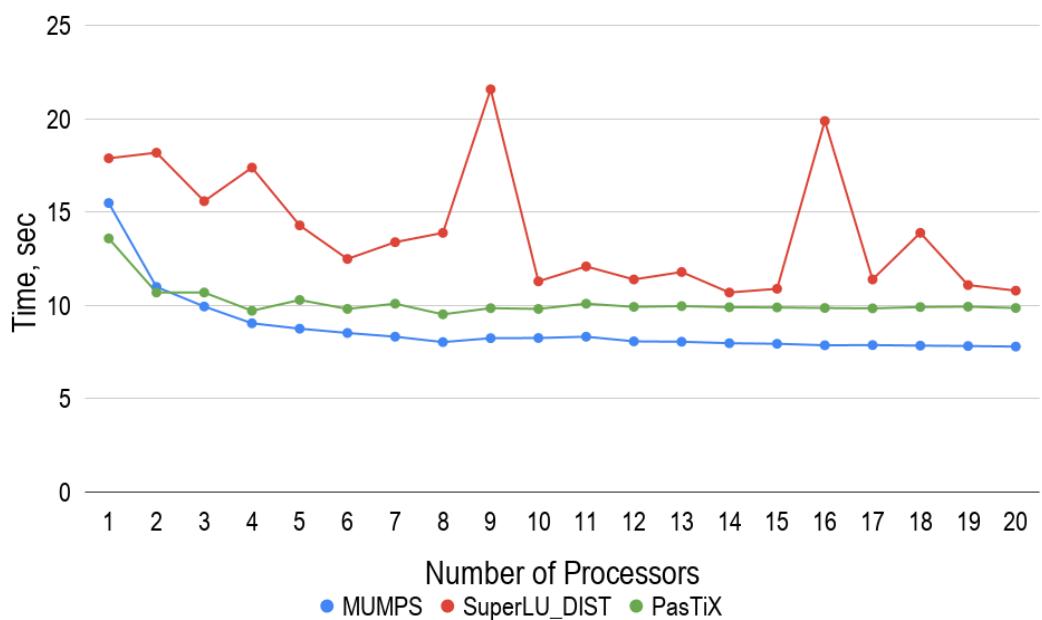


Figure 5.9.: Comparisons of parallel performance of MUMPS, PaStiX and SuperLU\_DIST libraries during 5 point-stencil Poisson matrix (1000000 equations) factorizations

According to the initial and additional tests, MUMPS library showed the best parallel performance and scaling in contrast to the other solvers. No abnormal behavior during its operation was detected. In some cases, it only required to increase a multiplicative factor of estimated working space which was used to hold frontal matrices and factors  $L$  and  $U$  in memory. PaStiX was the second fastest solver according to the results of testing. However, it was often considerably slower than MUMPS. At the same time, SuperLU\_DIST showed the worst results. Additionally, as it was mentioned above, we experienced some technical problems during operation of this library.

A literature review showed quite contradictory results and conclusions. For example, Gupta, Koric, and George, in [25], came to nearly the same inference with respect to MUMPS, as we did, comparing parallel performance of WSMP, MUMPS and SuperLU\_DIST libraries using their matrix set. However, Kwack, Bauer, and Koric showed, in [31], that SuperLU\_DIST spent the least amount of time on solving systems of linear equations in contrast to the other solvers used in their work. It is needless to say that both research groups used different matrix sets and hardware. Nevertheless, it reveals a quite important fact that a selection of a particular method and its implementation can depend heavily on a specific matrix set.

In this section, we have compared different sparse direct methods and their concrete implementations using their default parameter settings with regard to GRS matrix set. Based on the obtained results and literature review, MUMPS library is chosen for the rest of the study. In Section 5.3, we make an overview of the library and its specific traits.

### 5.3. Overview of MUMPS Library

Originally, MULTifrontal Massively Parallel sparse direct Solver (MUMPS) was a part of the PARASOL Project. The project was an ESPRIT IV long term research with the main goal to build and test a portable library for solving large sparse systems of equations on distributed memory systems [1]. An important aspect of the research was a strong link between the developers of the sparse solvers and industrial end users, who provided a range of test problems and evaluated the solvers [3]. Since 2000, MUMPS had continued as an ongoing project and the library have contained almost 5 main releases, at the moment of writing.

As it was mentioned in Section 5.2, MUMPS is an implementation of the multifrontal method. Therefore, MUMPS performs all three phases in sequence, namely: analysis,

numerical factorization and solution. The numerical factorization and solution phases were fully described in detail in Subsection 5.1.2.1. In this section, the analysis phase of MUMPS is examined since implementations of this phase often vary between libraries due to different performance considerations.

According to the library documentation, the analysis phases of MUMPS consists of several pre-processing steps:

1. Fill reducing reordering
2. Symbolic factorization
3. Scaling
4. Amalgamation
5. Mapping

1) To handle both symmetric and unsymmetric cases, MUMPS performs fill reducing reordering based on  $A + A^T$  sparsity pattern. The library provides numerous sequential algorithms for the reordering such as Approximate Minimum Degree (AMD) [2], Approximate Minimum Fill (AMF), Approximate Minimum Degree with automatic quasi-dense row detection (QAMD) [4], Bottom-up and Top-down Sparse Reordering (PORD) [43], Nested Dissection coupled with AMD (Scotch) [40], Multilevel Nested Dissection coupled with Multiple Minimum Degree (METIS) [29]. Additionally, MUMPS can work together with ParMETIS and PT-Scotch which are extensions of METIS and Scotch libraries for parallel execution, respectively. MUMPS also provides users with an option to select a fill-in reducing algorithm in run-time based on a matrix type, size and the number of processors [37].

2) Sparsity structures of factors  $L$  and  $U$  are computed during the symbolic factorization pre-processing step, based on permuted matrix  $A$  after fill-in reducing reordering. It gives the input information for the elimination tree building process. All computations are performed using an undirected graph  $G(A)$  associated with a matrix  $A$  at this step.

3) Rows and/or columns of matrix  $A$  can be scaled during either the analysis or factorization phase in order to improve numerical solution accuracy. As an additional consequence, this pre-processing step can result in more reliable estimations of required memory space and load balancing, performed during the analysis phase, due to a reduced amount of pivoting during numerical factorization. Different scaling approaches

are adopted in MUMPS, namely: diagonal, column or column-row scalings during the numerical factorization phase, see [37] for details; scalings based on works [14], [13] and [15] for the analysis phase.

4) During the amalgamation step, described in Subsection 5.1.2.1, sets of columns with the same off-diagonal sparsity pattern are group together to create denser nodes, also known as super-nodes. The process leads to restructuring the original elimination tree to an amalgamated one of super-nodes which is also know as the *assembly tree*. The main purpose of this step is to improve efficiency of dense matrix operations.

5) A host process, chosen by MUMPS, creates a pool of tasks where each task refers to partial factorization of a node i.e. a frontal matrix. Each node belongs to one of three different types according to a size of the frontal matrix that a node refers to. Type 1 and 2 nodes represent small- and medium-sized frontal matrices, respectively. Whereas a type 3 node represents the root node of an assembly tree i.e the largest frontal matrix. MUMPS uses different parallel computational strategies, that are explained below, in order to process nodes of different types with the aim of achieving better parallel performance. The host process distributes tasks among all available processes in such a way to achieve good memory and compute balances. Figure 5.10 shows an example of a process distribution in MUMPS.

Type 1 nodes are grouped in subtrees, according to the Geist-Ng algorithm [20], and each subtree is processed by a single process to avoid the finest task granularity, which can cause high communication overheads.

In case of type 2 nodes, the host process assigns each node to one process, called the *master*, which holds fully summed rows and columns of a node as well as performs threshold pivoting and partial factorization. During the numerical factorization phase, in run-time, a master process first receives symbolic information, describing contribution block structures, from its children. Then, the master collects information concerning the load balances of all other processes and decides, *dynamically*, which of them, *slaves*, are going to participate in the node factorization. After that, the master informs the chosen slaves that a new task has been allocated for them; maps them according to a 1D block column distribution and sends them the corresponding parts of the frontal matrix. Then, the slaves communicate with the children of the master process and collect the corresponding numerical values. The slaves are in charge of assembly and computations of the partly summed rows. The computational process is illustrated in Figure 5.18, Subsection 5.4.3.

The root node belongs to the 3rd type. The host *statically* assigns a master for the root, as it is in case of type 2 nodes, to hold all the indices describing the structure of its frontal matrix. Before factorization, the structure of the root frontal matrix is statically mapped onto a 2D grid of processes using a block cyclic distribution. This allows to determine, during the analysis phase, which process an entry of the root is assigned to. Hence, the original matrix entries and parts of the contribution blocks can be assembled as soon as they are available. Because of threshold pivoting, the master process collects indices for all delayed variables of its sons; builds the final structure of the root frontal matrix and broadcasts the corresponding symbolic information to all slaves. The slaves, in turn, adjust their local data structure and, right after this, perform numerical factorization in parallel.

It is important to mention that if the root node size is less than a certain computer depended parameter value, defined internally by MUMPS, the root node will be treated as a type 2 one, [37].

An example of static/dynamic scheduling i.e. process mapping, is depicted in Figure 5.10.

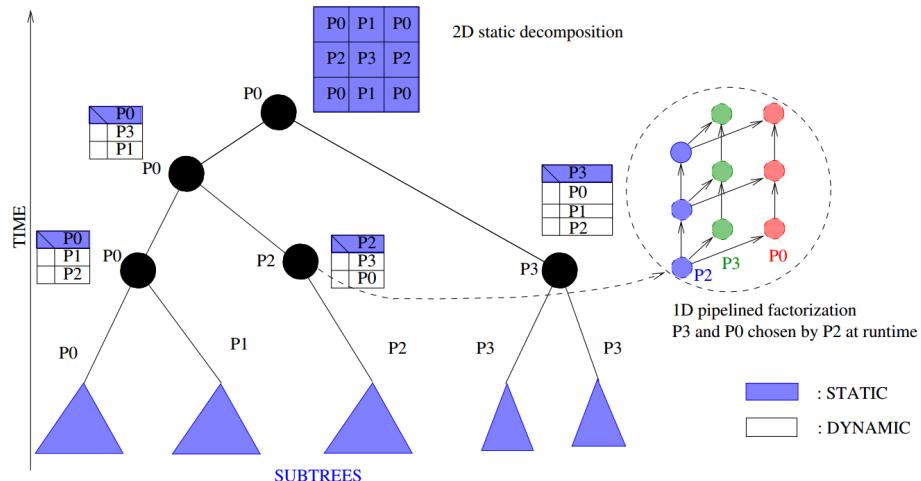


Figure 5.10.: An example of static and dynamic scheduling in MUMPS, [32]

## 5.4. Configuration of MUMPS Library

This section is organized as follows. Each subsection describes configuration of MUMPS with particular techniques, methods or libraries stated in the title of a subsection, e.g. Subsection 5.4.3, i.e. "*Optimized BLAS Implementations*", stands for "*Configuration of MUMPS Library with Optimized BLAS Implementations*".

### 5.4.1. Fill Reducing Reorderings

Fill reducing reordering is one of the first and the most important steps of sparse matrix factorization. As the name suggests, the step aims to reduce fill-in of  $L$  and  $U$  factors. However, it may have a strong and indirect impact on an elimination/assembly tree structure. As we discussed in Subsection 5.1.2.2, the structure defines tree-task parallelism as well as sizes of frontal matrices and, therefore, performance of the method.

MUMPS provides various algorithms for fill reducing reordering, as it was mentioned above. A detailed study and comparison of different methods were done by Guermouche, L'Excellent, and Utard, in [24], for sequential execution of the analysis phase. Guermouche, L'Excellent, and Utard noticed that trees generated by METIS and SCOTCH were rather wide (because of the global partitioning performed at the top), while the trees generated by AMD, AMF and PORD tend to be deeper. In addition, they observed two important things. Firstly, they noticed that both SCOTCH and METIS generated much better balanced trees in contrast to other methods. Secondly, according to their results, SCOTCH and METIS produced trees with bigger frontal matrices in contrast to those trees generated by other reordering techniques, [24].

In this subsection, we are going to investigate an influence of two different parallel fill reducing reordering algorithms provided by PT-Scotch and ParMETIS libraries on parallel performance of MUMPS. The algorithmic difference between the corresponding PT-Scotch and ParMETIS subroutines was mentioned in Section 5.3.

To perform testing, PETSc, MUMPS, PT-Scotch and ParMETIS libraries were downloaded, compiled, configured and linked together, using their default parameter settings. Tests were carried out using only flat-MPI mode on HW1 compute-cluster without any explicit process pinning. The results are shown in Figures 5.11 and 5.12 as well as in appendix C.

According to the results of testing, parallel performance of MUMPS can vary signifi-

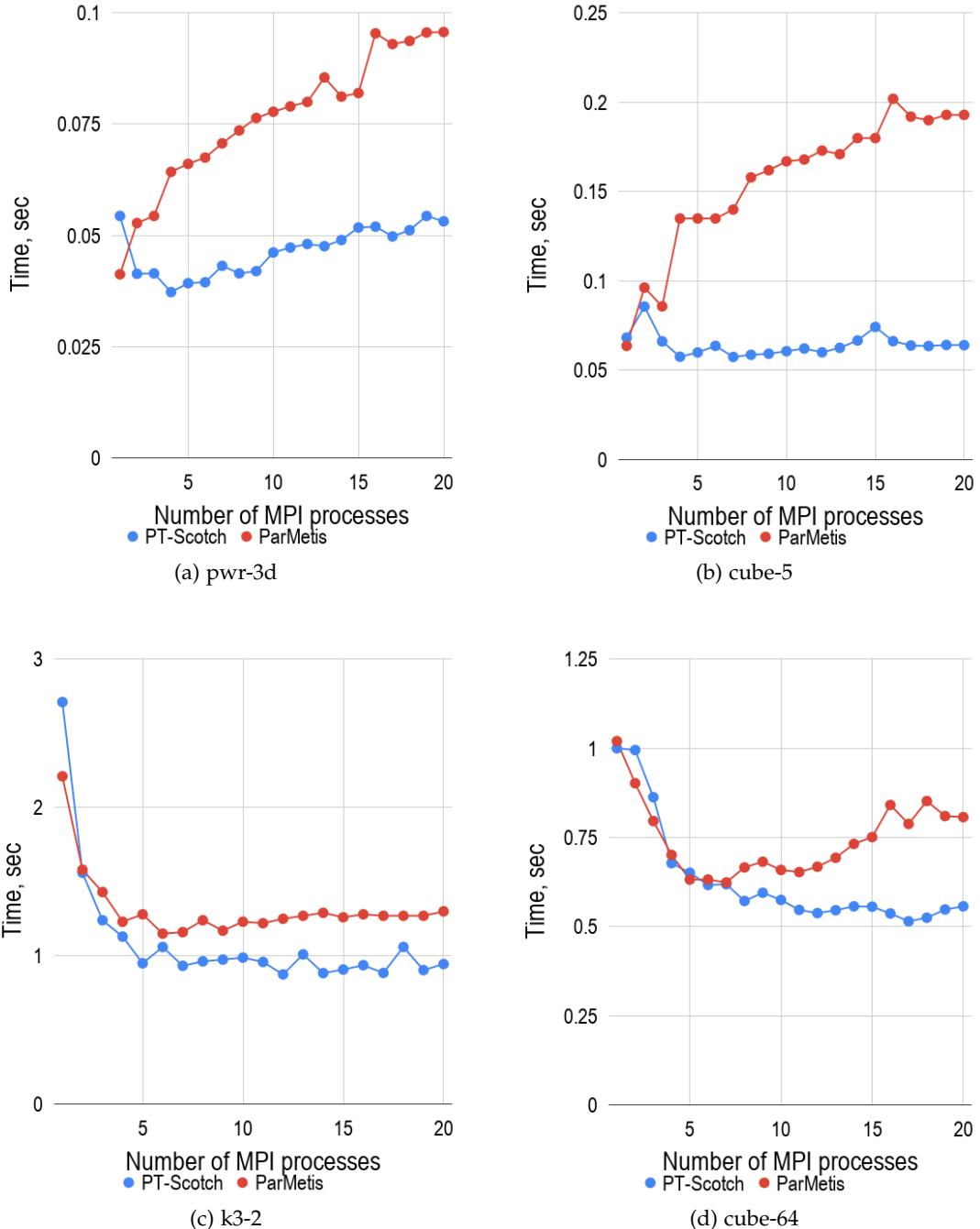


Figure 5.11.: An influence of different fill reducing algorithms on parallel factorizations of *pwr-3d*, *cube-5*, *k3-2* and *cube-64* matrices

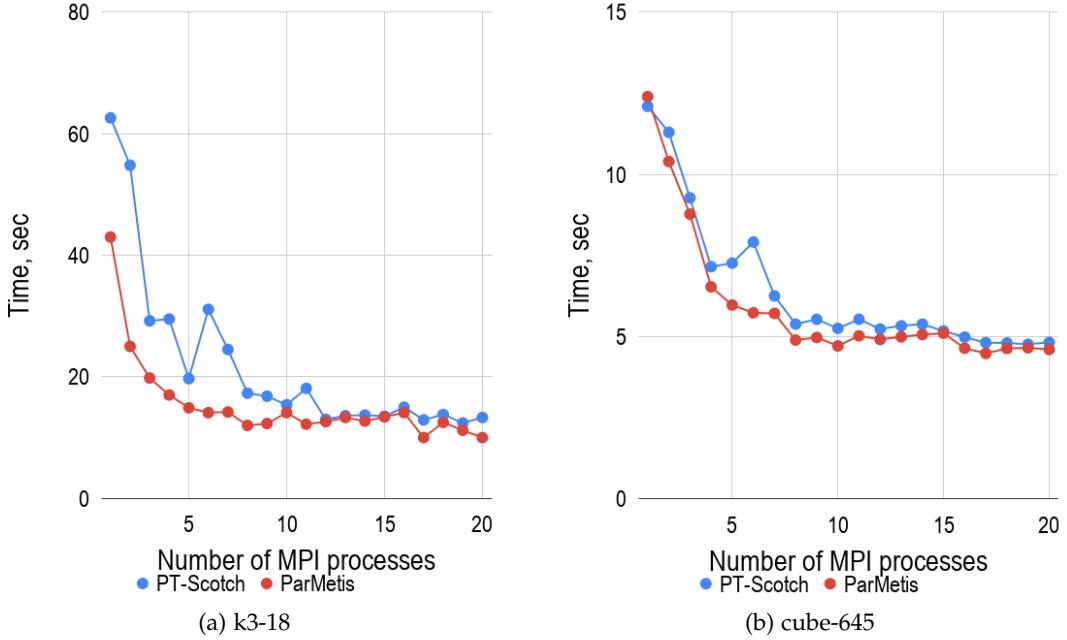


Figure 5.12.: An influence of different fill reducing algorithms on parallel factorizations of *k3-18* and *cube-645* matrices

cantly and very sensitive to applied fill-in reducing reordering algorithms. On average, difference in execution time between the algorithms achieves almost 15%. However, in some particular cases, *cube-5* and *pwr-3d*, the difference varies around 40-55%.

It is important to mention that both algorithms, PT-Scotch and ParMetis, are based on different heuristic approaches. It is relevant to assume that efficiency of a particular heuristic can be very sensitive to a matrix structure and size. This fact makes it difficult to predict in advance which algorithm is better to use for a specific case.

Considering results obtained using GRS matrix set, we can observe that PT-Scotch is the best choice for small- and medium-sized matrices, namely: *pwr-3d*, *cube-5*, *k3-2* and *cube-64* cases. Whereas, PerMetis tends to work better for relatively big systems, such as *k3-18* and *cube-645*. *However, we keep in mind that the number of GRS test-cases may be not enough to make such conclusion and, therefore, the matrix set must be extended considerably for a future study.*

During the testing, we noticed that applications of ParMetis to small systems of equations showed a strong negative effect on parallel performance of MUMPS. The results showed that factorization time of *pwr-3d* and *cube-5* matrices grew with the increase of the number of processing units.

A simple profiling showed two important things. Firstly, numerical factorization time and time spent on the analysis phase had approximately the same order in case of sequential execution i.e. 1 MPI process. Secondly, while numerical factorization time were barely decreasing with increase of the number of processing elements, time spent on the analysis phase significantly grew. Therefore, the slow-down of MUMPS in case of these two test-cases mainly came from overheads of the analysis phase.

A careful investigation revealed that the analysis phase contained several peaks at points where the MPI processor count was equal to a power of two. We assumed the cause could result from either fill reducing reordering or process mapping steps. However, a detailed profiling and tracing of the analysis phase, which are out of the scope of this study, are required in order to give the exact answer. The results of profiling are shown in Figure 5.13.

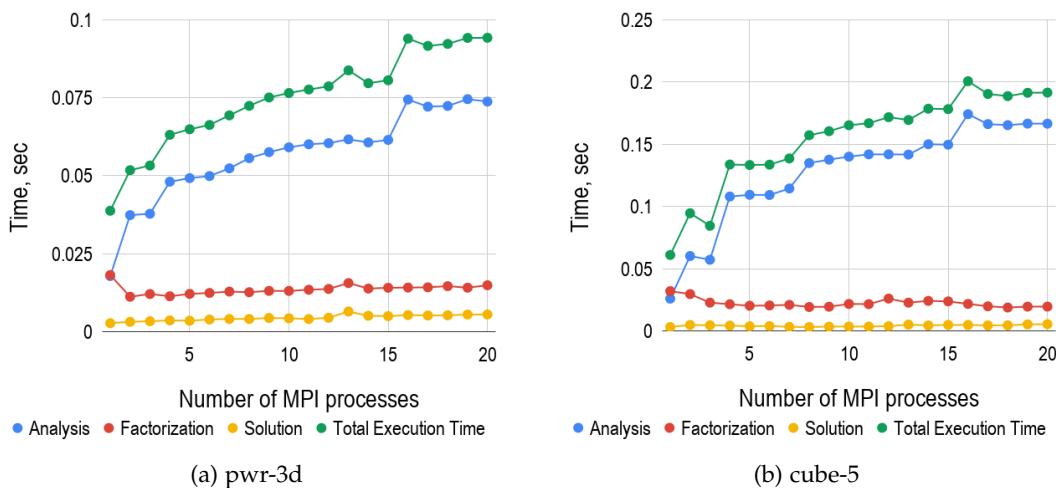


Figure 5.13.: Profiling of MUMPS-ParMetis configuration applied to parallel factorizations of relatively small matrices

In this subsection, we have presented an influence of two different fill-in reducing algorithms on parallel performance of MUMPS. We have observed that a correct choice

of an algorithm can lead to a significant improvement in terms of the overall parallel execution time. We have shown there is no a single algorithm that performs the best for all test-cases. At the moment of writing, we have come to the conclusion that there is no an indirect metric to predict the best algorithm in advance for a specific system of equations. Sometimes PT-Scotch and ParMetis can result in nearly the same performance as it was, for example, in case of *CurlCurl\_3* and *cant* matrices, see appendix C. Therefore, from time to time, it can be quite difficult to decide which package to use even with available flat-MPI test results. At the end, we have assigned each test-case to a specific fill reducing reordering method based on results of the conducted experiments and our subjective opinion. The results are summarized it in Tables 5.7 and 5.8.

Matrix Name	Ordering	<i>n</i>	<i>nnz</i>	<i>nnz / n</i>
cube-5	PT-Scotch	9325	117897	12.6431
cube-64	PT-Scotch	100657	1388993	13.7993
cube-645	ParMetis	1000045	13906057	13.9054
k3-2	PT-Scotch	130101	787997	6.0568
k3-18	ParMetis	1155955	7204723	6.2327
pwr-3d	PT-Scotch	6009	32537	5.4147

Table 5.7.: Assignment of GRS matrices to specific fill-in reducing algorithms

Matrix Name	Ordering	<i>n</i>	<i>nnz</i>	<i>nnz / n</i>
cant	ParMetis	62451	4007383	64.1684
consph	PT-Scotch	83334	6010480	72.1252
memchip	PT-Scotch	2707524	13343948	4.9285
PFlow_742	PT-Scotch	742793	37138461	49.9984
pkustk10	PT-Scotch	80676	4308984	53.4110
torso3	ParMetis	259156	4429042	17.0903
x104	PT-Scotch	108384	8713602	80.3956
CurlCurl_3	PT-Scotch	1219574	13544618	11.1060
Geo_1438	ParMetis	1437960	63156690	43.9210

Table 5.8.: Assignment of SuiteSparse matrices to specific fill-in reducing algorithms

From now onwards, assignments mentioned in Tables 5.7, 5.8 are going to be used without explicitly referring to it.

### 5.4.2. MPI Process Pinning

Due to intensive and complex manipulations with frontal and contribution matrices, one can assume that MUMPS belongs to a group of memory bound applications. In this case, memory access becomes a bottleneck. A common strategy to improve performance of a memory bound computer program running on distributed-memory machines is to distribute MPI processes equally across all available NUMA domains within a compute node. Given the fact that each NUMA domain possesses its own system bus, this strategy allows to reduce conjunction of memory traffic by balancing data requests equally among the memory channels.

However, due to the fact that MUMPS uses both task and data parallelism as well as a complex task scheduling, it becomes difficult to state which process pinning strategy is better to use i.e. *close* or *spread*, described in Chapter 4.

Therefore, a couple of flat-MPI tests were carried out using both GRS and SuiteSparse matrix sets in order to investigate an influence of different pinning strategies on MUMPS parallel performance. For this group of tests, MUMPS was ran with the default parameter settings but with a specific fill-in reducing algorithm assigned to each test-case according to Tables 5.7 and 5.8. The tests were performed on both HW1 and HW2 machines. A comparison between different hardware also allows to investigate an influence of different numbers of independent system buses within a compute-node on parallel performance of MUMPS since HW1 and HW2 machines have 2 and 4 NUMA domains, respectively. Results are shown in Figures 5.14, 5.15, 5.16 and in appendix D. Each graph depicts the total execution time of MUMPS spent on a test-case i.e. time spent on the analysis, factorization and solution phases.

The tests revealed that, in the general case, the *spread*-pinning strategy performed better for both machines. On average, the strategy allows to reduce run-time by approximately 5.5% and 13.8% for HW1 and HW2 machines, respectively. The main performance gain can be observed in the middle range of the MPI process count i.e. the range between 2 and 12 MPI processes, where performance curves of *spread* and *close* strategies noticeably deviate. On the other hand, the difference becomes less and less prominent while the process count is reaching either its maximum or minimal values. In these cases, the difference between process distributions of both strategies becomes less noticeable as well. As an extreme example, the points where the process count is equal to 1 and 20 show the same performance, in case of HW1 machine which possesses only 20 cores in a compute-node, because the points basically represent exactly the same process distributions.

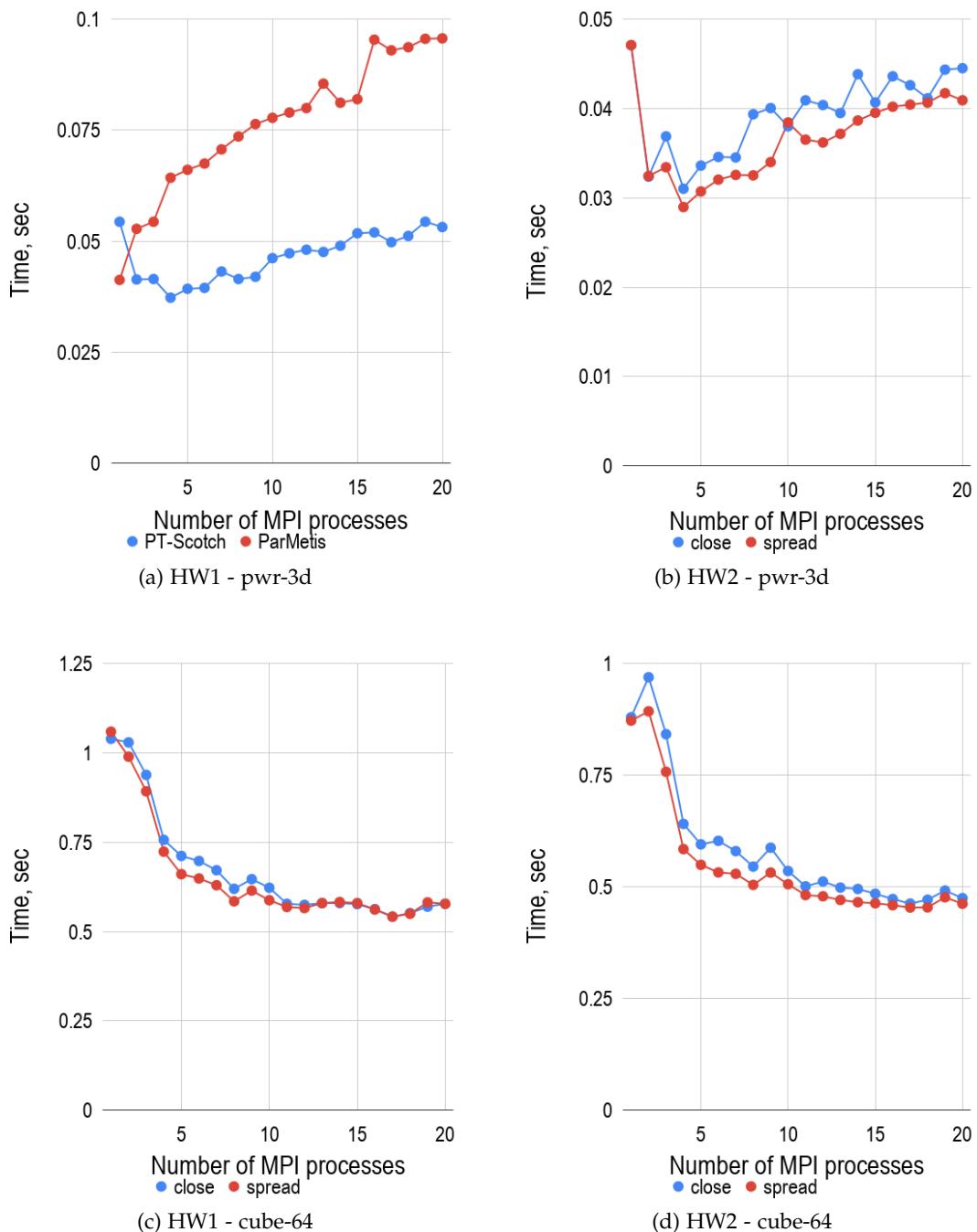


Figure 5.14.: Comparisons of *close* and *spread* pinning strategies applied to parallel factorizations of *pwr-3d* and *cube-64* matrices

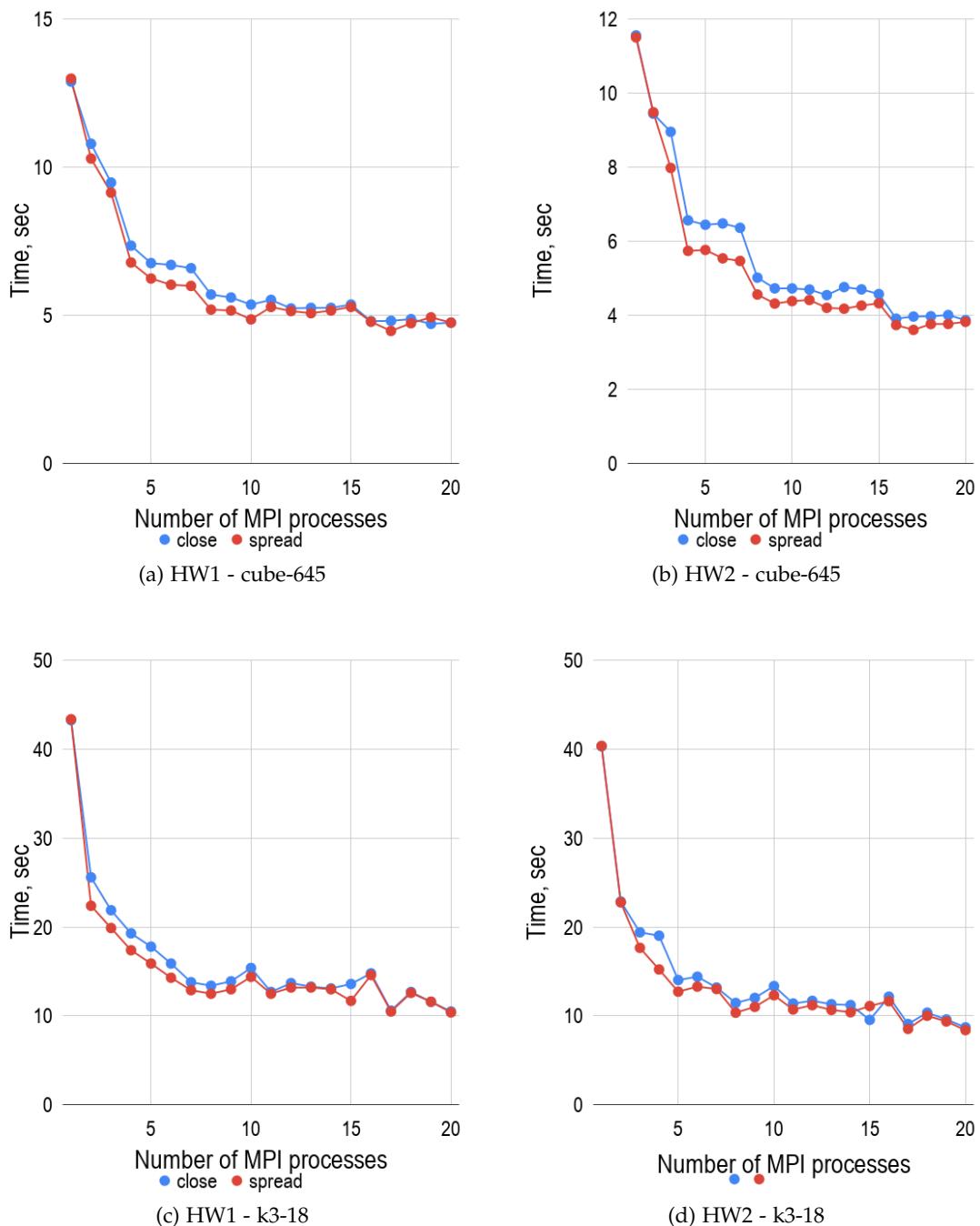


Figure 5.15.: Comparisons of *close* and *spread* pinning strategies applied to parallel factorizations of *cube-645* and *k3-18* matrices

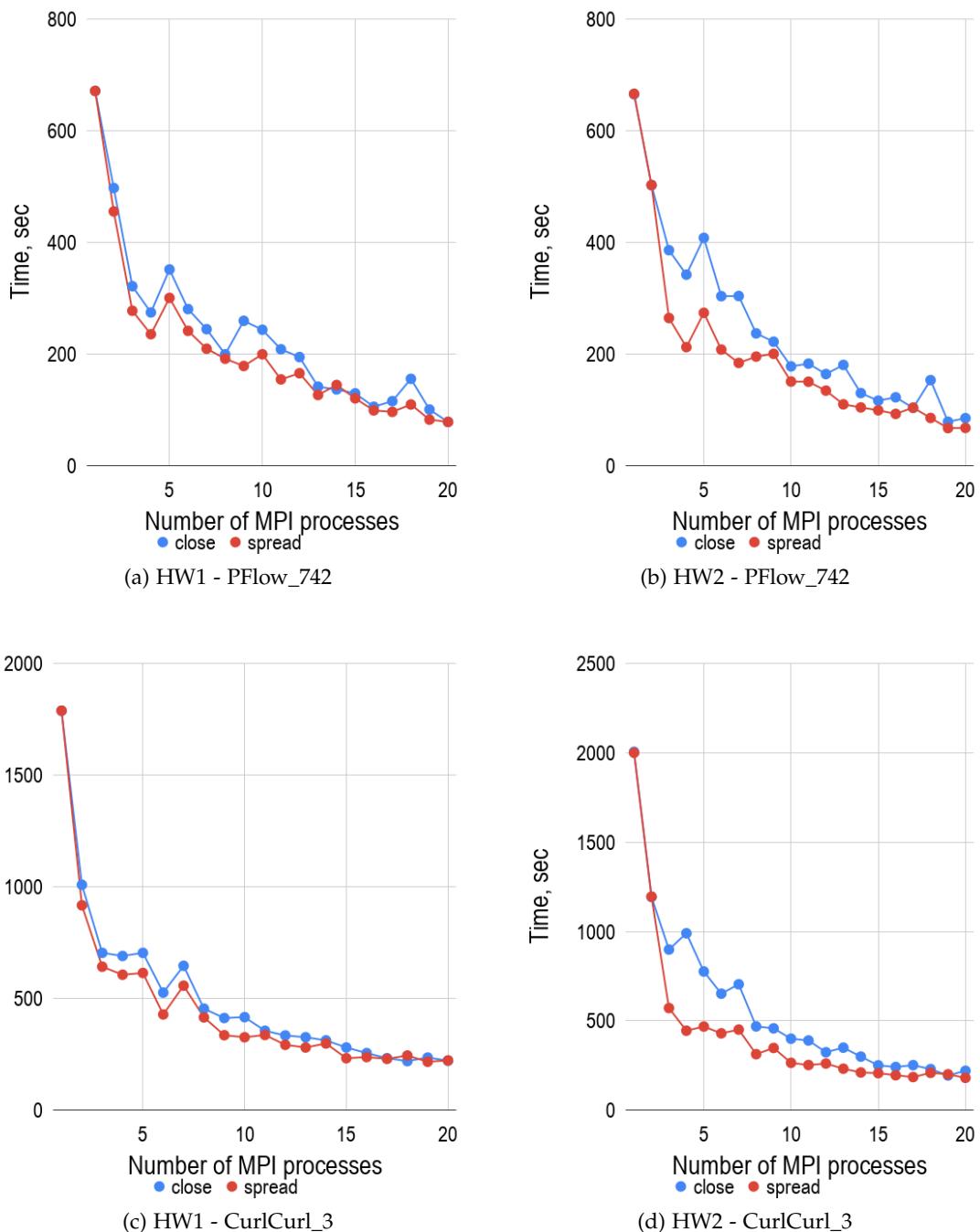


Figure 5.16.: Comparisons of *close* and *spread* pinning strategies applied to parallel factorizations of *PFlow\_742* and *CurlCurl\_3* matrices

It is also important to investigate performance gain around saturation points i.e. points after which further increase of the MPI process count results in either stagnation or drop of computer program speed-up. It is worth pointing out that, in our case, it becomes difficult to decide where saturation points locate because of jagged behavior of speed-up curves. For this reason, a careful analysis of each performance graph was performed based on values of speed-up, efficiency and our subjective opinion. The results are summarized in Tables 5.9 and 5.10 where each table is organized as follows. Each row of a table contains five fields and provides information about parallel performance of MUMPS at the saturation point of a test-case relatively to the *spread* pinning strategy. The first one is the name of a test-case. The second, fourth and fifth ones show values of the MPI process count, speed-up and parallel efficiency at the saturation point, respectively. The third field shows a gain in parallel factorization time of the *spread* pinning strategy over the *close* one at this point, in percent.

HW1					HW2				
Matrix Name	MPI	Gain w.r.t "close", %	Speed up	Efficiency	MPI	Gain w.r.t "close", %	Speed up	Efficiency	
pwr-3d	4	11.594	1.386	0.347	4	6.616	1.626	0.406	
cube-5	4	8.261	1.139	0.285	4	10.640	1.156	0.289	
cube-64	8	5.645	1.812	0.226	8	7.521	1.729	0.216	
cube-645	6	9.985	2.152	0.359	8	9.078	2.521	0.315	
k3-2	7	7.788	2.899	0.414	8	9.947	3.298	0.412	
k3-18	8	6.716	3.472	0.434	8	9.567	3.896	0.487	

Table 5.9.: Comparisons of MUMPS parallel performance at the saturation points in case of factorization of GRS matrix set

A study of Tables 5.9 and 5.10 reveals that HW2 machine performs slightly better in contrast to HW1 one with respect to parallel performance around the saturation points. This results are different from the overall performance gain mentioned above, however, they reflect the same trend. Additionally, it can be clearly observed that increase of NUMA domains always results in improving efficiency and speed-up of MUMPS.

In this subsection, we have shown an influence of different MPI process distributions and the number of NUMA domains on MUMPS parallel performance. We have observed that application of the *spread* process distribution is always advantageous together with increase of the number of NUMA domains.

HW1					HW2				
Matrix Name	MPI	Gain w.r.t "close", %	Speed up	Efficiency	MPI	Gain w.r.t "close", %	Speed up	Efficiency	
cant	8	7.914	3.297	0.412	8	12.437	3.407	0.426	
consph	15	0.110	6.147	0.410	15	2.409	6.667	0.444	
CurlCurl_3	19	8.051	8.249	0.434	20	17.908	11.039	0.552	
Geo_1438	13	21.609	4.548	0.350	ROM	ROM	ROM	ROM	
memchip	9	11.290	4.299	0.477	9	11.102	4.213	0.468	
PFlow_742	19	17.921	8.106	0.427	20	20.469	9.798	0.490	
pkustk10	17	-0.664	3.872	0.228	17	-1.108	4.036	0.237	
torso3	18	5.607	8.149	0.453	19	6.028	9.493	0.499	
x104	6	9.537	1.789	0.298	6	7.829	1.763	0.294	

Table 5.10.: Comparisons of MUMPS parallel performance at the saturation points in case of factorization of SuiteSparse matrix set, where ROM stands for Run Out of Memory

The result of this study can be relevant for energy-efficient parallel computing where strong requirements to program efficiency are applied. This fact usually forces the user to reduce the process count and go away from the saturation point in order to keep values of efficiency around **0.7-0.8**. In this case, performance of MUMPS can be improved by **15-20%** in contrast to a straightforward process pinning i.e. *close* strategy.

Taken into account results of testing, *spread*-pinning has been chosen for the rest of the study. This process distribution can be easily achieved by means of some advanced OpenMPI command line options, for example *-rank-by* and *-bind-to*, as following:

```
1 mpiexec --rank-by numa --bind-to core -n $num_proc $executable_name
$parameters
```

Listing 5.3: An example of setting *spread* process pinning using advanced OpenMPI command line options

#### 5.4.3. Optimized BLAS Implementations

To perform column eliminations of fully summed blocks of type 2 nodes, MUMPS intensively calls GEMM, TRSM and GETRF subroutines [33] which are parts of BLAS and LAPACK libraries, see Figures 5.17 and 5.18 as an example. Additionally, MUMPS

## 5. Configuration of a sparse linear solver

---

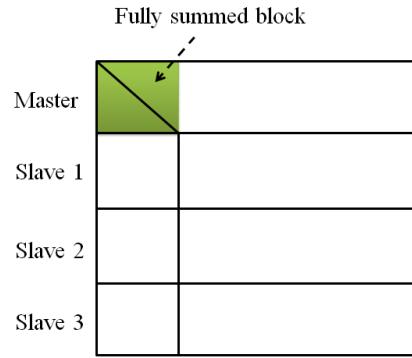


Figure 5.17.: One dimensional block column distribution of a type 2 node in MUMPS

uses ScaLAPACK library subroutines to perform parallel factorization of the root [37], according to the procedure described in Section 5.3.

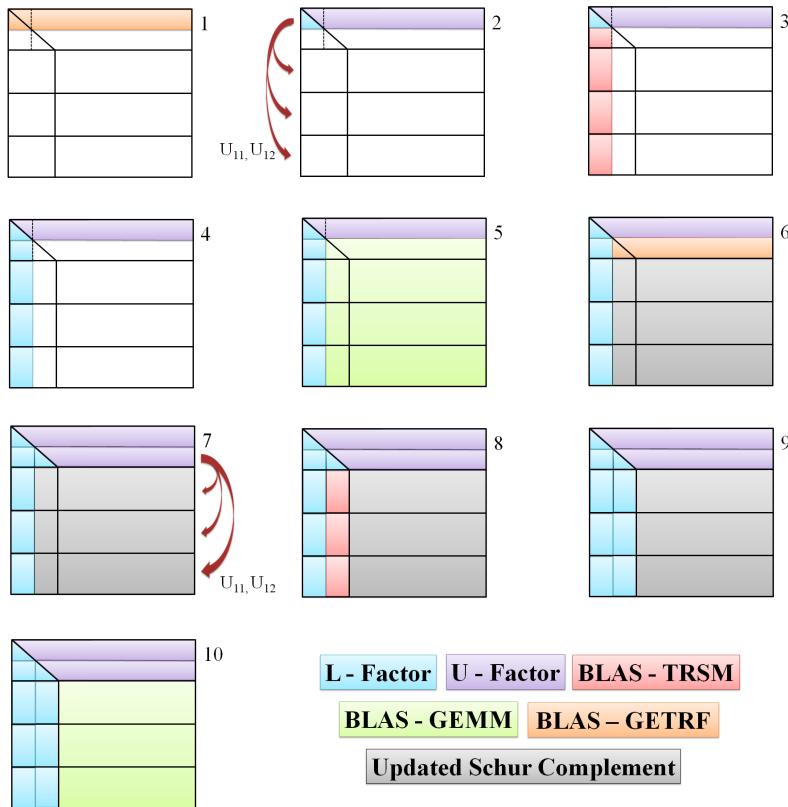


Figure 5.18.: An example of type 2 node factorization implemented in MUMPS

BLAS, LAPACK and ScaLAPACK libraries originate from the Netlib project which is a repository of numerous scientific computing software maintained by AT&T Bell Laboratories, the University of Tennessee, Oak Ridge National Laboratory and other scientific communities, [38].

The goal of BLAS library is provision of high efficient implementations of common dense linear algebra kernels achieved by high rates of floating point operations per memory access, low cache and Translation Lookaside Buffer (TLB) miss rates.

Both LAPACK and ScaLAPACK are designed in such a way so that as much as possible computations are performed by calling BLAS subroutines. Figure 5.19 represents software dependencies between the libraries. Hence, this software architecture allows to achieve high computational performance for operations such as LU, QR, SVD decompositions, triangular solve, etc., on modern computers and distributed-memory machines by an efficient implementation of BLAS library. However, the Netlib BLAS implementation is written for an abstract general-purpose central processing unit where hardware parameters are based on market statistics. Therefore, it is not possible to achieve the maximum possible performance on specific hardware.

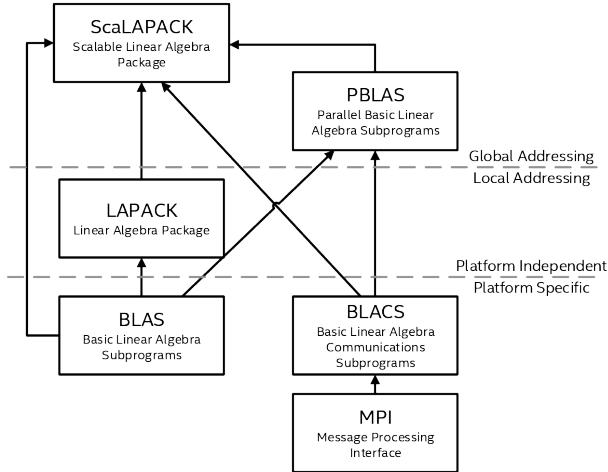


Figure 5.19.: Software dependencies between Netlib libraries, [26]

There exist special-purpose, hardware-specific implementations of BLAS developed by hardware vendors i.e. IBM, Cray, Intel, AMD, etc., as well as open-source tuned implementations such as ATLAS, OpenBLAS, etc. To achieve full compatibility, the developers consider the Netlib implementation as a standard, or a reference, and thus

overwrite all its subroutines with additional tuning and optimization. This approach makes it easy to replace one BLAS implementation by another one by substituting the corresponding object files during the linking stage. As a result, the source code of an application which calls BLAS, LAPACK or ScaLAPACK subroutines remains the same without a need to perform any source code modifications.

Table 5.11 shows commercial and open-source tunned BLAS implementations available on the market today.

Name	Description	License
Accelerate	Apple's implementation for macOS and iOS	proprietary license
ACML	BLAS implementation for AMD processors	proprietary license
C++ AMP	Microsoft's AMP language extension for Visual C++	open source
ATLAS	Automatically tuned BLAS implementation	open source
Eigen BLAS	BLAS implemented on top of the MPL-licensed Eigen library	open source
ESSL	optimized BLAS implementation for IBM's machines	proprietary license
GotoBLAS	Kazushige Goto's implementation of BLAS	proprietary license
HP MLIB	BLAS implementation supporting IA-64, PA-RISC, x86 and Opteron architecture	proprietary license
Intel MKL	Intel's implementation of BLAS optimized for Intel Pentium, Core, Xeon and Xeon Phi	proprietary license
Netlib BLAS	The official reference implementation on Netlib	open source
OpenBLAS	Optimized BLAS library based on GotoBLAS	open source
PDLIB/SX	BLAS library targeted to the NEC SX-4 system	proprietary license
SCSL	BLAS implementations for SGI's Irix workstations	proprietary license
Sun Performance Library	Optimized BLAS and LAPACK for SPARC, Core and AMD64 architectures under Solaris 8, 9, and 10 as well as Linux	proprietary license

Table 5.11.: Commercial and open source BLAS libraries [46]

Among all libraries listed in Table 5.11 there were only four available in HW1 machine environment, namely: Netlib BLAS, Intel MKL, OpenBLAS and ATLAS. However, installation of ATLAS requires to switch off dynamic frequency scaling, also called CPU throttling, to allow ATLAS configuration routines to find the best loop transformation parameters for specific hardware. In order to turn off CPU throttling, one has to reboot the entire machine and make appropriate changes in Basic Input/Output System (BIOS). This fact made ATLAS library not suitable for the study and we excluded it from our primary list of candidates. Moreover, during installation, one has to explicitly specify the number of OpenMP threads that are going to be forked once a BLAS subroutine is called. This means there is no way to change the number of threads per MPI process in run-time without re-installation of the library. Thus, only 3 versions of MUMPS-PETSc (linked with Netlib BLAS, Intel MKL and OpenBLAS) library were compiled, installed and tested using both GRS and SuiteSparse matrix sets and 1 thread per MPI process i.e. flat-MPI mode. Results of testing were obtained on HW1 machine and are represented in Figures 5.20, 5.21 and appendix E.

The tests show that OpenBLAS outperforms both Netlib and Intel MKL libraries in case of GRS matrix set. On average, OpenBLAS is about **13%** faster than the default Netlib implementation and approximately **21%** faster than Intel MKL library. It is interesting to notice that Intel MKL library turns out to be slower than the default Netlib BLAS implementation for small- and medium-sized GRS matrices in almost **52%** and **2%**, respectively. At the same time, both tuned libraries, OpenBLAS and Intel MKL, show significant performance gain in comparison to the standard Netlib BLAS implementation in case of SuiteSparse matrix set. The libraries reduce the execution time by almost **50%** on an average. In opposite to GRS matrix set, it turns out that Intel MKL is often faster than OpenBLAS for almost all test-cases from SuiteSparse matrix set. However, the difference between them is negligibly small. The result of the comparison are summarized in Tables 5.12 and 5.13.

It can be clearly observed from the tables that test-cases derived from GRS matrix set demonstrate insignificant improvements in execution time in contrast to the tests generated with SuiteSparse matrix set. This may be explained by relatively small numbers of type 2 nodes in assembly trees resulted from GRS test-cases. In this case, the trees are mainly formed with the root and type 1 nodes. As it was mentioned in Section 5.3, type 1 nodes are grouped in subtrees and each subtree is processed by a single MPI process. According to the documentation, it is not clear whether MUMPS calls BLAS subroutines while processing a type 1 node. Even if it is a case performance of BLAS can be limited because of small sizes of frontal matrices of such nodes.

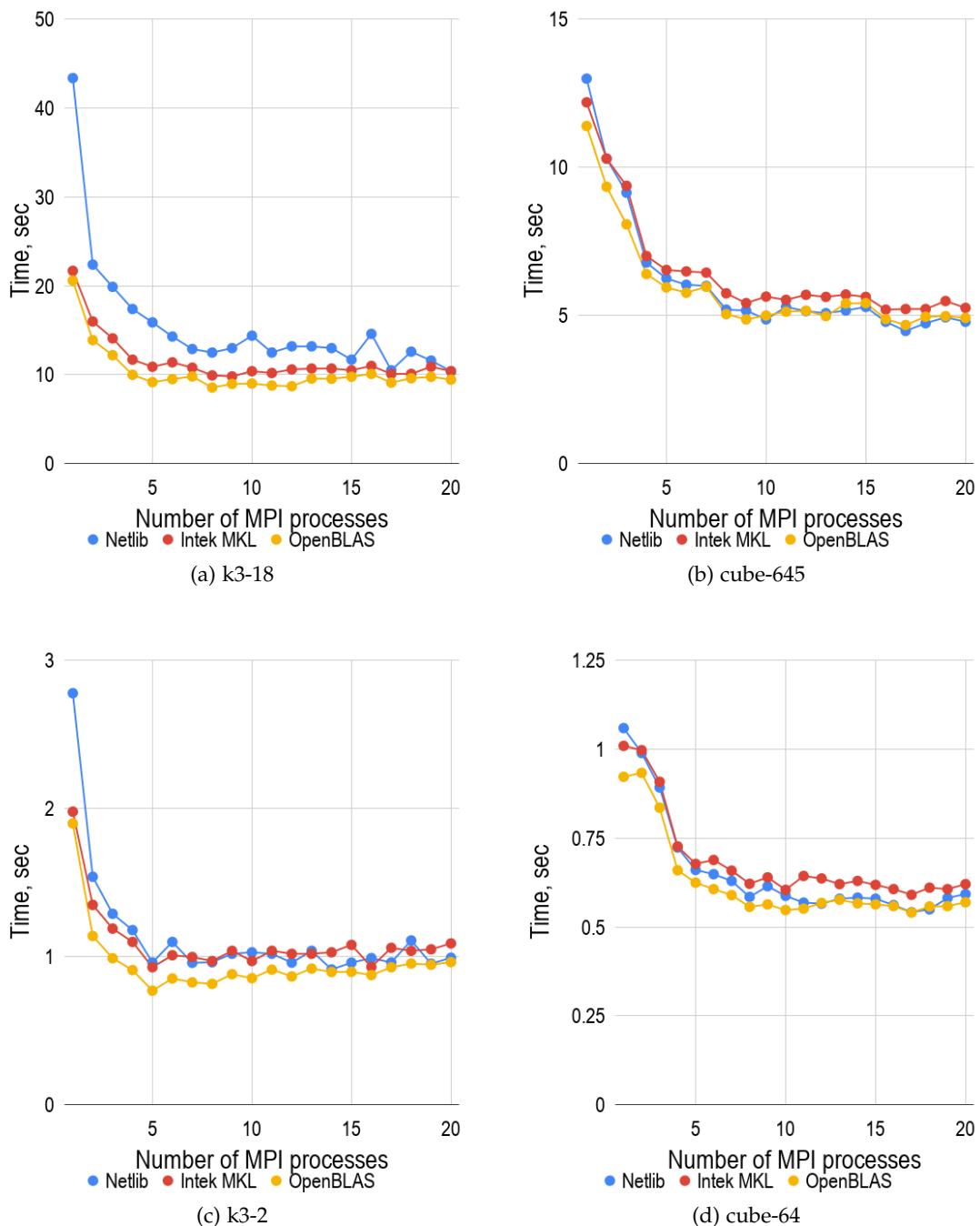


Figure 5.20.: Comparisons of parallel factorization of GRS matrix set performed on HW1 machine using MUMPS solver linked to different BLAS implementations

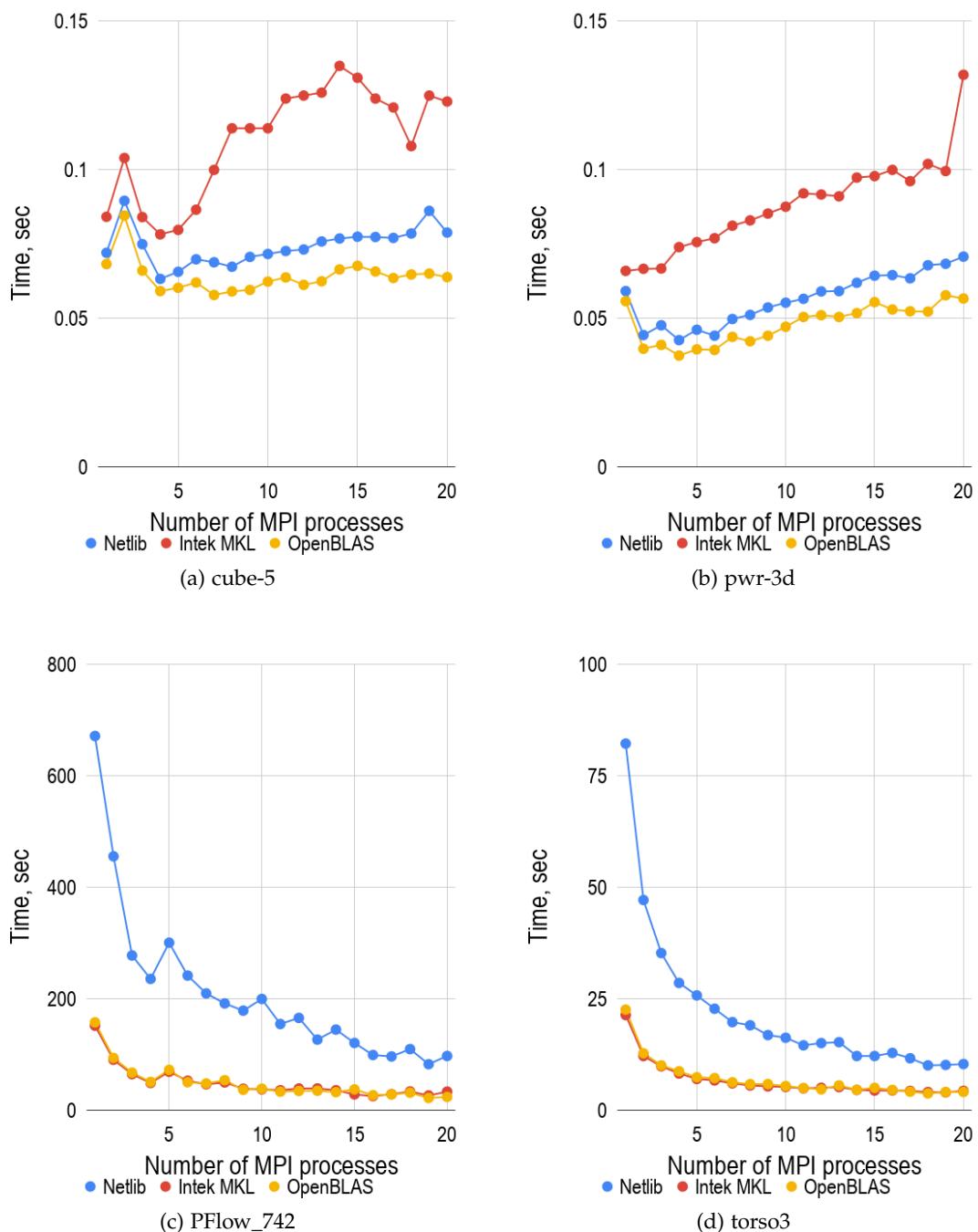


Figure 5.21.: Comparisons of parallel factorizations of GRS and SuiteSparse matrix sets performed on HW1 machine using MUMPS solver linked to different BLAS implementations

5. Configuration of a sparse linear solver

---

Matrix Name	Average performance gain of OpenBLAS relatively to Netlib %	Average performance gain of IntelMKL relatively to Netlib %	Average performance gain of OpenBLAS relatively to Intel MKL %
pwr-3d	14.607	-56.249	44.695
cube-5	13.569	-47.797	39.931
cube-64	4.385	-5.483	9.323
cube-645	1.897	-7.474	8.702
k3-2	13.906	0.833	13.057
k3-18	29.914	21.03	11.29

Table 5.12.: Comparisons of different MUMPS-BLAS configurations applied to GRS matrix set

Matrix Name	Average performance gain of OpenBLAS relatively to Netlib %	Average performance gain of IntelMKL relatively to Netlib %	Average performance gain of OpenBLAS relatively to Intel MKL %
cant	26.981	25.964	1.233
consph	67.617	68.252	-2.327
CurlCurl_3	78.804	79.37	-3.371
Geo_1438	83.106	83.565	-2.857
memchip	6.066	-6.909	11.883
PFlow_742	75.574	74.943	1.416
pkustk10	35.089	34.536	0.502
torso3	66.185	66.988	-2.837
x104	41.82	41.936	-0.445

Table 5.13.: Comparisons of different MUMPS-BLAS configurations applied to SuiteSparse matrix set

We assume that, in the general case, lack of type 2 nodes in an assembly tree can be due to an inefficient amalgamation process of the corresponding elimination tree resulted from the matrix sparsity pattern.

Based on the obtained results, comparisons between different matrix sets and our reasoning, we presume that ATHLET generates linear systems resulting in such trees where type 1 nodes predominate over the others. We can assume it is due to specifics of numerical spacial and time integration explained in Section 2.1.

In this subsection, we have discussed where and how MUMPS utilizes BLAS, LAPACK and ScaLAPACK libraries. We have compared two tuned BLAS implementations with the baseline, Netlib BLAS, using two different matrix sets. We have shown the overall statistics of the obtained results and come to the conclusion that MUMPS-OpenBLAS configuration is the best one for GRS matrix set. Additionally, we have given reasoning for a noticeable difference between results obtained from different matrix sets as well as we have talked about probable specifics of linear systems generated by ATHLET.

#### 5.4.4. Hybrid MPI/OpenMP Computing

As it was mentioned in Section 5.3, the development of MUMPS began in 1996 when message-passing programming paradigm dominated in parallel computing. Therefore, the library originally was designed only for distributed-memory machines.

In 2010, Chowdhury and L'Excellent published their first experiments and some issues, in [10], of exploiting shared memory parallelism in MUMPS. The authors showed that it was possible to achieve some improvements in multicore systems using multi-threading, given a purely MPI application. However, later L'Excellent and Sid-Lakhdar mentioned, in [33], that adaptation of the existing code for NUMA architecture was still a challenge because of memory allocation, memory affinity, thread pinning and other related issues.

In spite of an advantage of natural data locality of message-passing applications, a general motivation for switching to a hybrid mode, a mixed MPI/OpenMP process/thread distribution, is to reduce communication overheads between MPI processes. According to the profiling results obtained by Chowdhury and L'Excellent, MUMPS contained four main regions of shared-memory parallelization, namely:

1. BLAS Level 1, 2, 3 operations during both factorization and solution phases
2. Assembly operations, where contribution blocks of children nodes are assembled at the parent level

3. Copying contribution blocks during stacking operations
4. Pivot search operations

Almost all customized BLAS libraries, for example Intel MKL and OpenBLAS, are multi-threaded and can efficiently work in shared-memory environment. Hence, parallelization of region 1 can be achieved by linking a suitable BLAS library whereas regions 2, 3 and 4 can be multi-threaded by inserting appropriate OpenMP directives above the corresponding loop statements.

A detailed review of works [33] and [10] reveals that, in general, a pure OpenMP or mixed MPI/OpenMP strategy can reduce run-time of MUMPS. On average, factorization time is reduced by **14.3%** and in some special cases improvements reach about **50.4%**, according to analysis performed on data provided in the papers. However, at the same time, the results also show that sometimes a flat-MPI mode can significantly outperform other hybrid mixed strategies.

By and large, the results show two important aspects. Firstly, performance of a specific strategy depends heavily on a resulting assembly tree and thus on a matrix sparsity pattern and applied fill reducing reordering. Secondly, it is not possible to guess in advance which strategy gives the best parallel performance without detailed information about the tree structure and computational cost per node. L'Excellent and Sid-Lakhdar showed that performance of a particular mode dependeds on a ratio of large and small fronts. For example, they noticed that more threads per MPI process resulted in better parallel performance in case of high ratios. On the other hand, they observed the absolutely opposite result with relatively small ratios. Unfortunately, L'Excellent and Sid-Lakhdar did not provide any quantitative measure for the notion of small and large ratios in [33].

It is also interesting to notice that parallelization of region 1 using a multi-threaded BLAS library gives the most of the parallel performance improvement for mixed or pure OpenMP strategies, according to analysis of results from [33]. Whereas, multi-threading of regions 2, 3, 4 has only a small positive effect i.e it reduces numerical factorization run-time by only **0.66%** on an average.

This outcome is expected because BLAS subroutines, especially level 3, re-use data stored in caches as much as possible and thus achieve high ratios of floating point operations per memory access which is essential for efficient multi-threading. Meanwhile, regions 2, 3, 4 mainly perform initialization of variables, data movements and

executions of *if-statements* which always result in low computational intensity.

We have to admit that both works, [10] and [33], are relatively old and the analysis above may be not complete and full. Because MUMPS is a dynamic developing project, we can expect that adaptation of shared-memory parallelization in MUMPS has been significantly advanced since that time. Since the 4th release of MUMPS library, the developers have persistently recommended to use only hybrid strategies e.g. *one MPI process per socket and as many threads as the number of cores*, [37].

As an initial test, we compared an influence of both Intel MKL and OpenBLAS libraries on parallel performance of MUMPS using GRS matrix set only. In order to pin OpenMP threads in a correct way, without any conflicts between them, the following OpenMP environment variables were set as follows:

- OMP\_PLACES=cores
- OMP\_PROC\_BIND=spread

During the testing, we found that sometimes execution time of MUMPS-OpenBLAS configuration abnormality increased. For instance, in case of parallel factorization of matrix *cube-645*, the increase reached almost **450%** in contrast to the pure sequential execution.

Multiple conflicts between application and system threads were observed using *htop* software as an interactive process viewer. Figure 5.23 shows a snapshot taken during factorization of matrix *k3-18* running with 1 MPI process and 20 threads.

It is difficult to state what exactly caused such behavior. However, Chowdhury and L'Excellent also reported about the same problem using GotoBLAS (OpenBLAS). They assumed that GotoBLAS created and kept some threads active even after the main threads returned to the calling application which could lead to interference with threads created in other OpenMP regions [10]. For this reason, we decided to use only Intel MKL library for the rest of the study because there were no such thread-conflicts detected during operation of MUMPS-Intel MKL configuration.

Only common mixed MPI/OpenMP strategies were tested in order to check an influence of shared-memory parallelism on parallel performance of MUMPS as well as to limit an amount of testing. The following strategies were chosen: 20 MPI - 1 thread (flat-MPI), 10 MPI - 2 threads, 4 MPI - 5 threads, 2 MPI - 10 threads, 1 MPI - 20 threads (flat-OpenMP). The tests were conducted on both HW1 and HW2 machines

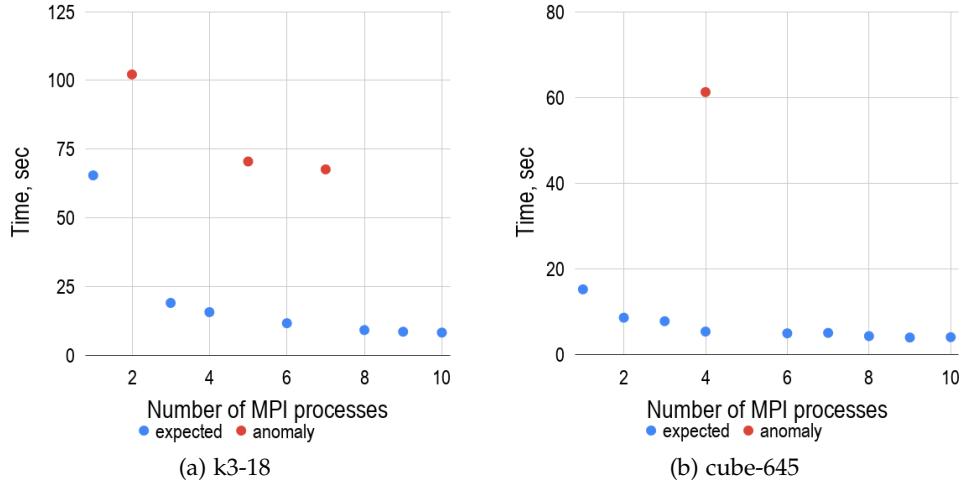


Figure 5.22.: Anomalies of parallel executions of MUMPS-OpenBLAS configuration during factorizations of large-sized GRS matrices running with 2 OpenMP threads per MPI process

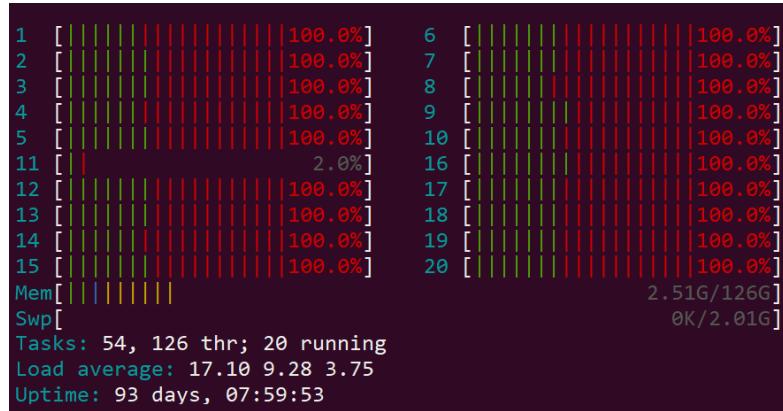


Figure 5.23.: Thread conflicts of MUMPS-OpenBLAS configuration detected during parallel factorization of matrix  $k3-18$ , where green - application threads, red - system threads

with the aim of checking whether results would be consistent between different hardware running under different operating and environment settings. Results of testing are represented in Tables 5.14 5.15, 5.16 and 5.17 where numerical values are given in seconds.

Matrix Name	20 MPI 1 thread	10 MPI 2 threads	4 MPI 5 threads	2 MPI 10 threads	1 MPI 20 threads	Gain w.r.t. flat-MPI
k3-18	<b>12.520</b>	12.630	14.010	18.020	19.170	-
k3-2	1.341	<b>1.250</b>	1.470	1.671	2.052	1.073
cube-645	<b>6.585</b>	6.859	8.552	12.010	14.080	-
cube-64	0.756	<b>0.749</b>	0.874	1.178	1.354	1.010
cube-5	0.181	0.132	<b>0.104</b>	0.126	0.117	1.744
pwr-3d	0.130	0.114	0.0972	<b>0.077</b>	0.109	1.691

Table 5.14.: Compassions of different hybrid MPI/OpenMP modes used for parallel factorization of GRS matrix set on HW1

Matrix Name	20 MPI 1 thread	10 MPI 2 threads	4 MPI 5 threads	2 MPI 10 threads	1 MPI 20 threads	Gain w.r.t. flat-MPI
k3-18	8.558	<b>7.819</b>	8.165	11.330	14.320	1.095
k3-2	1.168	<b>0.788</b>	0.956	1.131	1.651	1.482
cube-645	5.735	<b>4.859</b>	6.069	9.360	11.040	1.180
cube-64	0.805	<b>0.541</b>	0.664	0.947	0.918	1.490
cube-5	0.241	0.121	<b>0.093</b>	0.129	0.126	2.582
pwr-3d	0.234	0.095	0.098	<b>0.070</b>	0.094	3.341

Table 5.15.: Compassions of different hybrid MPI/OpenMP modes used for parallel factorization of GRS matrix set on HW2

According to analysis of obtained results and flat-MPI performance graphs from Subsection 5.4.3, we have noticed that an optimal hybrid MPI/OpenMP mode locates near the saturation point of the corresponding flat-MPI test. Generally speaking, a location of the saturation point is specific for each matrix and, therefore, there is no way to predict a mode in advance. However, having known the point, an amount of testing can be considerably reduced by searching around and applying different mixed MPI/OpenMP strategies.

The results show that average performance gain is around 2.1% in case of GRS matrix set for HW1 hardware, excluding small test-cases such as *cube-5* and *pwr-3d* from statistics. We consider these two scenarios, *cube-5* and *pwr-3d*, as specific ones because their execution time with 20 MPI processes using flat-MPI mode is originally slower in contrast to the sequential execution and, therefore, it is relevant to assume

## 5. Configuration of a sparse linear solver

---

Matrix Name	20 MPI 1 thread	10 MPI 2 threads	4 MPI 5 threads	2 MPI 10 threads	1 MPI 20 threads	Gain w.r.t. flat-MPI
cant	1.400	<b>0.990</b>	1.050	1.605	2.019	1.414
consph	3.495	<b>2.652</b>	3.015	3.706	3.714	1.318
memchip	<b>7.470</b>	9.080	13.301	20.198	45.800	-
PFlow_742	26.802	24.204	<b>21.897</b>	30.389	54.501	1.224
pkustk10	<b>0.748</b>	0.879	0.972	1.459	1.280	-
torso3	<b>3.922</b>	4.285	4.642	5.603	8.144	-
x104	<b>1.597</b>	1.644	2.024	3.208	2.167	-
CurlCurl_3	49.250	44.120	<b>39.909</b>	43.311	63.001	1.234
Geo_1438	478.101	234.697	<b>151.603</b>	157.697	158.102	3.154

Table 5.16.: Compassions of different hybrid MPI/OpenMP modes used for parallel factorization of SuiteSparse matrix set on HW1

Matrix Name	20 MPI 1 thread	10 MPI 2 threads	4 MPI 5 threads	2 MPI 10 threads	1 MPI 20 threads	Gain w.r.t flat-MPI
cant	2.128	<b>0.955</b>	1.011	1.577	2.058	2.229
consph	3.840	<b>2.852</b>	3.111	3.695	3.897	1.346
memchip	<b>7.811</b>	7.816	9.811	15.160	31.969	-
PFlow_742	24.190	29.241	<b>19.686</b>	27.530	55.431	1.230
pkustk10	1.373	<b>0.904</b>	1.022	1.421	1.403	1.520
torso3	4.733	<b>4.080</b>	4.483	5.648	8.217	1.160
x104	2.676	<b>1.597</b>	2.025	3.204	2.133	1.676
CurlCurl_3	39.890	<b>34.579</b>	38.620	41.171	67.760	1.154
Geo_1438	ROM	ROM	ROM	ROM	ROM	ROM

Table 5.17.: Compassions of different hybrid MPI/OpenMP modes used for parallel factorization of SuiteSparse matrix set on HW2, where ROM stands for Run Out of Memory

that the improvement came only from reducing of the MPI process count. At the same time, much optimistic results were obtained from experiments conducted on HW2 machine where performance gain reached almost 31% for the same test-cases.

Results obtained with SuiteSparse matrix set demonstrate much better performance improvements from hybrid parallel computing obtained on both hardware. On average, execution time improves by more than 15% running tests on HW1 and approximately by 41% on HW2, excluding *Geo\_1438* from the statistics. The best result was obtained

exactly in case of *Geo\_1438* test-case on both machines where execution time dropped about **3 times** for all hybrid modes in contrast to the corresponding flat-MPI one. We assume it may occur because of a high ratio of large and small fronts of this particular test-case.

According to the outcomes of testing, we have observed a negligible improvement in MUMPS parallel performance from the application of the multi-threaded Intel MKL BLAS library to GRS matrix set. Such unimpressive results can be explained with the same reasoning given in Subsection 5.4.3 i.e. lack of type 2 nodes. Moreover, in case of GRS matrix set, parallel efficiency drops significantly probably due to inefficient utilization of additional processing elements i.e cores. However, at the same time, results obtained using SuiteSparse matrix set have shown an advantage of hybrid parallel computing, especially in case of *Geo\_1438* matrix factorization.

These contradictory results obtained from two different matrix sets second our reasoning about specifics of linear systems generated by ATHLET software. Again, we presume that assembly trees resulted from GRS matrices are mostly formed with subtrees filled with type 1 nodes where each subtree is processed by a single MPI process. Hence, parallel factorization of GRS matrices mainly gets benefit from MPI parallelization that can be clearly observed from the results.

In this subsection, we have discussed how MUMPS adopts hybrid parallel computing. As it is in case of fill reducing reordering algorithm selection, Subsection 5.4.1, it is not possible to find an optimal mixed MPI/OpenMP strategy in advance without performance testing. We have come to the conclusion that flat-MPI mode is the best one for GRS matrix set and provided our reasoning for that. Generally speaking, there are 3 reasons to use this mode in our case. Firstly, the mode always resulted in more efficient hardware utilization. Secondly, MUMPS-Intel MKL configuration running with optimal hybrid MPI/OpenMP strategies can deteriorate performance gain obtained with MUMPS-OpenBLAS flat-MPI configuration, shown in Subsection 5.4.3. Finally, efficient utilization of flat-MPI strategy only demands to find an optimal MPI process count i.e the saturation point on a performance graph. Hence, it leads to a significant reduction of testing due to a reduced number of parameters which are needed to be taken into account.

## 5.5. Results

Figures 5.24 and 5.25 show comparisons of MUMPS parallel performance before and after applications of the optimal MUMPS-MPI parameter settings, found in Subsections 5.4.1, 5.4.2, 5.4.3, and 5.4.4, to GRS matrix set. Results labeled as *default* were obtained using the fill reducing reordering algorithm provided by ParMetis library because it had been used by ATHLET users before the current study.

On average, factorization time is reduced by **51.4%** for small-sized linear systems, *cube-5* and *pwr-3d*. As it was expected, the most significant performance gain mainly comes from a correct choice of a fill reducing reordering algorithm. Moreover, the application of PT-Scotch to these systems of equations results in a drastic change of strong scaling behavior, see Figures 5.24a and 5.24b, which allows to reduce execution time by approximately **17%** in contrast to the sequential execution of MUMPS running with the default parameters.

Execution time spent on factorizations of medium-sized systems, such as *cube-64* and *k3-2*, drops in **1.4** times on an average. We have noticed that strong scaling of *cube-64* matrix factorization considerably improves. Additionally, the application of PT-Scotch to *cube-64* test-case results in shifting of the optimal MPI process count, the saturation point, from 5 to 10 and, as a result, it reduces execution time of parallel factorization. The application of all optimal settings results in reduction of execution time around the corresponding saturation points by almost **31%** on average for this type of GRS matrices.

Improvements in parallel factorization of large-sized GRS systems comes only from optimal processes pinning and configuration of MUMPS with OpenBLAS library because of usage of the same fill reducing reordering algorithm, namely: ParMetis, according to the assignment Table 5.7. On average, performance increased by almost **20%** in case of *k3-18* test-case and only by **1.3%** for *cube-645* one. This difference in results can be explained by the fact that the assembly tree of *cube-645* test-case may lack type 2 nodes. However, the saturation points of both test-cases are shifted towards lower values of the MPI process count which result in a considerably improvement of hardware utilization. For example, a detailed study of *k3-18* performance graph, Figure 5.25b, shows the optimal value of the MPI process count decreases from 17 to 8 and, at the same time, execution time drops by almost **19%**. These two effects result in almost **13%** jump of parallel efficiency. The same trend can be observed for *cube-645* test-case as well.

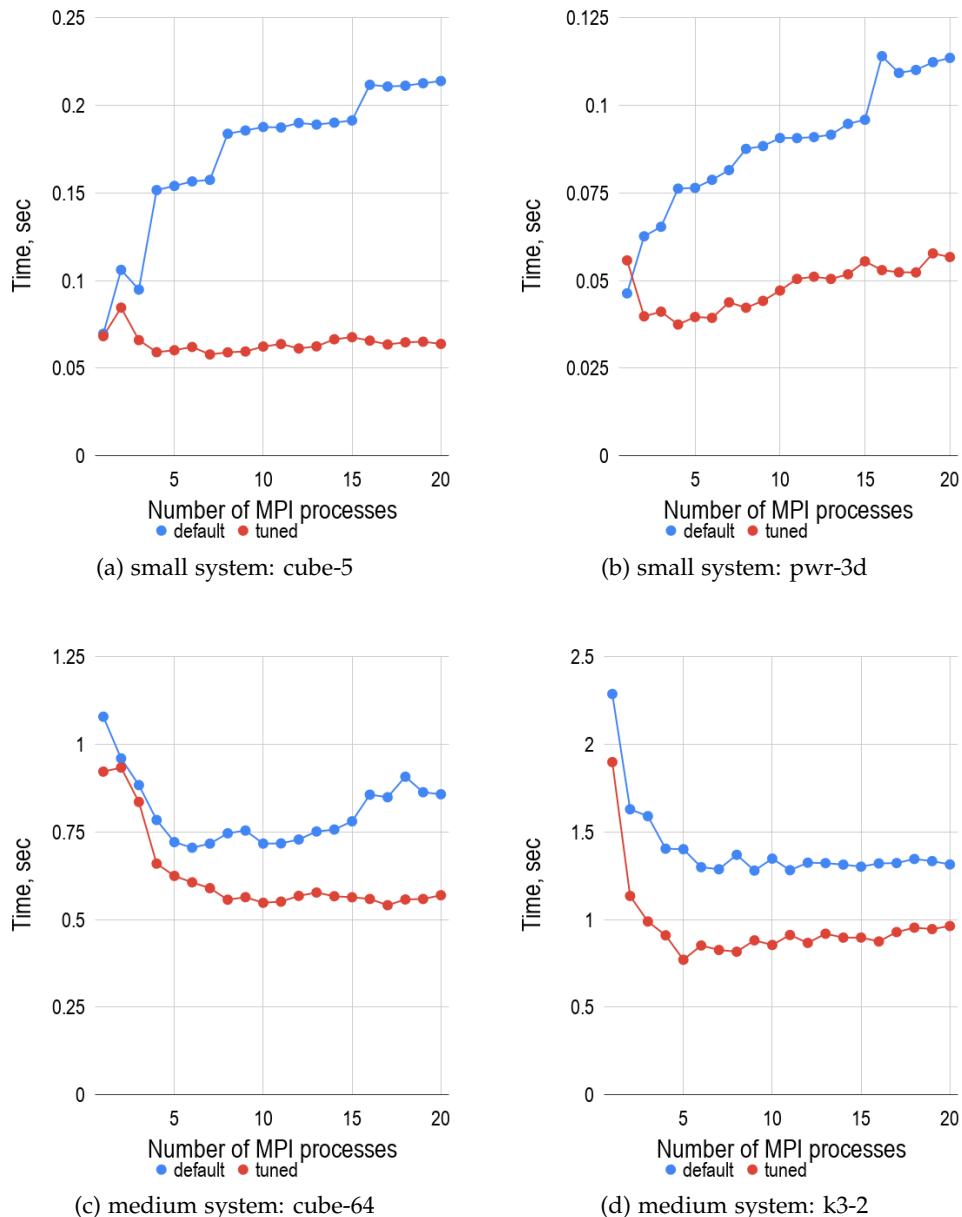


Figure 5.24.: Comparisons of parallel factorizations of small- and middle-sized GRS matrices between applications of the default and optimal MUMPS configurations

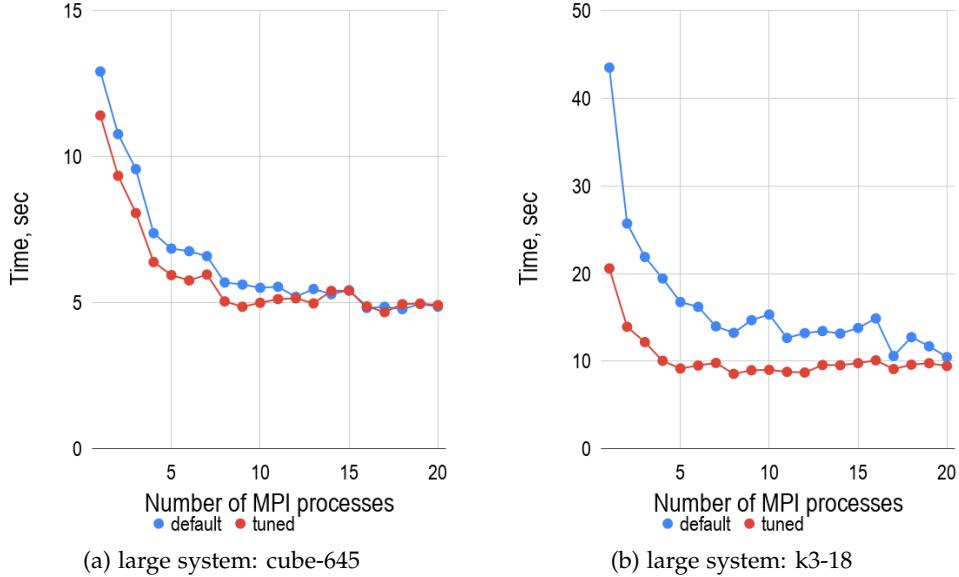


Figure 5.25.: Comparisons of parallel factorizations of large-sized GRS matrices between applications of the default and optimal MUMPS configurations

By and large, in this subsection we have shown that applications of the optimal parameter settings to MUMPS lead to total accumulative improvements in both factorization time and hardware utilization.

## 5.6. Conclusion

In this chapter, we have examined different types of sparse linear solvers applied to linear systems generated by ATHLET software resulting from numerical integration of thermo-hydraulic computations. We have come to the conclusion that, in spite of better scalability and parallel efficiency of iterative methods due to efficient data-based parallelism, direct sparse linear solvers are well suitable for this purpose because of their robustness, see Subsection 5.1.3.

In Section 5.2, we tested different direct sparse solvers, namely: MUMPS, SuperLU\_DIST and PaStiX. MUMPS showed better parallel performance among the others according to the results of testing and, therefore, was chosen for the following study where we mainly focused on performance tuning of the library.

We have shown in subsequent subsections there have been four main sources of library parameter tuning, namely:

1. correct selection of a fill reducing reordering algorithm
2. distribution of MPI processes among multiple NUMA domain within a compute node
3. configuration of MUMPS with an optimal, tuned BLAS library implementation
4. execution of MUMPS with optimal hybrid MPI/OpenMP process/thread distributions

Testing was performed using two different matrix set, GRS and SuiteSparse, on two different computer-clusters, HW1 and HW2, see Chapter 4, in order to check consistency of obtained results. In this section, we give most general conclusions relevant to only GRS matrix set and HW1 cluster as targets of the study. The reader can become familiar with detailed conclusions relative to both matrix sets and hardware given at the end of each subsection that we are going to reference to.

1. In Subsection 5.4.1, it has been shown that parallel performance of MUMPS is quite sensitive to applied fill-in reducing reordering algorithms. A correct choice of the algorithm can lead to a significant improvement in execution time and strong scaling behavior. We have noticed that MUMPS performs factorizations of small- and medium-sized matrices faster using PT-Scotch library whereas large-sized problems tend to get benefit from the algorithm provided by ParMetis. We assume that the obtained conclusion may be inaccurate due to a small size of GRS matrix set. At the moment of writing, we have not found any indirect metric to predict a correct choice of an algorithm beforehand. Hence, we encourage ATHLET users to perform similar testing described in the subsection before running thermo-hydraulic simulations on distributed-memory machines to achieve better performance of parallel computations.
2. In Subsection 5.4.2, an influence of different process pinning strategies on MUMPS parallel performance has been investigated. The tests have shown that an equal distribution of MPI process among all available NUMA domains always results in additional performance gain.
3. In Subsection 5.4.3, we tested MUMPS configured with 3 different implementations of BLAS library, namely: Netlib, OpenBLAS and Intel MKL. The results have shown

the application of OpenBLAS library always results in better parallel performance.

4. In Subsection 5.4.4, we have investigated an impact of various MPI/OpenMP process/thread distributions on parallel factorizations of GRS matrices within a compute-node. We have observed that multi-threading of OpenBLAS library in MUMPS leads to multiple thread conflicts which sometimes result in significant slow-down of the solver. Results obtained with MUMPS-Intel MKL configuration have demonstrated a negligible improvement in solver execution time resulting in a significant parallel efficiency drop, probably due to inefficient usage of additional processing elements utilized by forked Intel MKL threads. At the end, we have concluded that flat-MPI mode is the best one for matrices generated by ATHLET software.

In Section 5.5, we have studied the overall impact of introduced configuration changes found in Subsections 5.4.1, 5.4.2, 5.4.3 and 5.4.4. Testing shows the changes result in a positive accumulative effect leading to considerable improvements of both factorization time and hardware utilization.

During the study, we have noticed that optimal values of the MPI process count lay within the range between 1 and 4 in case of small-sized GRS matrices and 4 and 8 for middle- and large-sized problems. An exact value is impossible to predict beforehand and, therefore, it always demands individual, problem-specific testing.

## 6. Improvement of ATHLET-NUT Communication

### 6.1. Jacobian Matrix Compression

The main goal of Jacobian matrix compression is minimization of a number of non-linear function evaluations which are usually quite expensive computational operations. The minimization is performed by means of efficient treatment of non-zero entries of a sparse matrix. The problem is also known as matrix partitioning.

In the general case, a finite difference method can be used to compute a Jacobian matrix approximation in the following way:

$$\frac{1}{\epsilon}(F(y + \epsilon e_k) - F(y)) \approx J(y)e_k, \quad 1 \leq k \leq N \quad (6.1)$$

where  $F : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is a non-linear function;  $e_k \in \mathbb{R}^N$  is the  $k$ th coordinate unit vector,  $\epsilon$  is a small step size.

Equation 6.1 does not exploit Jacobian matrix sparsity and thus such estimation of the Jacobian matrix requires  $N$  function evaluations.

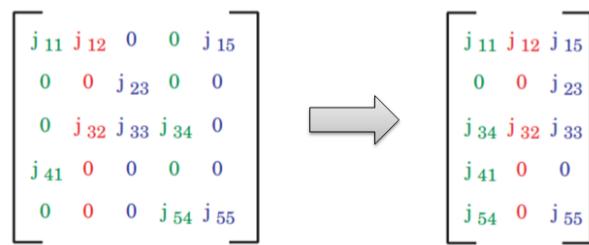


Figure 6.1.: An example of matrix coloring and compression, [19]

A compression algorithm is based on a notion of *structurally orthogonal* columns i.e. columns which do not share any non-zero entry in a common row. Figure 6.1 shows

an example of matrix compression where each color denotes independent *structurally orthogonal* columns.

Having obtained a compressed form of the Jacobian, another set of vectors  $d \in \mathbb{R}^N$ , also known as seed vectors, can be used to perform function perturbations instead of unit vectors  $e_k$ . A seed vector  $d$  has 1's in components corresponding to the indices of columns in a structurally orthogonal group of columns, and zeros in all other components [19]. By differencing the function  $F$  along the vector  $d$ , one can simultaneously determine the nonzero elements in all of these columns through one additional function evaluation at  $F(y + d)$  [19].

It is obvious the algorithm requires to partition a matrix into the fewest amount of groups, colors, in order to achieve the most of efficiency. It means it is a NP-hard problem and, therefore, a heuristical approach is required. Gebremedhin, Manne, and Pothen, in [19], conducted one of the most recent studies in this field and summarized different matrix partitioning algorithms proposed over the last 20 years. Currently, a Jacobian matrix compression algorithm has been successfully implemented in NUT by means of the corresponding built-in PETSc subroutines. The algorithm is used by ATHLET via the corresponding NUT interface.

Figure 6.2 shows an illustrative example of an efficient matrix partitioning where an initial 100 by 100 Jacobian matrix is transformed into its 100 by 28 compressed form using 28 distinct colors. It can be clearly observed from the figure that column vector lengths of the compressed Jacobian form are gradually decreasing. Figure 6.3 provides a detailed and clear view in the problem, using data from Figure 6.2 as an example, where bars represent the corresponding column lengths.

According to the ATHLET-NUT coupling design, each column is transferred to NUT by means of the synchronous 3-way handshake procedure, described in Section 2.3, immediately after its evaluation. Thus, Figure 6.3 determines the communication pattern during the Jacobian matrix transfer for the example shown in Figure 6.2.

Code Listings 6.1 and 6.2 represent the default implementation of a compressed Jacobian matrix transfer between ATHLET and NUT. This code is used as a baseline for the remaining part of the study.

All **code listings**, presented in this part of the study, are written in **pseudocode** and intended for convenience of reading. The aim is to show and display the main ideas skipping non-relevant parts of the actual source code. The **pseudocode** is a mixture of

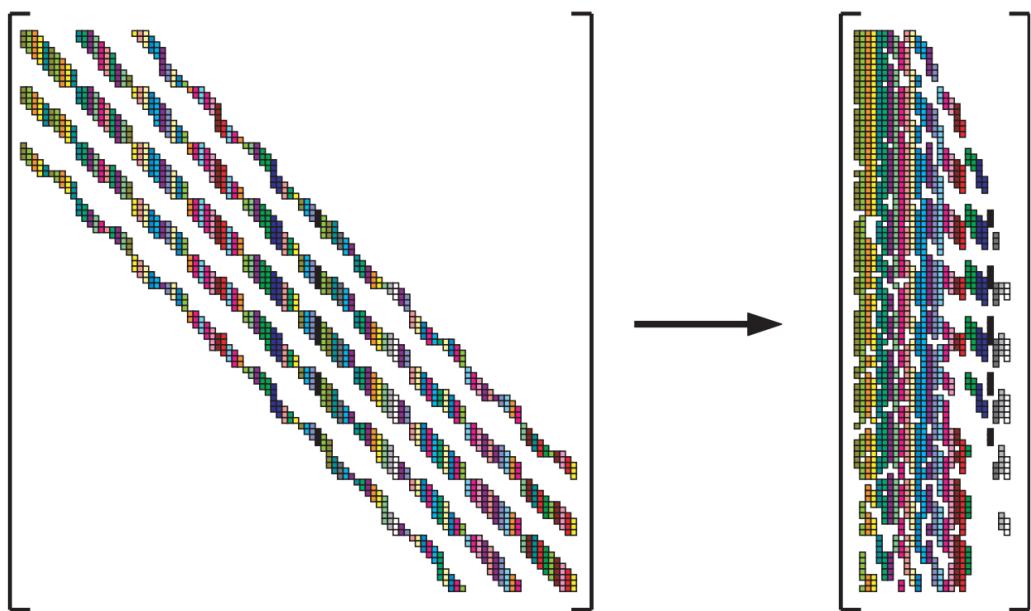


Figure 6.2.: An example of an efficient Jacobian matrix partitioning, [19]

several programming languages, namely: **Python**, **C/C++**, **Fortran**, **(MPI)**.

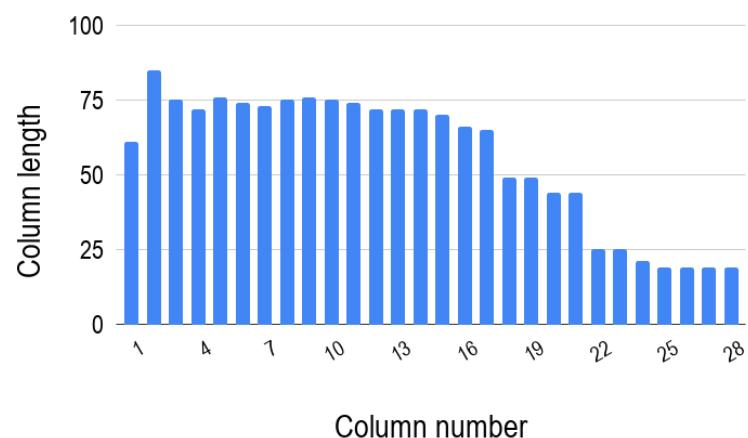


Figure 6.3.: A column-length distribution of the example depicted in Figure 6.2

## 6. Improvement of ATHLET-NUT Communication

---

```
1 # GIVEN PARAMETERS:  
2 # acomm – the athlet communicator  
3 # acomm_id – athlet identification number  
4 # y – known vector  
5 # N – problem size  
6 # COO – compressed matrix coordinate format  
7  
8 eps = 1e-4  
9 center = f(y)  
10 column = zeros(N)  
11  
12 # compute Jacobian and send it to NUT column-by-column  
13 for seed_vector in seed_vectors:  
14  
15     # compute the next column  
16     vector = evaluate_jacobian(f, seed_vector, center, eps)  
17  
18     length = perturbed_vector.length  
19     signal = [encode("add_to_jacobian"), acomm_id]  
20  
21     # perform 3-way handshake  
22     MPI_Send(signal, 2, int, acomm.head, acomm)  
23  
24     # broadcast jacobian column length  
25     MPI_Bcast(length, 1, int, acomm.head, acomm)  
26  
27     # broadcast jacobian column  
28     MPI_Bcast(vector.data, length, COO, acomm.all, acomm)
```

Listing 6.1: Pseudocode of the original ATHLET-NUT coupling; ATHLET part

## 6. Improvement of ATHLET-NUT Communication

---

```
1 # N - problem size
2 # J - allocated distributed jacobian matrix
3 # COO - compressed matrix coordinate format
4 nut_running = True
5
6 while nut_running:
7     if rank in heads:
8
9         # receive request
10        MPI_Recv(signal, 2, int, NUT_WORLD.any_client, NUT_WORLD)
11
12        comm = my_comm_list[signal[1]]
13        if (comm not None):
14            # posses resources
15            MPI_Bcast(signal, 2, int, comm.all, comm)
16        else:
17            continue
18
19    else:
20        MPI_Recv(signal, 2, int, NUT_WORLD.any_head, NUT_WORLD)
21
22    # decode request
23    comm = my_comm_list[signal[1]]
24    if (comm not None):
25        request = decode(signal[0])
26
27    case(request):
28        ...
29        if (request == "exit"):
30            # break while loop
31            nut_running = False
32
33        if (request == "add_to_jacobian"):
34            # receive jacobian column length
35            MPI_Recv(length, 1, int, comm.client, comm)
36
37            # receive row jacobian column
38            MPI_Recv(elements, length, COO, comm.client, comm)
39
40            for i in range(0, length):
41                if (local_min < elements[i].row < local_max):
42                    J.insert(elements[i])
43
44        ...
```

Listing 6.2: Pseudocode of the original ATHLET-NUT coupling: NUT part

## 6.2. Accumulator Concept

A simple concept, named *accumulator*, has been proposed to improve MPI communication during Jacobian matrix transfers preserving the current ATHLET-NUT architecture and coupling.

The concept represents two arrays of length  $2L$  where the first one, called *accumulator*, is used for accumulation of Jacobian matrix elements, stored in the compressed coordinate sparse matrix format, till the critical array length equaled to  $L = F \cdot N$ ; where  $N$  is the size of the underlying Jacobian matrix and  $F$  is a so-called capacity factor. Once the current array length of *accumulator* exceeds its critical length, the accumulated data are moved to *send buffer* by means of a simple swap of pointers, *ACC\_PTR* and *SEND\_BUFF\_PTR*, see Figure 6.4. Having swapped the pointers and reset control variables, the accumulation process can be immediately resumed together with an immediate instantiation of the corresponding non-blocking broadcast operation with respect to *send buffer* content.

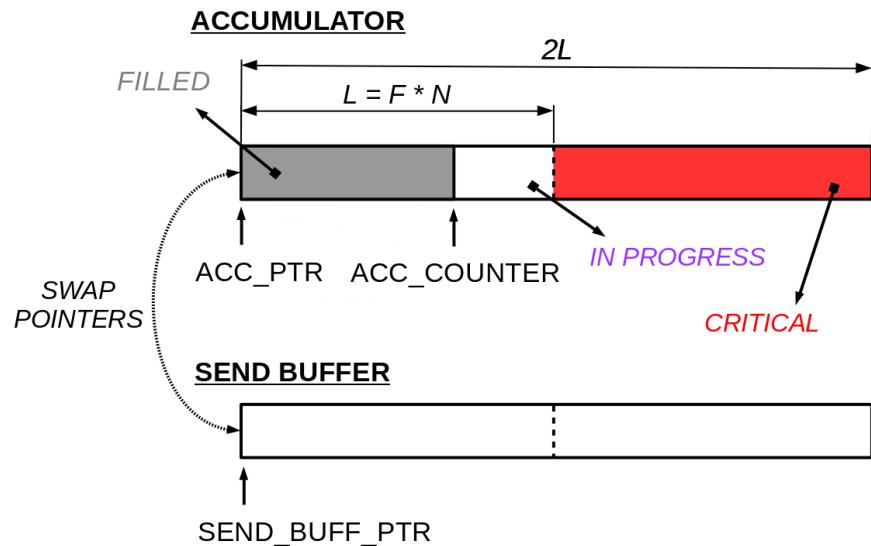


Figure 6.4.: *Accumulator* concept

The second array part of *accumulator*, also called the critical part, is used for safe placement of data surplus without any extra program checks and manipulations. Additionally, this event triggers a signal for a regular pointer swap and, therefore, the subsequent non-blocking data transfer.

The factor  $F$ , depicted in Figure 6.4, can be used by the user for two purposes. Mainly, it allows the user to adjust *send buffer* length  $L$  till the point of saturation of physical interconnection bandwidth, see Figure 6.5 as an example, and thus achieve efficient resource utilization. Additionally, it reduces an amount of handshakes, i.e. an amount of resource acquisition requests, between NUT and a client. The default value of the factor is equal to 1, however, we insistently recommend to increase the value via the corresponding environment variable for small-sized problems and operation of NUT in a multi-client mode.

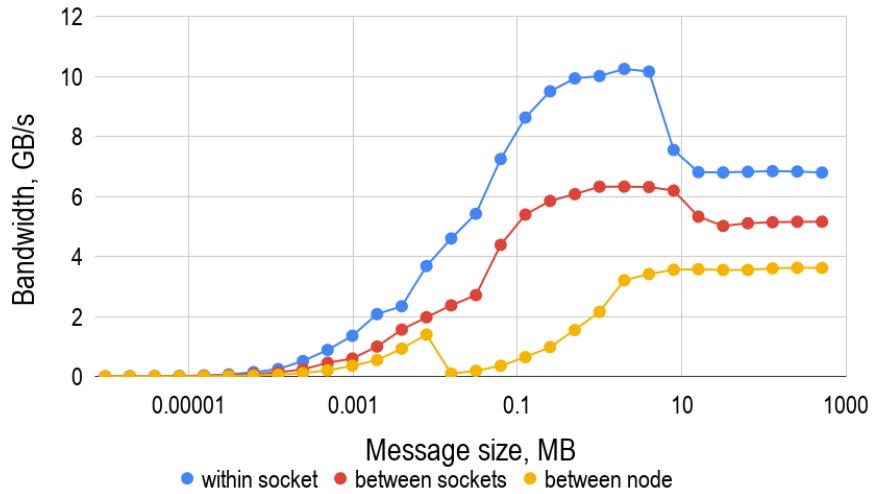


Figure 6.5.: Technical characteristics of HW1 hardware interconnection

Figure 6.6 depicts an application of the *accumulator* algorithm to the example represented in Figure 6.3 with the following parameters:  $N = 100$  and  $F = 1$ . It can be clearly observed the algorithm reduces the number of transfers from 28 to 12. Additionally, the average column length, excluding the last one, jumps from 56 to 131. By and large, the algorithm allows to transform an original distribution shape to a more or less rectangular one which, in turn, allows to transfer a matrix in approximately equal chunks.

Before ATHLET can send a request to NUT to start solving Systems 3.1 it has to be certain that the entire Jacobian matrix has been transferred to the NUT side. For this reason, the last column transfer is done by means of the corresponding blocking MPI operation. It means ATHLET gets blocked only during the last column transfer and MPI gives the execution control back only when the last piece of data has been successfully distributed among NUT processes.

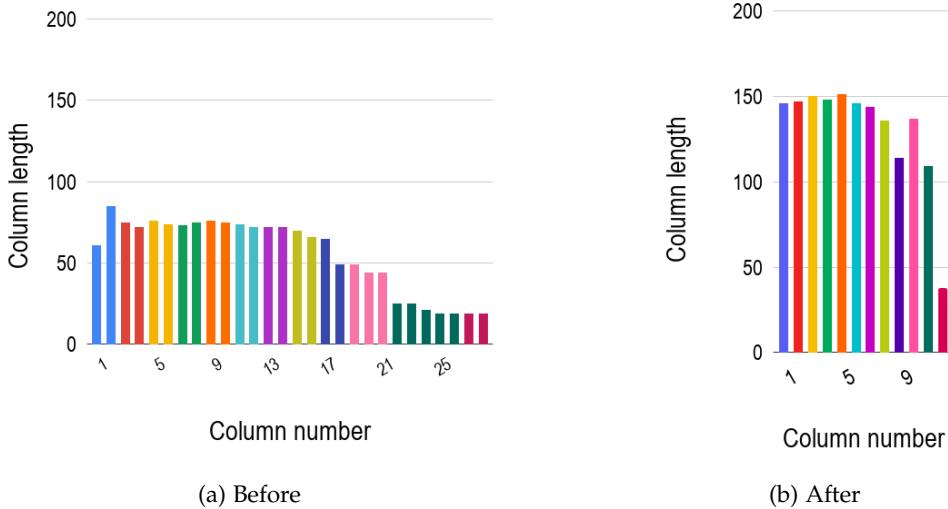


Figure 6.6.: An application of the *accumulator* concept to the example depicted in Figure 6.3, with  $N = 100$  and  $F = 1$

### 6.3. Benchmark and Test Data

ATHLET is a dedicated industrial CFD package meant for simulation of thermal-hydraulic circuits in various nuclear power plant facilities. Besides the main part, i.e. the solver, it includes some pre-processing steps that allow the user to conveniently set up different simulation parameters, computational mesh, output data, etc.

Testing of new concepts and ideas directly in ATHLET can be quite cumbersome, computationally expensive and inconvenient. Therefore, a dedicated benchmark has been developed to test the *accumulator* concept.

The benchmark fully replicates all basic ideas of the original ATHLET implementation and the new data transfer concept. It focuses only on compressed Jacobian matrix transfers and, therefore, does not include any compute-expensive operations such as non-linear function perturbations with seed vectors. The approach allowed to sufficiently speed up time of development, comparison and testing which, in turn, helped in designing the final concept described in Section 6.2.

In order to mimic the real run-time ATHLET-NUT behavior during Jacobian matrix updates, a few communication patterns were recorded in ATHLET and played in the benchmark. The recordings helped to generate column vectors with the lengths corresponding to that in the recordings, filled with random numbers. Figure 6.7 shows an example of a part of *cube-64* communication pattern used in the study, where COO stands for compressed coordinate format. As it can be observed, the pattern includes both full and partial Jacobian updates, described in Section 2.1.

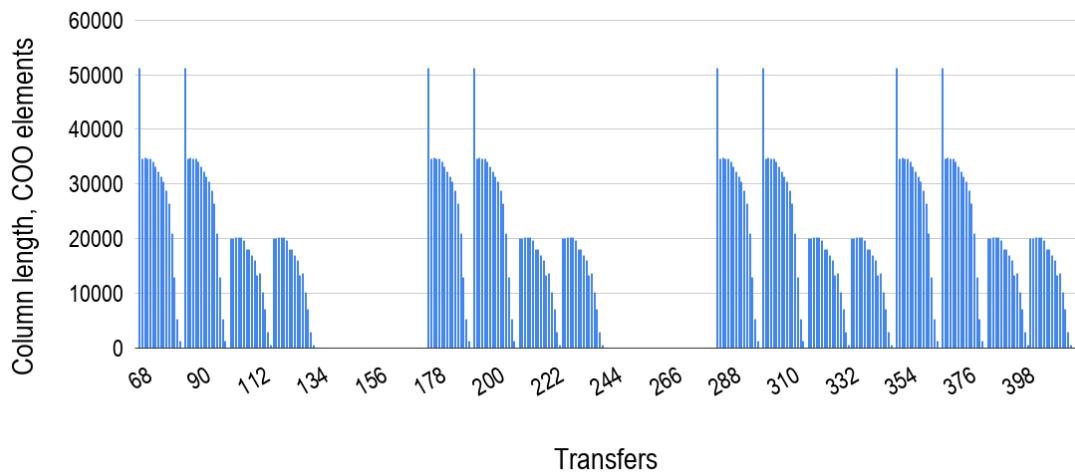


Figure 6.7.: A part of *cube-64* communication pattern

According to the *accumulator* concept, described in Section 6.2, the main changes take place only on the client side and hence the server side remains unchanged which follows the original idea of the least code modifications. Code Listing 6.3 represents an additional auxiliary class used for data accumulation. Pseudocode of the benchmark client side is in Listing 6.4.

## 6. Improvement of ATHLET-NUT Communication

---

```
1 # problem_size – given Jacobian matrix size
2 # COO – compressed matrix coordinate format
3 class Accumulator:
4     constructor(problem_size, acomm, acomm_id):
5         private:
6             N = problem_size; comm = acomm; id = acomm_id
7             signal = [encode("add_to_jacobian"), id]
8             is_allocated = false; is_non_blocking_op_called = false
9             send_buffer = []
10            factor = int(read_enviroment_variable("CNUT_ACC_SIZE"))
11            if factor == None:
12                factor = 1
13                permissible_size = factor * N
14            public:
15                accumulator = []
16
17            def allocate_accumulator():
18                if is_allocated == false:
19                    accumulator = allocate(2 * permissible_size, type(COO))
20                    send_buffer = allocate(2 * permissible_size, type(COO))
21                    is_allocated = true
22
23            def deallocate_accumulator():
24                if is_allocated == true:
25                    deallocate(accumulator); deallocate(send_buffer)
26                    is_allocated = false
27
28            def commit():
29                if accumulator.size > permissible_size:
30                    swap(accumulator.pointer, send_buffer.pointer)
31                    if is_non_blocking_op_called == true:
32                        MPI_Wait()
33                    # perform 3-way handshake
34                    MPI_Send(signal, 2, int, comm.head, comm)
35                    # send data
36                    MPI_Ibcast(send_buffer.size, 1, int, comm.head, comm)
37                    MPI_Ibcast(send_buffer.data, send_buffer.size, COO, comm.all, comm)
38                    is_non_blocking_op_called = true
39                    accumulator.content.reset("to_beginning")
40
41            def finalize():
42                if is_non_blocking_op_called == true:
43                    MPI_Wait()
44                    MPI_Send(signal, 2, int, comm.head, comm)
45                    MPI_Bcast(accumulator.size, 1, int, comm.head, comm)
46                    MPI_Bcast(accumulator.data, accumulator.size, COO, comm.all, comm)
47                    is_non_blocking_op_called = false
48                    accumulator.content.reset("to_beginning")
```

Listing 6.3: Pseudocode of an auxiliary *Accumulator* class

```

1 # GIVEN PARAMETERS:
2 # acomm – the athlet communicator
3 # acomm_id – athlet identification number
4 # N – problem size
5 # recording – data structure that holds a recorded communication pattern
6 # COO – compressed matrix coordinate format
7
8 if global_counter == 0:
9     container = Accumulator.constructor(N, acomm, acomm_id)
10    container.allocate_accumulator()
11    ++global_counter
12    file = open("benchmark_results.txt", "w")
13
14 for column in recording:
15
16     time_start = MPI_Wtime()
17
18     # charge accumulator
19     for i in range(column.length):
20         element = generate_random_coo_element()
21         container.accumulator.add(element)
22
23     # instantiate non-blocking data broadcast
24
25     container.commit()
26     time_end = MPI_Wtime() - time_start
27     file.write(column.length, time_end)
28
29 # transfer the remainder and synchronize
30 time_start = MPI_Wtime()
31     container.finalize()
32 time_end = MPI_Wtime() - time_start
33 file.write(column.length, time_end)

```

Listing 6.4: Pseudocode of a modified client side of the benchmark

## 6.4. Results

The benchmark was ran on HW1 compute-cluster where the client and server parts were distributed in three different ways, namely: within a socket, in two separate sockets of a node and in two separate nodes. Nodes of HW1 cluster are connected via an *infiniband* interconnect with the characteristics shown in Figure 6.5. In order to estimate an effect of pure data accumulation, the benchmark, Listing 6.3, was modified to use only blocking MPI operations i.e. MPI\_Bcast. We denote the main benchmark as *BM1* and the modified one as *BM2* to distinguish and separately explain effects of

non-blocking data transfers and pure data accumulation.

Figure 6.8 represents results of the benchmarks obtained using *cube-64* communication pattern. The client and server parts of the code were distributed in different sockets within the same node. Factor  $F$  was equal to 1.

Figure 6.8a shows that *accumulator* approach results in more than **6** times drop, from **344** to **51**, of the total number of data transfers and resulting resource acquisitions, within the range depicted on the graphs. According to *BM2* benchmark, the accumulation effect reduces the run-time by almost **9%** by means of more efficient utilization of intra-node interconnection. The obtained results also demonstrate that overall accumulative effect of both accumulation and non-blocking data transfers reduces the run-time of *BM1* benchmark in more than **26%**. Table 6.1 summarizes results obtained for all three client-server distributions within the same range of the recorded communication pattern displayed in Figure 6.7.

Benchmark name	BM2, %	BM1, %
within a socket	7.61	13.84
between sockets	9.04	26.26
between nodes	-2.06	3.20

Table 6.1.: Time reduction of data transfers with respect to the original implementation in case of execution of *cube-64* communication pattern

It turns out that *BM2* benchmark is slower than the original ATHLET approach in approximately **2%** in case of inter-node communication. However, non-blocking data broadcasts, according to *BM1* benchmark, help to alleviate the slow-down and achieve almost **3%** of improvement.

Unimpressive results of non-blocking inter-node communication can be explained by specifics of the benchmark design. In particular, time spent on generation of random matrix elements was not enough to overlap time spent on non-blocking data transfers in case of *cube-64* test-case, see Figure 6.9. Thus, the execution control was probably suspended by MPI library at each subsequent call of *MPI\_Wait()* function.

The slow-down resulting from pure data accumulation could be explained by automatic MPI protocol switching, namely: *Eager* and *Rendezvous* [34]. The protocols are

## 6. Improvement of ATHLET-NUT Communication

---

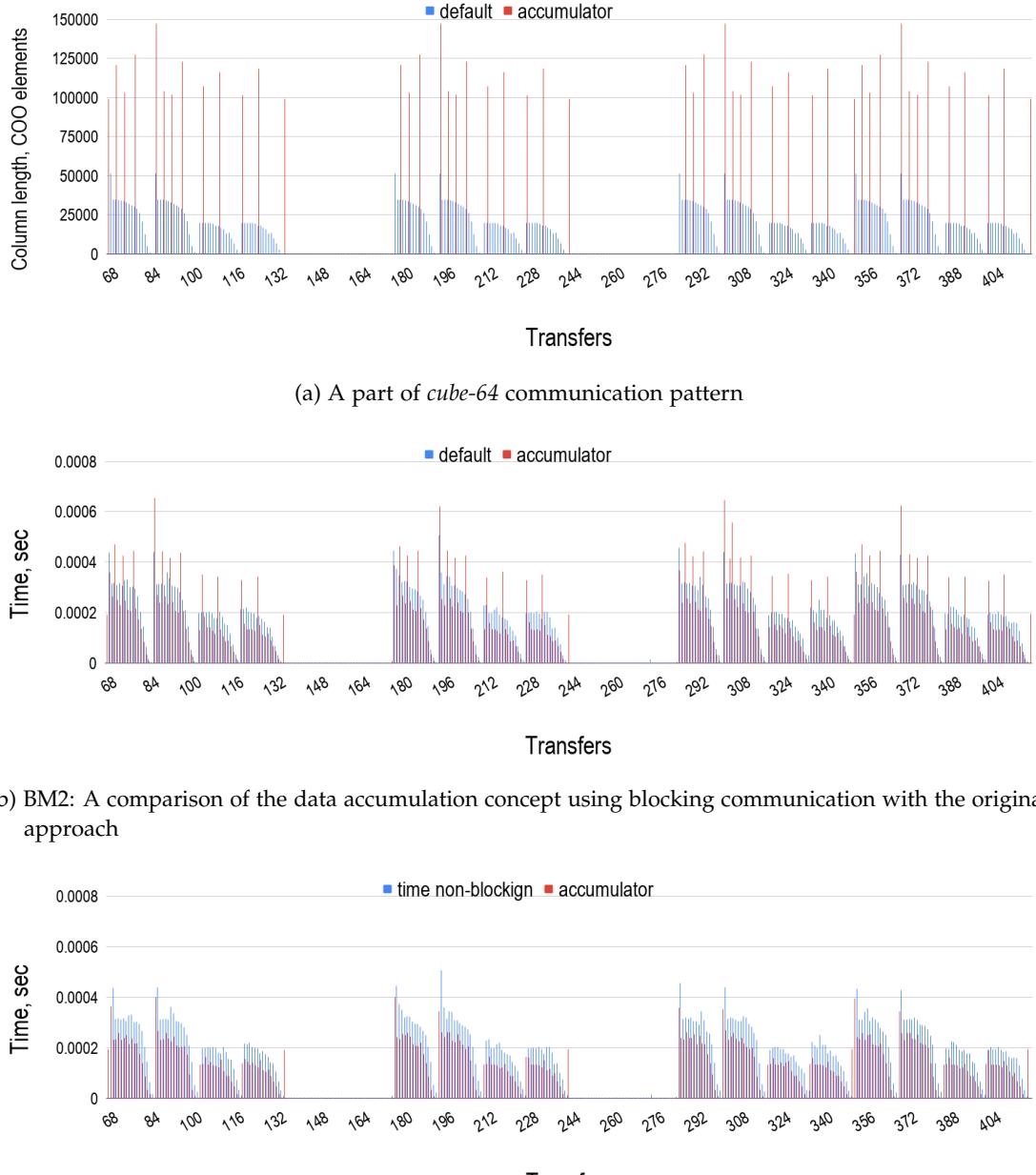


Figure 6.8.: Comparisons of the benchmarks running a recorded part of *cube-64* communication pattern between two sockets of a node

## 6. Improvement of ATHLET-NUT Communication

---

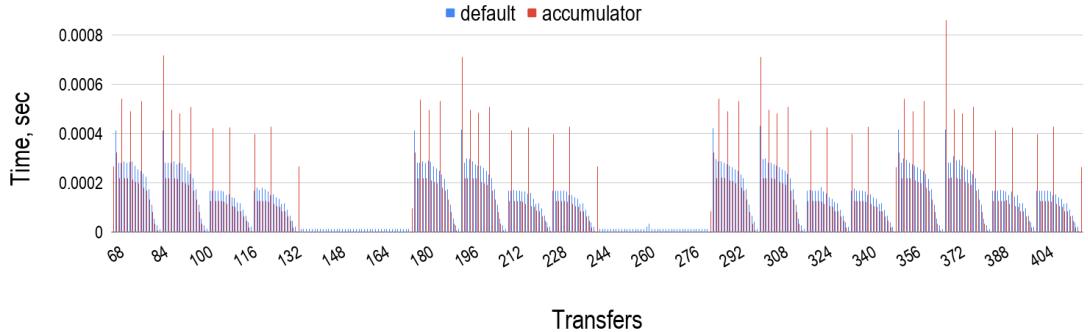


Figure 6.9.: A comparison of *BM1* benchmark with the original ATHLET-NUT implementation running a recorded part of *cube-64* communication pattern between two compute-nodes

dedicated to small and large message transfers, respectively, where a quantitative measure of the message size is defined by a concrete implementation of the MPI standard, however, it can be controlled through dedicated environment variables.

Similar results were observed for *cube-645* test case where the number of equations was approximately  $10^6$  and the average compressed Jacobian column length reached around  $1.7 \cdot 10^5$  elements. In case of inter-node communication, *BM1* benchmark again showed performance degradation by **6.35%** whereas non-blocking data broadcasts improved run-time by **23.21%**. Such performance jump, from **-6.35%** to **23.21%**, can be explained by the fact that time spent on generation of random elements was enough to hide the corresponding data transfers and overheads.

Ideas, expressed in *BM1* benchmark, Listings 6.3 and 6.4, were successfully implemented in NUT, namely: in the client side of NUT located in ATHLET. Several simulation scenarios were taken for the final verification and performance testing, namely: *cube-64*, *k3-2* and *pwr-3d*. Verification of the modified code did not detect any deviations of numerical results from the original implementation. Additionally, all tests showed considerable improvements in communication time. As an example, time spent on communication between ATHLET and NUT during compressed Jacobian transfers decreased by **66.17%**, **76.03%** and **42.55%** for intra-socket, intra-node and inter-node client-server process allocations, respectively, for *pwr-3d* scenario, taking it as the most representative simulation test-scenario known in GRS. However, the overall improvement of applied changes achieved only **0.14%** on average, regardless of a client-server allocation. Profiling showed the communication part of the original

implementation took around **0.24%** of the total time spent on matrix evaluations and transfers. This fact explains this negligible overall performance gain, resulted from the source code modification, that was observed in all conducted final tests.

## 6.5. Conclusion

In this part of the study, we have designed and implemented the *accumulator* concept for efficient transfers of sparse compressed Jacobian matrices between ATHLET and NUT. The concept is rather simple and did not require drastic changes of the existing software design and architecture. In spite of simplicity, the concept allows to significantly reduce communication time i.e. by almost **60%**. The overall performance gain comes from three different sources:

1. efficient utilization of interconnection
2. reduced number of handshakes and, as a result, a reduced amount of NUT process synchronizations
3. overlaps of communications with computations

The study has shown that non-blocking data transfers are the main source of the performance gain. Efficient bandwidth utilization can additionally give **7-9%** of improvement when applications work within the same compute-node.

One can experience slight slow-down from pure data accumulation in case of inter-node communication due to probable MPI protocol switching. However, as it has been shown, it is always compensated by means of communication/computation overlaps.

The final tests have shown the concept does not give a considerable overall improvement because the computation part takes almost **99.8%** of the total execution time of the corresponding part of the source code. However, results may be much better in case of multi-client operation of NUT, especially when clients are sharing common NUT processes; a reduced number of data transfers results in a reduced amount of handshakes which are always accompanied by the resource acquisition mechanism, described in Section 2.3. Unfortunately, it is difficult to design and prepare a set of valid tests to verify this statement.

By and large, verification of the modified code has not detected any deviations in numerical results. The new concept has always resulted in a slight overall performance

---

*6. Improvement of ATHLET-NUT Communication*

---

gain. The study has also shown the main bottleneck is, indeed, the non-linear function evaluation.

It is worth mentioning that only the sequential ATHLET code, capable to run only in a single core, was available for this study. However, there exists a parallel version of ATHLET multi-threaded with OpenMP. Therefore, the results can be even better because of a reduced fraction of execution time spent on non-linear function evaluations. This fact also shows that performance tuning of ATHLET is constantly in progress and is being done in parallel among several departments of GRS, covering different areas of the program source code.

# **Appendices**

## A. Sparsity Patterns of Matrix Sets

#### *A. Sparsity Patterns of Matrix Sets*

---

*Sparsity patterns of the GRS matrix set are not available for any online publication. Please, contact GRS representatives to get access to the data.*

*A. Sparsity Patterns of Matrix Sets*

---

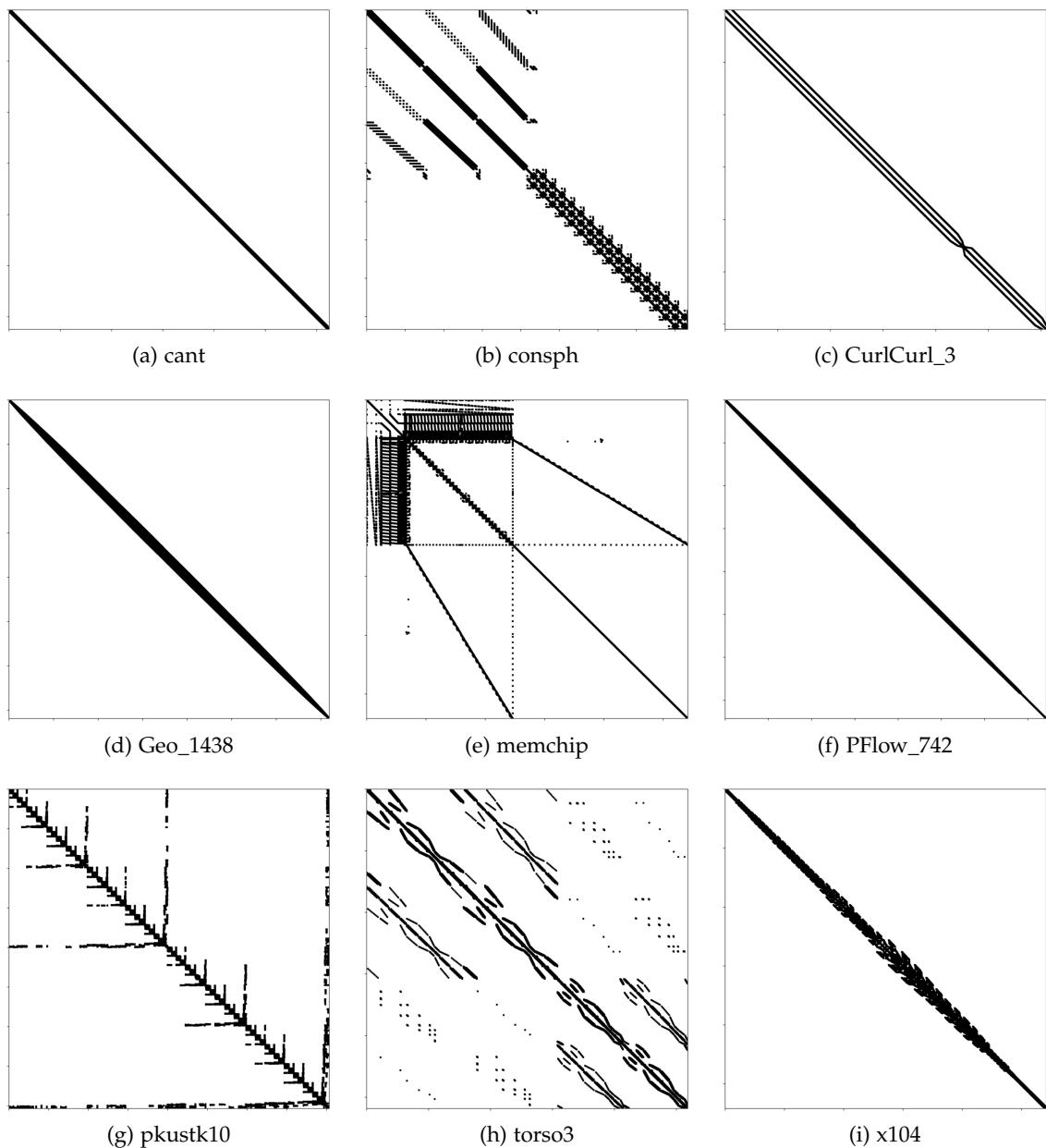


Figure A.1.: Sparsity patterns of SuiteSparse matrix set

## **B. Selection of a Sparse Direct Linear Solver**

MPI	MUMPS	PaStiX	SuperLU
1	4.58E-02	5.60E-02	4.64E+00
2	4.31E-02	5.14E-02	1.89E+00
3	4.51E-02	5.28E-02	1.22E+00
4	4.61E-02	5.64E-02	9.13E-01
5	4.92E-02	5.97E-02	7.70E-01
6	5.37E-02	6.14E-02	6.04E-01
7	5.42E-02	6.51E-02	crashed
8	5.41E-02	6.60E-02	4.81E-01
9	5.69E-02	6.84E-02	4.35E-01
10	5.86E-02	7.22E-02	4.08E-01

MPI	MUMPS	PaStiX	SuperLU
11	5.93E-02	8.97E-02	crashed
12	6.07E-02	9.20E-02	3.61E-01
13	6.26E-02	8.25E-02	crashed
14	6.28E-02	9.75E-02	crashed
15	6.43E-02	1.03E-01	3.05E-01
16	6.55E-02	1.05E-01	2.99E-01
17	6.61E-02	9.46E-02	crashed
18	6.73E-02	1.24E-01	2.65E-01
19	6.84E-02	1.14E-01	crashed
20	7.02E-02	1.32E-01	2.60E-01

Table B.1.: Comparisons of parallel performance of *pwr-3d* matrix factorizations using MUMPS, PaStiX and SuperLU\_DIST libraries with their default parameter settings

MPI	MUMPS	PaStiX	SuperLU
1	1.55E+02	6.44E+01	time-out
2	6.28E+01	4.84E+01	time-out
3	5.06E+01	5.02E+01	time-out
4	4.17E+01	4.50E+01	time-out
5	2.52E+01	3.98E+01	time-out
6	2.58E+01	4.29E+01	time-out
7	2.65E+01	4.30E+01	time-out
8	2.59E+01	3.73E+01	time-out
9	1.95E+01	4.08E+01	time-out
10	1.91E+01	3.81E+01	time-out

MPI	MUMPS	PaStiX	SuperLU
11	1.77E+01	3.75E+01	time-out
12	1.60E+01	3.58E+01	time-out
13	1.42E+01	3.59E+01	time-out
14	1.45E+01	3.57E+01	time-out
15	1.47E+01	3.52E+01	time-out
16	1.41E+01	3.45E+01	time-out
17	1.54E+01	3.31E+01	time-out
18	1.52E+01	3.31E+01	time-out
19	1.52E+01	3.16E+01	time-out
20	1.38E+01	3.15E+01	time-out

Table B.2.: Comparisons of parallel performance of *k3-2* matrix factorizations using MUMPS, PaStiX and SuperLU\_DIST libraries with their default parameter settings

MPI	MUMPS	PaStiX	SuperLU
1	1.52E+01	1.61E+01	crashed
2	1.13E+01	1.13E+01	crashed
3	1.00E+01	1.03E+01	crashed
4	9.29E+00	1.05E+01	crashed
5	8.85E+00	9.84E+00	crashed
6	8.43E+00	8.99E+00	crashed
7	8.64E+00	9.69E+00	crashed
8	8.70E+00	9.12E+00	crashed
9	8.91E+00	8.94E+00	crashed
10	8.76E+00	9.26E+00	crashed

MPI	MUMPS	PaStiX	SuperLU
11	8.62E+00	9.09E+00	crashed
12	8.53E+00	8.92E+00	crashed
13	8.44E+00	9.13E+00	crashed
14	8.52E+00	9.00E+00	crashed
15	8.54E+00	9.19E+00	crashed
16	8.56E+00	9.05E+00	crashed
17	8.65E+00	9.12E+00	crashed
18	8.62E+00	8.96E+00	crashed
19	8.66E+00	9.30E+00	crashed
20	8.66E+00	9.16E+00	crashed

Table B.3.: Comparisons of parallel performance of *cube-645* matrix factorizations using MUMPS, PaStiX and SuperLU\_DIST libraries with their default parameter settings

## C. Fill Reducing Reorderings

### C. Fill Reducing Reorderings

---

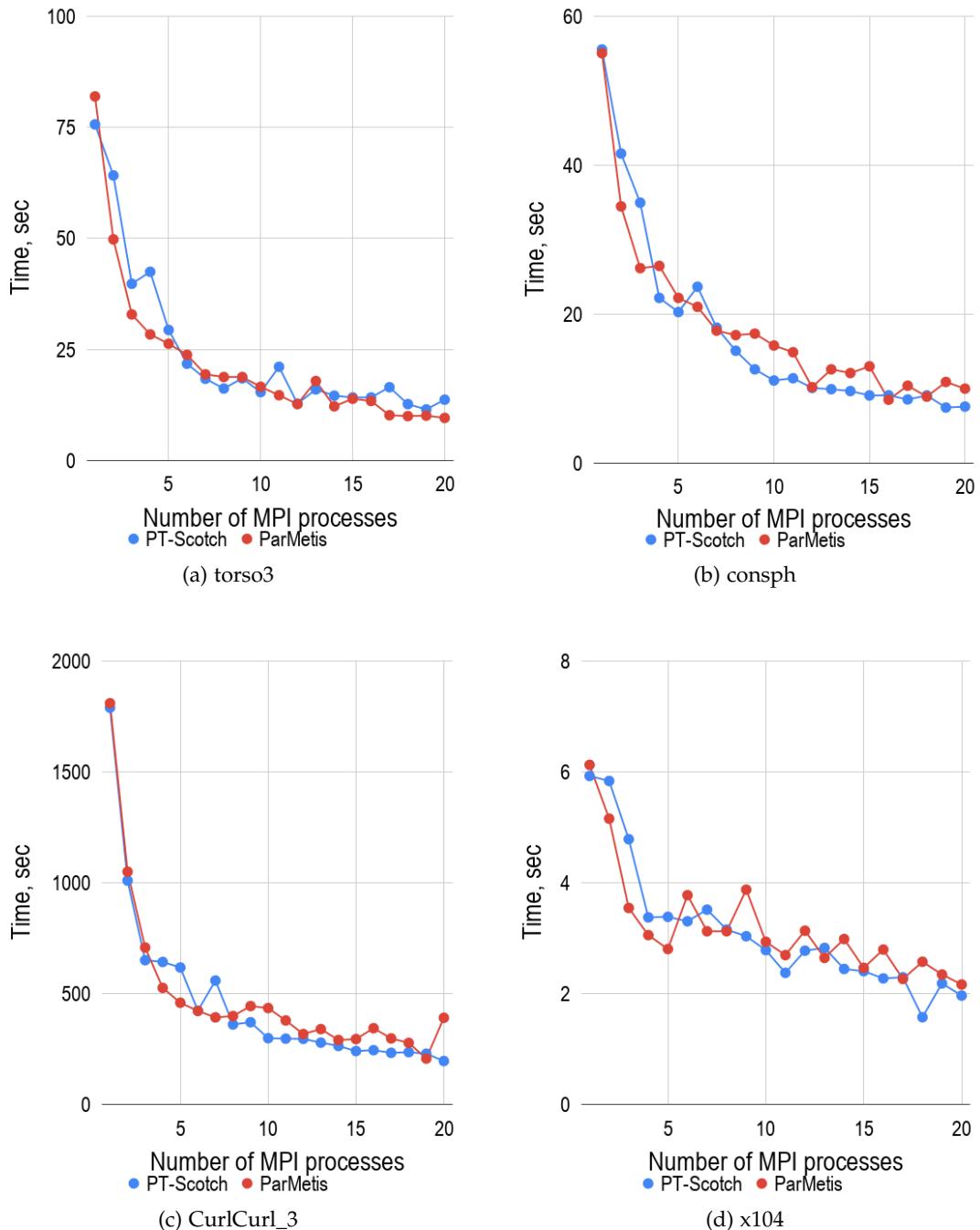


Figure C.1.: An influence of different fill reducing algorithms on parallel factorizations of *torso3*, *consph*, *CurlCurl\_3* and *x104* matrices

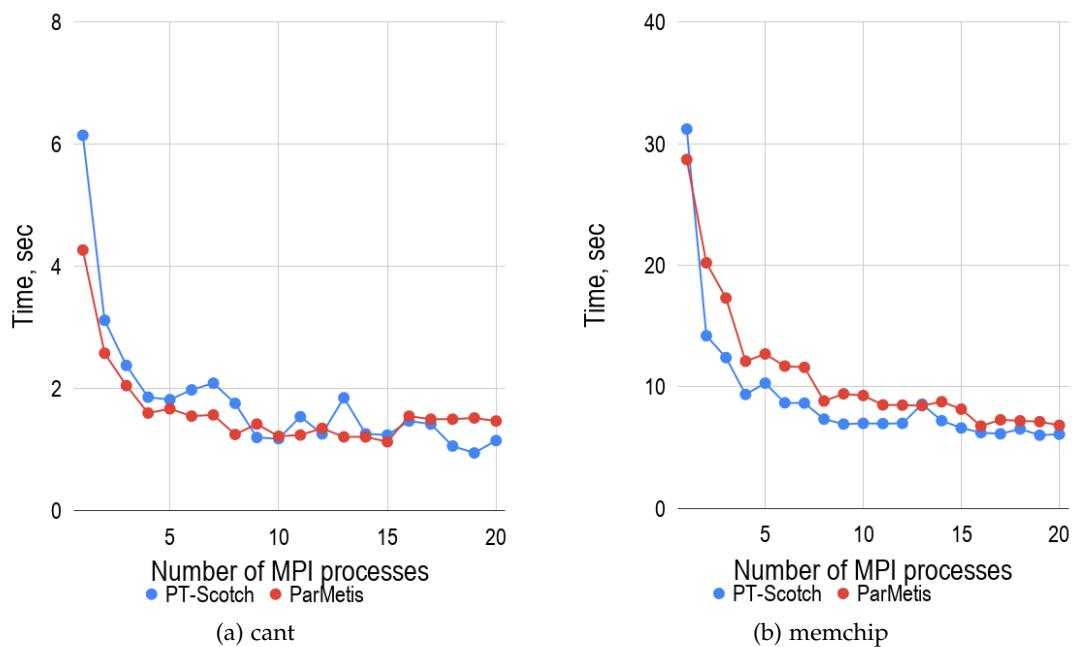


Figure C.2.: An influence of different fill reducing algorithms on parallel factorizations of *cant* and *memchip* matrices

## D. MPI Process Pinning

#### D. MPI Process Pinning

---

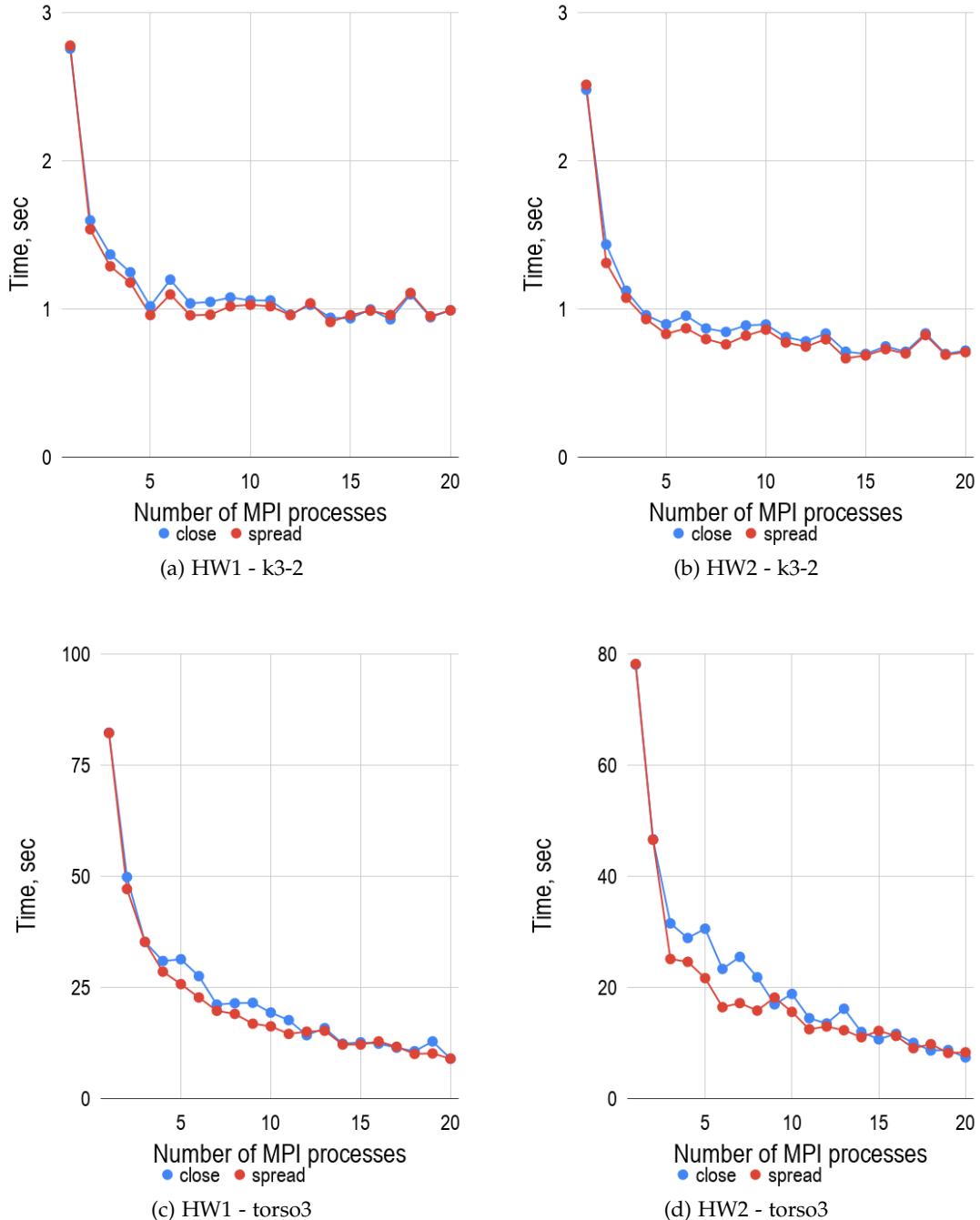
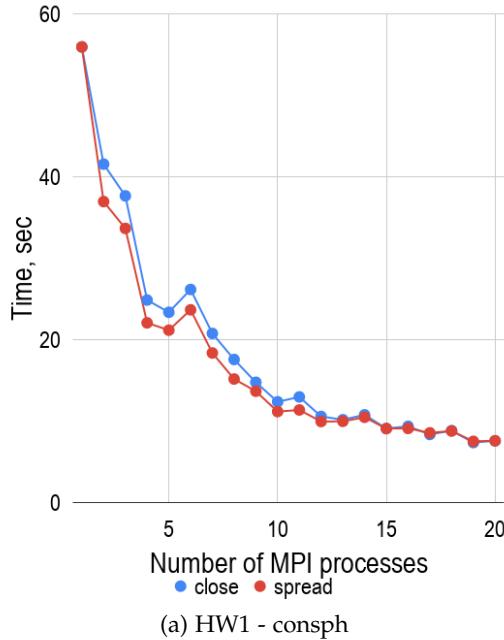


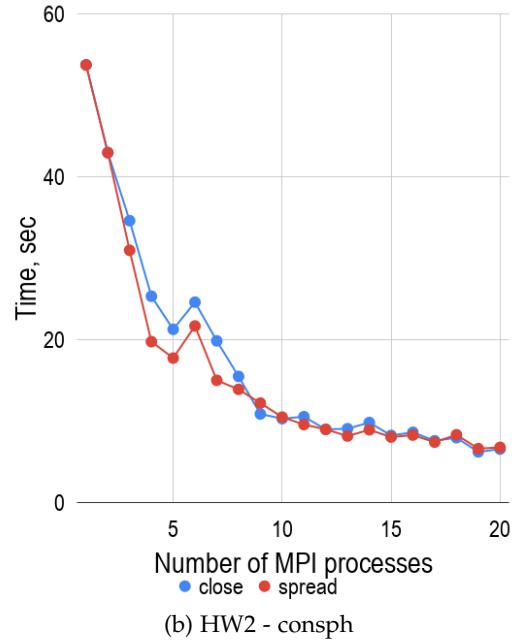
Figure D.1.: Comparisons of *close* and *spread* pinning strategies applied to parallel factorizations of *cube-64* and *torso3* matrices

#### D. MPI Process Pinning

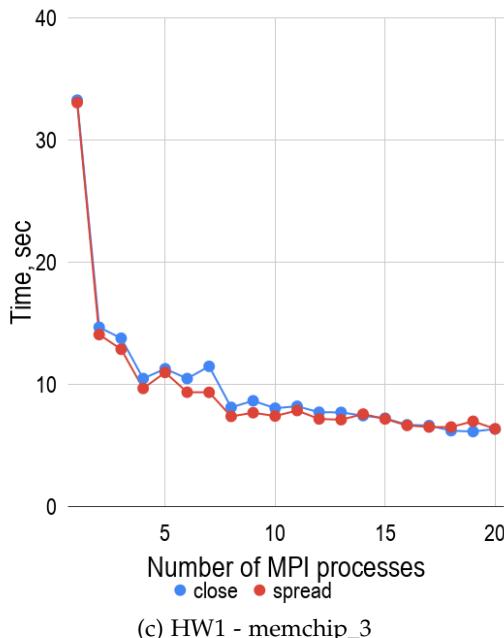
---



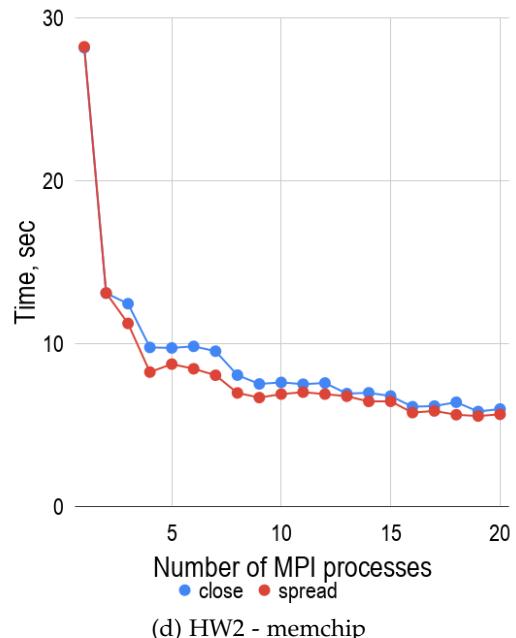
(a) HW1 - consph



(b) HW2 - consph



(c) HW1 - memchip\_3



(d) HW2 - memchip

Figure D.2.: Comparisons of *close* and *spread* pinning strategies applied to parallel factorizations of *consph* and *memchip* matrices

## **E. Optimized BLAS Libraries**

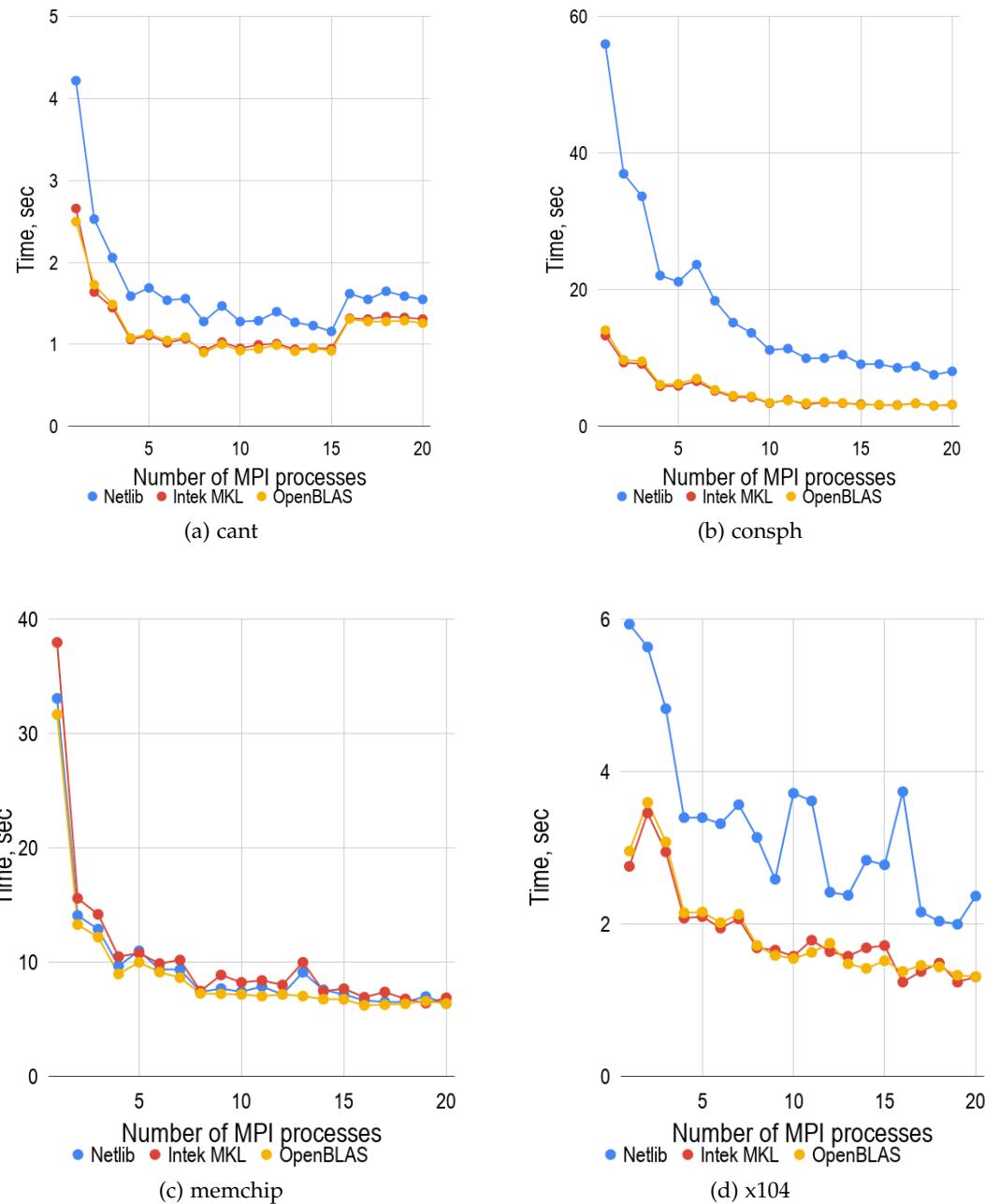


Figure E.1.: Comparisons of parallel factorizations of *cant*, *consph*, *memchip* and *x104* matrices performed on HW1 machine using MUMPS solver linked to different BLAS implementations

### E. Optimized BLAS Libraries

---

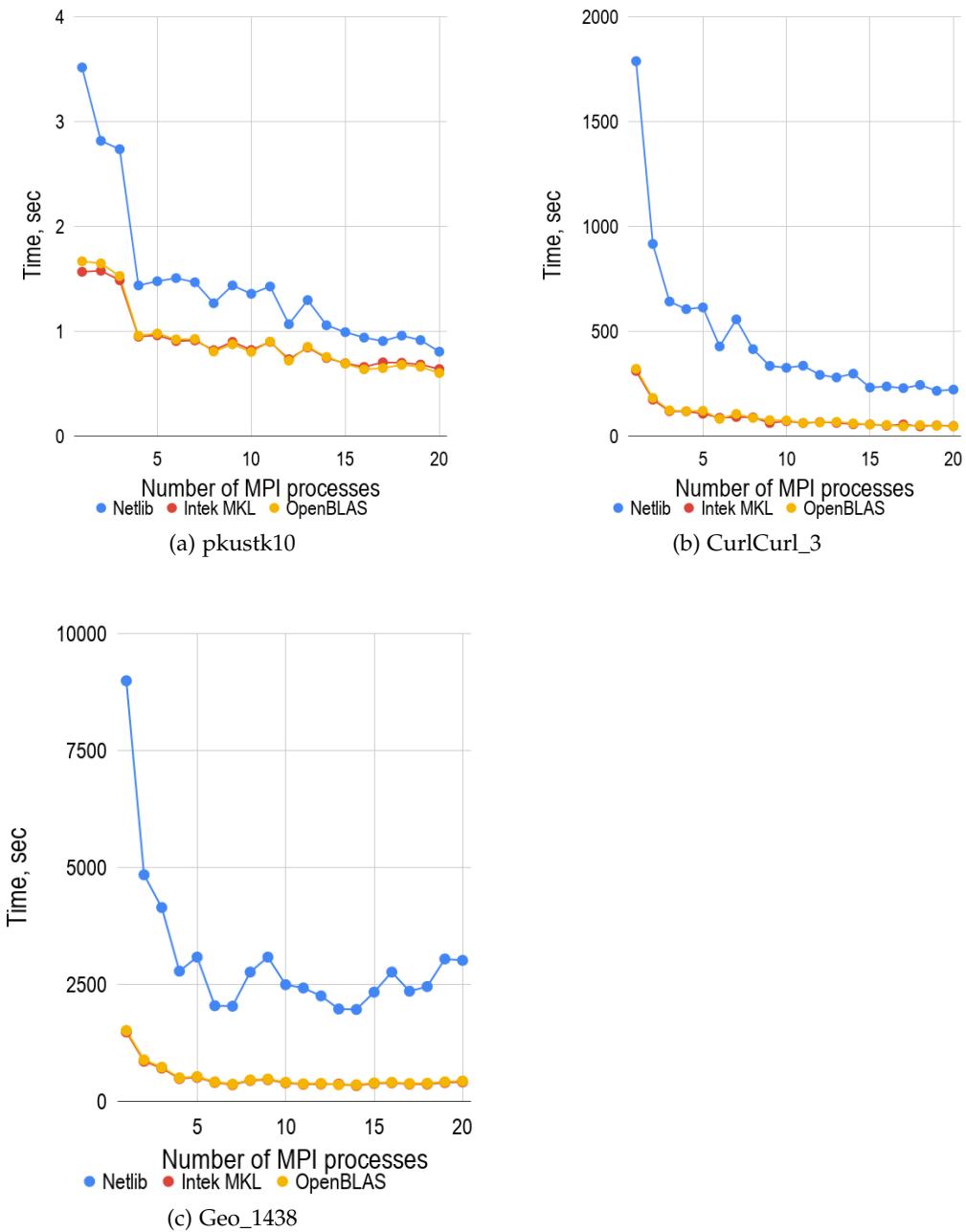


Figure E.2.: Comparisons of parallel factorizations of *pkustk10*, *CurlCurl\_3* and *Geo\_1438* matrices performed on HW1 machine using MUMPS solver linked to different BLAS implementations

# Bibliography

- [1] P. R. Amestoy, I. S. Duff, J.-Y. L'Excellent, and P. Plecháč. "PARASOL: An Integrated Programming environment for Parallel Sparse Matrix Solvers". In: *PINEAPL Workshop, A Workshop on the Use of Parallel Numerical Libraries in Industrial End-user Applications*. CERFACS, Toulouse, France, 1998.
- [2] P. R. Amestoy, T. A. Davis, and I. S. Duff. "An approximate minimum degree ordering algorithm". In: *SIAM Journal on Matrix Analysis and Applications* 17.4 (1996), pp. 886–905.
- [3] P. R. Amestoy, I. S. Duff, J. Koster, and J.-Y. L'Excellent. "MUMPS: A multifrontal massively parallel solver". In: *ERCIM News* 50 (2001), pp. 14–15.
- [4] P. Amestoy. "Recent progress in parallel multifrontal solvers for unsymmetric sparse matrices". In: *Proceedings of the 15th World Congress on Scientific Computation, Modelling and Applied Mathematics, IMACS*. Vol. 97. 1997.
- [5] P. Arbenz. *Lecture notes of Numerical Methods for Solving Large Scale Eigenvalue Problems: Arnoldi and Lanczos algorithms*. 2018. URL: <http://people.inf.ethz.ch/arbenz/ewp/Lnotes/chapter10.pdf>.
- [6] M. Arioli, J. W. Demmel, and I. S. Duff. "Solving sparse linear systems with sparse backward error". In: *SIAM Journal on Matrix Analysis and Applications* 10.2 (1989), pp. 165–190.
- [7] H. Austregesilo, C. Bals, A. Hora, G. Lerchl, P. Romstedt, P. Schöffel, D. Von der Cron, and F. Weyermann. *ATHLET Mod 3.1A – Models and Methods*. distributed with ATHLET. Mar. 2016.
- [8] S. Balay, S. Abhyankar, M. Adams, J. Brown, P. Brune, K. Buschelman, L. Dalcin, A. Dener, V. Eijkhout, W. Gropp, D. Kaushik, M. Knepley, D. May, L. C. McInnes, R. T. Mills, T. Munson, K. Rupp, P. Sanan, B. Smith, S. Zampini, H. Zhang, and H. Zhang. *PETSc Web page*. <http://www.mcs.anl.gov/petsc>. 2018.
- [9] S. Balay, S. Abhyankar, M. Adams, J. Brown, P. Brune, K. Buschelman, L. Dalcin, A. Dener, V. Eijkhout, W. Gropp, et al. *PETSc Users Manual: Revision 3.10*. Tech. rep. Argonne National Lab.(ANL), Argonne, IL (United States), 2018.

## Bibliography

---

- [10] I. Chowdhury and J.-Y. L'Excellent. "Some experiments and issues to exploit multicore parallelism in a distributed-memory parallel sparse direct solver". PhD thesis. INRIA, 2010.
- [11] T. A. Davis and Y. Hu. "The University of Florida Sparse Matrix Collection". In: *ACM Trans. Math. Softw.* 38.1 (Dec. 2011), 1:1–1:25. issn: 0098-3500. doi: 10.1145/2049662.2049663.
- [12] I. S. Duff, R. G. Grimes, and J. G. Lewis. "Sparse Matrix Test Problems". In: *ACM Trans. Math. Softw.* 15.1 (Mar. 1989), pp. 1–14. issn: 0098-3500. doi: 10.1145/62038.62043.
- [13] I. S. Duff and J. Koster. "On algorithms for permuting large entries to the diagonal of a sparse matrix". In: *SIAM Journal on Matrix Analysis and Applications* 22.4 (2001), pp. 973–996.
- [14] I. S. Duff and J. Koster. "The design and use of algorithms for permuting large entries to the diagonal of sparse matrices". In: *SIAM Journal on Matrix Analysis and Applications* 20.4 (1999), pp. 889–901.
- [15] I. S. Duff and S. Pralet. "Strategies for scaling and pivoting for sparse symmetric indefinite problems". In: *SIAM Journal on Matrix Analysis and Applications* 27.2 (2005), pp. 313–340.
- [16] I. S. Duff and J. K. Reid. "The multifrontal solution of indefinite sparse symmetric linear". In: *ACM Transactions on Mathematical Software (TOMS)* 9.3 (1983), pp. 302–325.
- [17] European Commission. *Nuclear Energy: Safe nuclear power*. 2018. URL: <https://ec.europa.eu/energy/en/topics/nuclear-energy>.
- [18] R. Falgout, A. Cleary, J. Jones, E. Chow, V. Henson, C. Baldwin, P. Brown, P. Vassilevski, and U. Yang. "Hypre User's Manual, Version 2.7". In: *Center for Applied Scientific Computing (CASC), Lawrence Livermore National Laboratory, Livermore, CA* (2010).
- [19] A. H. Gebremedhin, F. Manne, and A. Pothen. "What color is your Jacobian? Graph coloring for computing derivatives". In: *SIAM review* 47.4 (2005), pp. 629–705.
- [20] G. Geist and E. Ng. "Task scheduling for parallel sparse Cholesky factorization". In: *International Journal of Parallel Programming* 18.4 (1989), pp. 291–314.
- [21] Gesellschaft für Anlagen und Reaktorsicherheit (GRS) gGmbH. *ATHLET 3.1A: Program Overview*. 2016. URL: <https://software.intel.com/en-us/mkl-developer-reference-c-overview-of-scalapack-routines>.

## Bibliography

---

- [22] Gesellschaft für Anlagenund Reaktorsicherheit (GRS) gGmbH. *Development and Validation of Simulation Codes*. URL: <https://www.grs.de/en/content/development-and-validation-simulation-codes>.
- [23] Gesellschaft für Anlagenund Reaktorsicherheit (GRS) gGmbH. *GRS – Safety for man and the environment*. URL: <https://www.grs.de/en/company>.
- [24] A. Guermouche, J.-Y. L'Excellent, and G. Utard. "On the memory usage of a parallel multifrontal solver". In: *Parallel and Distributed Processing Symposium, 2003. Proceedings. International*. IEEE. 2003, 8-pp.
- [25] A. Gupta, S. Koric, and T. George. "Sparse matrix factorization on massively parallel computers". In: *Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*. ACM. 2009, p. 1.
- [26] *Intel Math Kernel Library Documentation, Overview of ScalAPACK Routines*. 2017. URL: <https://software.intel.com/en-us/mkl-developer-reference-c-overview-of-scalapack-routines>.
- [27] B. M. Irons. "A frontal solution program for finite element analysis". In: *International Journal for Numerical Methods in Engineering* 2.1 (1970), pp. 5–32.
- [28] Jordan Hanania and Braden Heffernan and Jenden, James and Lefsrud, Nathan and Lloyd, Ellen and Stenhouse, Kailyn and Toor, Jasdeep and Donev, Jason. *Nuclear power plant - Energy Education*. 2019. URL: [https://energyeducation.ca/encyclopedia/Nuclear\\_power\\_plant](https://energyeducation.ca/encyclopedia/Nuclear_power_plant).
- [29] G. Karypis and V. Kumar. *MeTis: Unstructured Graph Partitioning and Sparse Matrix Ordering System, Version 4.0*. <http://www.cs.umn.edu/~metis>. University of Minnesota, Minneapolis, MN, 2009.
- [30] *Krylov subspace projection methods: Rayleigh-Ritz procedure*. 2009. URL: <http://web.cs.ucdavis.edu/~bai/Winter09/krylov.pdf>.
- [31] J. Kwack, G. Bauer, and S. Koric. "Performance test of parallel linear equation solvers on Blue Waters–Cray XE6/XK7 system". In: *Preceedings of the Cray Users Group Meeting (CUG2016), London, England*. 2016.
- [32] J.-Y. L'Excellent. "Multifrontal methods: parallelism, memory usage and numerical aspects". PhD thesis. Ecole normale supérieure de Lyon-ENS LYON, 2012.
- [33] J.-Y. L'Excellent and M. W. Sid-Lakhdar. "Introduction of shared-memory parallelism in a distributed-memory multifrontal solver". PhD thesis. INRIA, 2013.
- [34] Lawrence Livermore National Laboratory. *MPI Performance Topics: MPI Message Passing Protocols*. 2014. URL: <https://computing.llnl.gov/tutorials/mpi-performance/>.

## Bibliography

---

- [35] X. Li, J. Demmel, J. Gilbert, iL. Grigori, M. Shao, and I. Yamazaki. *SuperLU Users' Guide*. Tech. rep. LBNL-44289. <http://crd.lbl.gov/~xiaoye/SuperLU/>. Last update: August 2011. Lawrence Berkeley National Laboratory, Sept. 1999.
- [36] J. W. Liu. "The multifrontal method for sparse matrix solution: Theory and practice". In: *SIAM review* 34.1 (1992), pp. 82–109.
- [37] *Multifrontal massively parallel solver (MUMPS 5.1.2) users' guide*. 2017. URL: [http://mumps.enseeiht.fr/doc/userguide\\_5.1.2.pdf](http://mumps.enseeiht.fr/doc/userguide_5.1.2.pdf).
- [38] *Netlib Frequently Asked Questions*. 2006. URL: <http://www.netlib.org/misc/faq.html#2.1>.
- [39] *PaStiX User's manual*. 2013. URL: <https://gforge.inria.fr/docman/view.php/186/5707/pastixman.pdf>.
- [40] F. Pellegrini. "Scotch and libScotch 5.1 user's guide". In: (2008).
- [41] A. Pothen and S. Toledo. *Elimination Structures in Scientific Computing*. 2004.
- [42] Y. Saad and M. H. Schultz. "GMRES: A generalized minimal residual algorithm for solving nonsymmetric linear systems". In: *SIAM Journal on scientific and statistical computing* 7.3 (1986), pp. 856–869.
- [43] J. Schulze. "Towards a tighter coupling of bottom-up and top-down sparse matrix ordering methods". In: *BIT Numerical Mathematics* 41.4 (2001), pp. 800–841.
- [44] T. Steinhoff. "Singly implicit FilterRK methods for thermal-hydraulic simulations". ANODE. 2018.
- [45] Time for change. *Pros and cons of nuclear power*. 2007. URL: %20<https://timeforchange.org/pros-and-cons-of-nuclear-power-and-sustainability>.
- [46] Wikipedia contributors. *Basic Linear Algebra Subprograms — Wikipedia, The Free Encyclopedia*. [Online; accessed 14-December-2018]. 2018.
- [47] Wikipedia contributors. *Portable, Extensible Toolkit for Scientific Computation — Wikipedia, The Free Encyclopedia*. [Online; accessed 14-January-2019]. 2018.
- [48] C. J. Wu. "Partial Left-Looking Structured Multifrontal Factorization & Algorithms for Compressed Sensing". PhD thesis. UC Berkeley, 2012.
- [49] Xiaoye Li. *Direct Solvers for Sparse Matrices*. 2018. URL: <http://crd-legacy.lbl.gov/~xiaoye/SuperLU/SparseDirectSurvey.pdf>.

**Gesellschaft für Anlagen-  
und Reaktorsicherheit  
(GRS) gGmbH**

Schwertnergasse 1  
**50667 Köln**  
Telefon +49 221 2068-0  
Telefax +49 221 2068-888

Boltzmannstraße 14  
**85748 Garching b. München**  
Telefon +49 89 32004-0  
Telefax +49 89 32004-300

Kurfürstendamm 200  
**10719 Berlin**  
Telefon +49 30 88589-0  
Telefax +49 30 88589-111

Theodor-Heuss-Straße 4  
**38122 Braunschweig**  
Telefon +49 531 8012-0  
Telefax +49 531 8012-200

[www.grs.de](http://www.grs.de)