

DISCUSSION PAPER SERIES

IZA DP No. 16817

**Diversity and Discrimination  
in the Classroom**

Dan Anderberg  
Gordon B. Dahl  
Christina Felfe  
Helmut Rainer  
Thomas Siedler

FEBRUARY 2024

## DISCUSSION PAPER SERIES

IZA DP No. 16817

# Diversity and Discrimination in the Classroom

**Dan Anderberg**

*Royal Holloway University of London, CESifo,  
IFS*

**Gordon B. Dahl**

*UC San Diego, Norwegian School of Eco-  
nomics, NBER, Ifo Institute, CESifo, CEPR,  
CReAM and IZA*

**Christina Felfe**

*University of Konstanz, CESifo, CEPR and  
IZA*

**Helmut Rainer**

*University of Munich, Ifo Institute and  
CESifo*

**Thomas Siedler**

*University of Potsdam and IZA*

FEBRUARY 2024

Any opinions expressed in this paper are those of the author(s) and not those of IZA. Research published in this series may include views on policy, but IZA takes no institutional policy positions. The IZA research network is committed to the IZA Guiding Principles of Research Integrity.

The IZA Institute of Labor Economics is an independent economic research institute that conducts research in labor economics and offers evidence-based policy advice on labor market issues. Supported by the Deutsche Post Foundation, IZA runs the world's largest network of economists, whose research aims to provide answers to the global labor market challenges of our time. Our key objective is to build bridges between academic research, policymakers and society.

IZA Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

ISSN: 2365-9793

**IZA – Institute of Labor Economics**

Schaumburg-Lippe-Straße 5–9  
53113 Bonn, Germany

Phone: +49-228-3894-0  
Email: [publications@iza.org](mailto:publications@iza.org)

[www.iza.org](http://www.iza.org)

## ABSTRACT

---

# Diversity and Discrimination in the Classroom\*

What makes diversity unifying in some settings but divisive in others? We examine how the mixing of ethnic groups in German schools affects intergroup cooperation and trust. We leverage the quasi-random assignment of students to classrooms within schools to obtain variation in the type of diversity that prevails in a peer group. We combine this with a large-scale, incentivized lab-in-field-experiment based on the investment game, allowing us to assess the in-group bias of native German students in their interactions with fellow natives (in-group) versus immigrants (out-group). We find in-group bias peaks in culturally polarized classrooms, where the native and immigrant groups are both large, but have different religious or language backgrounds. In contrast, in classrooms characterized by non-cultural polarization, fractionalization, or a native supermajority, there are significantly lower levels of own-group favoritism. In terms of mechanisms, we find empirical evidence that culturally polarized classrooms foster negative stereotypes about immigrants' trustworthiness and amplify taste-based discrimination, both of which are costly and lead to lower payouts. In contrast, accurate statistical discrimination is ruled out by design in our experiment. These findings suggest that extra efforts are needed to counteract low levels of inclusivity and trust in culturally polarized environments.

**JEL Classification:** J15

**Keywords:** in-group bias, discrimination, diversity

**Corresponding author:**

Gordon B. Dahl  
Department of Economics  
University of California, San Diego  
9500 Gilman Drive  
La Jolla, CA 92093-0508  
USA  
E-mail: gdahl@ucsd.edu

---

\* We thank Samuel Bazzi, V. Bhaskar, Paul Niehaus, and David Yanagizawa-Drott, as well as seminar participants at several universities and conferences. We are grateful to the almost 20 research assistants and interns for their invaluable assistance with data collection. The project received generous financial support from the ifo Institute, University of Munich, University of St. Gallen, and University of Hamburg.

# 1 Introduction

Immigration has surged in recent decades, rendering societies increasingly diverse. In North America and Europe, the number of international migrants has risen by roughly 50% in the past 20 years, leading to an influx of individuals with culturally diverse backgrounds (McAuliffe and Triandafyllidou, 2021). Forecasts indicate a continued rise in immigration in the years to come (Hanson and McIntosh, 2016), which will perpetuate ongoing societal transformations. Sustaining social cohesion amid this heightened diversity requires trust and cooperation across national, religious, and ethnic divides.

Public education stands out as one of the few social institutions with the potential to foster such intergroup cooperation. Indeed, one argument for publicly provided schooling is to create a unified citizenry (Dee, 2004; Milligan et al., 2004). Schools could play a crucial role in promoting greater cross-cultural understanding by offering opportunities for youth from diverse backgrounds to form meaningful relationships across group boundaries. This, in turn, could break down prejudices and stereotypes, ultimately preparing them to thrive in and contribute positively to a diverse society (Allport, 1954; Pettigrew and Tropp, 2008; Wells et al., 2016; Tropp and Saxena, 2018).

As intuitively appealing as this view might be, it ignores that diversity may facilitate *or* obstruct intergroup cooperation depending on the form it takes. In settings dominated by a few large groups (polarization), the drive to establish cultural domination could grow strong, causing own-group attachment to increase and integration to decrease. In contrast, in environments where many small groups coexist (fractionalization), substantial benefits could be gained by unifying under a shared identity, which in turn would ultimately improve social cohesion.

The opposing forces of polarization and fractionalization have been explored in the context of civil conflict and nation-building in the developing world (e.g., Montalvo and Reynal-Querol, 2005; Bazzi et al., 2019). Our study brings the fractionalization-polarization paradigm to bear on the issue of intergroup cooperation in diverse classrooms in Germany. Our interest lies in whether the type of diversity that prevails in schools matters for in-group bias in trust and cooperation. Specifically, are adolescents embedded in polarized peer groups more likely to display own-group favoritism compared to adolescents with highly homogeneous or fractionalized peer groups? And what role does the cultural distance between majority and minority groups play in this?

Our setting is unique for three reasons. First, we study these forces during adolescence, a critical life-cycle stage where an individual's attitudes and biases are susceptible to socializing influences. Schools are a setting where youth interact with peers from diverse backgrounds on a daily basis, and hence create an environment where both intergroup conflict and cooperation could emerge (Allport, 1954; Lowe, 2021; Mousa, 2020). Second, we zoom in on the micro level of classrooms, whereas most of the existing literature has focused on nations and communities. Third, we study a highly developed country where immigration is the driving force for changes in diversity, a trend which will only accelerate across Europe due to demographic and economic pressures.

Estimating the impacts of polarization and fractionalization is challenging for two reasons. First,

data on in-group out-group cooperation among adolescents is scarce, and even when available, cannot readily be linked to individuals' peer groups. Second, the set of adolescents who interact with ethnically diverse peers is likely to be endogenous, resulting in selection bias. For example, youth embedded in ethnically diverse peer groups may be more open to diversity to begin with, while those exposed to an ethnically homogeneous learning environment may have parents who transmit prejudicial attitudes.

Our paper overcomes these challenges by using a unique design that combines two elements. First, on the data front, we run a large, incentivized lab-in-the field experiment in 220 classes spread across 57 German secondary schools. The experiment, based on an investment (or "trust") game, allows us to measure how native German students cooperate with in-group (other natives) versus out-group (immigrants) partners, respectively. In the game, the sender chooses how much money to transfer, this amount is multiplied by 3, and the receiver decides how much money to transfer back. Key to our experiment is that each participant is asked how they will play if paired with a native and an immigrant interaction partner, both of which are anonymous to the participant as they are drawn from a different school. Our main measure for in-group bias in cooperation is how much a sender chooses to transfer if they are paired with a native versus an immigrant.<sup>1</sup>

In addition, we ask students to fill out an extensive survey, allowing us to characterize classrooms in terms of their ethnic composition. Forty-six percent of students in our sample have an immigrant background, spanning in total 95 different countries of origin, with the largest groups coming from Turkey, Poland, and Russia. Indeed, Germany has the second largest number of international migrants behind the U.S. (McAuliffe and Triandafyllidou, 2021), providing a rich testing ground to study the effects of ethnic and cultural diversity.

Second, in terms of identification, we exploit variation in peer group diversity arising from students' quasi-random assignment to classes within schools. As we explain later, in Germany, all students have the right to receive a non-discriminatory education, and schools often have explicit rules to ensure this applies in the assignment of students to classrooms. Consistent with this, we conduct several empirical tests which provide strong support that native and immigrant students are randomly assigned to classrooms within a school. We first verify that, once school fixed effects are accounted for, individual and family background characteristics are not statistically significant predictors of classroom diversity. We also document that the actual distribution of immigrant peers closely matches the simulated distribution assuming random assignment.

Combining this variation in classroom diversity with our experimentally elicited data, we first explore how natives' in-group bias varies with the fraction of immigrants in a class, regardless of the immigrants' origin. We estimate an inverse-U shaped relationship: the in-group bias among natives initially widens as the fraction of immigrants increases, but eventually turns around and narrows again as the immigrant fraction continues to increase. The turning point occurs where

---

<sup>1</sup>Participants are asked to play the game as both sender and receiver. We use the decisions of receivers to study mechanisms.

there is close to an even split of natives versus immigrants in a class. In other words, when natives are the dominant majority, adding in extra immigrants to a classroom increases the in-group bias among natives. But once natives start to lose their majority status, adding in extra immigrants decreases the bias. At the turning point, a one-standard deviation change in the immigrant share (21 pp) shrinks the in-group out-group investment gap by 10% of a standard deviation.

In a second step, we investigate how the cultural background and relative size of immigrant subgroups matter. To that end, we no longer treat immigrants in a classroom as a single group, but consider cultural heterogeneity among them. For our main set of results, we utilize religious affiliation as the differentiating factor; specifically, we distinguish between immigrants who are Muslim and those who are not. This split allows us to compare immigrants who are culturally distant from their native peers versus those who are culturally more similar. We estimate that the highest level of in-group bias among natives occurs in classrooms characterized by cultural polarization, where a slight majority of natives coexists with a substantial minority of Muslim immigrants.

In sharp contrast, in classrooms exhibiting other forms of diversity, natives' tendency towards in-group favoritism is 32 to 41% of a standard deviation lower compared to culturally polarized classrooms. This encompasses scenarios with (i) *non-cultural polarization*, where a slight majority of natives coexist with a substantial minority of non-Muslim immigrants, (ii) *immigrant fractionalization*, where there is a slight majority of natives and immigrants are equally divided between Muslims and non-Muslims, (iii) *overall fractionalization*, where natives, Muslim immigrants, and non-Muslim immigrants are equally large groups, and (iv) a *native supermajority*. Several robustness checks substantiate that in-group bias among natives peaks when classrooms exhibit a high level of cultural polarization.

We next explore three possible mechanisms: accurate statistical discrimination, negative stereotypes (i.e., inaccurate statistical discrimination), and taste-based discrimination. An important rationale for transferring money in the investment game we had students play is *trust*. This refers to the sender's beliefs about the recipient's trustworthiness and their expectation of generous reciprocation. In our setting, accurate statistical discrimination would be a correct belief among natives that the immigrants they play with are less trustworthy than natives. We set up our experiment so that any accurate statistical discrimination would be constant across classroom types—hence, by design it cannot explain the heightened in-group favoritism observed within culturally polarized classrooms versus other classroom configurations.

An alternative mechanism is that natives in polarized classrooms form negative stereotypes. The idea is that natives in polarized classrooms could rationally develop more mistrust towards the immigrants in their classroom because of limited or adversarial interactions with them. If natives inaccurately generalize their classroom experiences, they could develop negative stereotypes, leading them to place less trust in *all* immigrants compared to natives. Consequently, they would rationally favor fellow natives over immigrants in general, given their biased beliefs. We present two pieces of evidence consistent with culturally polarized classrooms fostering negative stereotypes. First, within

culturally polarized classrooms, natives are treated less favorably by immigrants than by fellow natives relative to other classroom types. This conclusion is based on findings for the decisions of native and immigrant classmates when they play as receiver and are asked how much money they would return to a native sender. Second, we use questions from our survey which capture respondents' trust towards natives and immigrants in the general population. We estimate that the in-group out-group trust gap among natives reaches its peak in culturally divided classrooms. The gap is between 35 to 50% of a standard deviation higher when contrasted with non-culturally polarized, immigrant fractionalized, overall fractionalized, and native majority classrooms.

A third possible mechanism for in-group bias is taste discrimination. That is, natives could simply prefer to give money to fellow natives over immigrants. We can explore the role of taste-based discrimination since we utilized the strategy-vector method in our experiment. This design choice implies the choices of students when playing the role of receiver in our investment game are akin to the choices made in a standard dictatorship game. We find some evidence that taste-based discrimination is more pronounced in culturally polarized classrooms, but with a peak that occurs when there is a somewhat larger share of natives versus Muslim immigrants relative to our in-group bias estimates. As an additional exercise, we use a survey-based measure of anti-immigrant sentiment related to the labor market as a proxy for taste discrimination. This measure of anti-immigrant sentiment peaks in culturally polarized classrooms, with estimates which are significantly different from the other four classroom types.

As first pointed out by Becker (1957), discrimination can be either profit maximizing or costly, depending on whether it reflects (accurate) statistical discrimination or taste-based discrimination, respectively. More recently, researchers have argued that inaccurate statistical discrimination (Bohren et al., forthcoming) is also costly, as individuals will make decisions based on incorrect beliefs. Since accurate statistical discrimination is ruled out by design, the implication is that natives in culturally polarized classrooms should experience payoff losses due to their in-group bias. The reason is that the remaining explanations of negative stereotypes (i.e., inaccurate beliefs) and taste discrimination, both of which we find evidence for, are costly. In line with this prediction, we find that payoff losses peak in culturally polarized classrooms. Payoff losses are 0.33 to 0.46 standard deviations higher in culturally polarized classrooms compared to classes characterized by non-cultural polarization, immigrant fractionalization, overall fractionalization, and a native supermajority.

To help interpret our findings, we outline a simple model grounded in the idea that in-group bias arises from limited meaningful intergroup interactions (Allport, 1954). The model integrates a social interactions model à la Brock and Durlauf (2001) into a Hotelling-type structure (Hotelling, 1929) with multiple interacting types. It features individuals of different types actively deciding who to identify and socially engage with. In a version of the model tailored to our empirical setting, individuals fall into three types: natives, culturally close immigrants, and culturally distant immigrants. They face the choice of *either* joining an “inclusive” group with a common shared

identity *or* preserving their innate cultural identity in an “exclusive” group composed solely of others of their own type. This decision entails a trade-off between the desire to be part of a large group and a preference for one’s own cultural identity. The identity of the inclusive group shifts endogenously with the distribution of types in the population. Individual heterogeneity is captured through additive extreme value-distributed preferences for inclusive versus exclusive group membership. Due to a social multiplier, small shifts in the population type distribution can lead to large changes in inclusivity. Specifically, consistent with our empirical evidence, comparative statics show a low willingness of natives to engage with immigrants in the inclusive group under cultural polarization, as opposed to scenarios involving non-cultural polarization, fractionalization, or a native supermajority.

Our study contributes to several strands of research. Diversity and its consequences for social cohesion have long been studied in economics. Examples of earlier work include Alesina and La Ferrara (2000, 2002) and Luttmer (2001), while Alesina and Tabellini (forthcoming) provide an up-to-date survey. Recent evidence indicates that migration-induced diversity can either hinder or improve social cohesion depending on the context. In France, it has been shown to strain social bonds among neighbors and reduce the quality of local public goods (Algan et al., 2016). However, in the U.S., the outmigration of millions of African Americans from the South to the North in the mid-20th century was associated with increased support among whites for racial equality (Calderon et al., 2023). To our knowledge, only one prior study has used individual-level data to disentangle the effects of diversity into its polarization and fractionalization components. Bazzi et al. (2019) examine a resettlement program in Indonesia which relocated millions of ethnically diverse migrants. In cases where the program led to fractionalized communities, individuals developed shared identities, whereas in polarized communities, attachment to one’s own ethnic group remained strong. Like Bazzi et al. (2019), we hone in on the different components of diversity, but study its consequences in a wealthy country where migration-induced diversity is only going to become more pervasive in the years ahead. In contrast to a long-standing argument that diversity is less divisive in rich nations (Horowitz, 1985), we find that polarization and fractionization matter for social cohesion in a way akin to the developing world.

Our study also relates to recent work examining social cohesion in school and university settings. For example, Alan et al. (2021) study a perspective-taking intervention in the aftermath of the influx of Syrian refugees into Turkish elementary schools, showing that it contributed to fewer physical conflicts among peers, less social exclusion, and improved inter-ethnic social ties. Boucher et al. (2022) examine a classroom intervention where 5-year old Turkish and Syrian refugee children were randomly brought into contact, discovering that such exposure resulted in an increased formation of interethnic friendships. Rao (2019) exploits a policy change in India, where poor students were integrated into elite private schools, finding that economically disadvantaged classmates makes wealthier students more prosocial, generous, and egalitarian. Carrell et al. (2019) and Corno et al. (2022) leverage the random assignment of roommates in university settings in the U.S. and South Africa to demonstrate that students become more empathetic and open to forming friendships

with members of the social groups to which their roommates belong. These studies suggest that intergroup contact promotes integration. In contrast, our study highlights that classroom contact alone might not be sufficient, but in fact can be counteracted by own-group attachment in culturally polarized classrooms.

As emphasized in the literature, peer groups are central in the socialization of adolescents, playing a key role in shaping their attitudes and behaviors (Brown, 2011). A large literature studies peer effects for outcomes as diverse as education, crime, drug use, and teenage pregnancy (see, e.g., Kremer and Levy, 2008; Bifulco et al., 2011; Sacerdote, 2011; Lavy et al., 2012; Ohinata and Van Ours, 2013; Brenøe and Zölitz, 2020; Figlio et al., forthcoming). Our paper adds to this literature by studying how the type of diversity in a peer group matters (polarization versus fractionalization) and by examining intergroup cooperation and trust.

In our own previous work (Felfe et al., 2021), we utilized the experimental data we collected from German schools to investigate a completely different question for a different group. Specifically, we studied how immigrants' prosocial behavior was affected by a reform which granted them citizenship from birth. We also used the survey data we collected to show that the same birthright citizenship reform had unintended consequences for the well-being of immigrant girls (Dahl et al., 2022).

The remainder of the paper is organized as follows. We first explain our setting, the experiment, and the survey. In Section 3, we discuss our empirical design and how we model classroom diversity. Section 4 presents our results on how in-group bias varies both with the amount and type of diversity in a classroom, and Section 5 explores possible mechanisms. The following section outlines a simple model which is consistent with our empirical findings. The final section concludes.

## 2 Study Design

### 2.1 Setting

We study how diversity affects in-group bias in the context of classrooms in Germany, which have varying mixes of native and immigrant students. We received permission to enter classrooms and collect data from all students completing their final year of secondary school for 57 schools (222 classrooms). These students were either in 9<sup>th</sup> or 10<sup>th</sup> grade, and hence mostly 15 or 16 years old. In the German school system, when students enter secondary school in fifth grade, they are assigned to a class. Typically these students remain together in the same class until the end of compulsory school, making their exposure to classroom peers long term.

We collected data in two distinct phases. From June 2 and July 15, 2015, we visited 31 schools across five cities in the German state of Schleswig Holstein (SH), where we gathered data from all 122 ninth-grade classes. From October 19 to November 16, we surveyed students in 100 tenth-grade

classes located in 26 schools in two cities within the state of North Rhine Westphalia (NRW).<sup>2</sup>

Our data collection consisted of two main components. First, we implemented an incentivized lab-in-the-field experiment to measure native and immigrant students' willingness to cooperate with in-group versus out-group partners. Second, we administered a survey to collect comprehensive family background information, including details on the students' ethnicity and religious affiliation.

The research was conducted during regular class periods, when students would otherwise be engaged in standard classroom learning. The regular teacher either remained inconspicuous at the back of the classroom or left the room. Meanwhile, members of the research team and/or trained research assistants introduced and supervised the study. Students were seated at their usual desks, with mobile privacy screens placed between them. The study was administered using traditional paper and pencil methods and spanned two consecutive class periods, each lasting 45 minutes. Whether the experiment or the survey was administered first was randomized on a daily basis.

Our study was planned and executed in close collaboration with the education ministries of the federal states of Schleswig-Holstein and North Rhine Westphalia. Prior to commencing our research, we sought approval from these ministries and obtained consent from school principals. These school principals, in turn, communicated with parents to inform them about the upcoming study and allow them to opt out. Out of the 4,634 students present at school on the day of the study, only 44 parents (<1%) chose not to have their child participate. At the beginning of each session, we emphasized to the students that their participation was entirely voluntary and anonymous. To ensure anonymity, every student received a unique, anonymous identity. 154 students (3.5%) chose not to take part in our study, leaving us with 4,436 students. Of those, 342 did not provide basic survey information essential for our analysis, namely, the respondent's gender, whether their parents were born in Germany, their parents' countries of birth, and the respondent's religious affiliation. Thus, our baseline sample comprises 4,094 students, distributed across 222 classes in 57 schools.

## 2.2 Measuring Classroom Diversity

To measure classroom diversity, we first use the survey data we collected to categorize students into two distinct groups: (i) native students, defined as those with both parents born in Germany, and (ii) students with an immigrant background, defined as those with at least one parent born abroad.<sup>3</sup> In our baseline sample of 4,094 students, 2,216 (54%) are natives, while 1,878 (46%) are

---

<sup>2</sup>In both federal states, a school year starts in August/September and ends in June/July. One important difference is that compulsory schooling lasts 9 years in SH, but 10 years in NRW. Thus, in our sample, students from the two states are from the same school starting cohort. Our secondary schools span various types: 10 general schools ("Hauptschule"); 8 intermediate schools ("Realschule"); 29 comprehensive schools without the final years of grammar school education ("Gesamtschule ohne gymnasiale Oberstufe"); 8 comprehensive schools with the final years of grammar school education ("Gesamtschule mit gymnasialer Oberstufe"); and 2 grammar schools or high schools ("Gymnasium").

<sup>3</sup>There are 41 students who are Muslim and whose parents were both born in Germany. We classify these observations as students with an immigrant background, as they are likely third generation Turkish immigrants. Many Turks came to Germany during the guest worker program of the 1960s, never gained Turkish citizenship, and

immigrants.<sup>4</sup>

The main goal of our analysis is to examine how the in-group biases of native German adolescents are affected by the diversity of their school peer group, and the role played by cultural distance. We explore two dimensions of peer group diversity: (i) intergroup diversity, which is based on the share of immigrants in a classroom, irrespective of cultural background; and (ii) within-immigrant diversity, stemming from the cultural heterogeneity and relative sizes of immigrant subgroups within a classroom.

To model intergroup diversity, we compute the proportion of classmates with an immigrant background as the leave-one-out share.<sup>5</sup> On average, native students are in classrooms where 38% of their peers are immigrants, with a standard deviation of 21% (see Table 1). Figure 1 shows two histograms: one for the number of classes with different shares of immigrant peers (Panel A) and one for the number of native student observations in classes with different shares of immigrant peers (Panel B). As we will show in the next section, even after we take out school fixed effects, there is wide variation in the fraction of immigrants in a class.

To study the role of within-immigrant diversity, instead of treating immigrants in a classroom as a homogeneous group, we consider their cultural heterogeneity. For our main set of results, we focus on a binary categorization which separates immigrants into those who are culturally close to versus far from their native peers. In the context of Germany, a natural criterion for assessing cultural distance is religion. Specifically, we differentiate between immigrants with a Muslim background and those with a non-Muslim background. Table 1 reveals that, on average, 19% of a native’s classmates have a Muslim immigrant background (std. dev.=18%) and 19% have a non-Muslim immigrant background (std. dev.=11%). As we will show in the next section, even after we take out school fixed effects, there is wide variation in the fraction of Muslim versus non-Muslim immigrant peers in a classroom. In a robustness check, we will alternatively define cultural distance using linguistic differences between the German language and the country-of-origin language for an immigrant’s parents. Furthermore, we also explore the impact of classroom polarization versus fractionalization by decomposing classrooms into a native group and 11 distinct immigrant subgroups based on their parents’ country or region of origin.

---

remained Muslim. These Turkish students would likely be identified by natives as culturally distant immigrants. If we alternatively exclude these 41 observations, the empirical results are virtually identical.

<sup>4</sup>In the general population, the share of individuals age 15-20 with an immigrant background was 34% as of 2018 (Federal Statistical Office of Germany, 2009); our share is higher because we purposefully targeted cities with many immigrants.

<sup>5</sup>To compute this share and other measures of classroom diversity, we use survey information from all of the 4,094 students in our baseline sample, of which 2,216 are natives. 53 native students failed to complete the first stage of our experiment. Two classes of the 222 classes in the baseline sample only have immigrants, and hence are not used in our analysis. Therefore, our estimation sample contains 2,163 native students spread over 220 classes. Summary statistics correspond to the 220 classrooms attended by these 2,163 native students.

## 2.3 Measuring In-Group Bias

**The Experiment.** To measure in-group bias, we used a modified version of the investment game originally introduced by Berg (1995). This choice was informed by the expanding body of research on discrimination between real social groups, which can be traced back to Fershtman and Gneezy’s (2001) seminal work.

The investment game is a two-player scenario, with one participant acting as the sender (first-mover) and the other as the receiver (second-mover). Both players start with an initial endowment, in our case, 5 euros. The sender makes the initial decision regarding how much of their endowment to send to the receiver,  $x \in [0, 5]$ , with the constraint that they can only send in 50 cent increments. We then triple the amount sent to the receiver. The receiver, now possessing an amount of  $5 + 3x$  euros, decides how much to send back to the sender. They can send an amount  $y \in [0, 5 + 3x]$  in 10 cent increments. Consequently, the sender exits the game with  $5 - x + y$  euros, while the receiver exits with  $5 + 3x - y$  euros.<sup>6</sup>

In our experiment, we utilized the strategy method, which meant that each participant had to make decisions both as the first-mover and as the second-mover. We first had participants assume the role of first-movers. Importantly, first-movers had to factor in the gender and migration background of their potential interaction partners when making investment decisions. This was achieved by letting first-movers decide the amounts they wished to allocate to a male with German parents ( $I_1$ ), a female with German parents ( $I_2$ ), a male with foreign parents ( $I_3$ ), and female with foreign parents ( $I_4$ ).<sup>7</sup>

After participants completed the initial stage of the investment game, they were asked to specify their expected back transfer from each of the four possible interaction partners. This expectation was recorded within a range of 0 to 20 euros, in increments of ten cents.

In the final stage of the investment game, participants assumed the role of second-movers, and we employed the contingent response method to elicit their back transfers (returns). For instance, on one decision sheet, participants were tasked with determining their back transfers to a male with German-born parents for each of the eleven possible investments (0, .5, 1, 1.5, ..., 5) made by a male with German-born parents as the first mover. Using a similar approach, we collected back payments for the other potential interaction partners. Participants were allowed to specify amounts between 0 and  $5 + 3x$  euros, in ten cent increments.

Prior to commencing the experiment, written instructions were provided to all students in the class, with an experimenter verbally going through the instructions with the students as well.

---

<sup>6</sup>Assuming self-interested preferences, the only subgame-perfect equilibrium has no investment and zero returns. In contrast, “full” cooperation, where the first-mover invests their entire endowment, maximizes the players’ joint payoff.

<sup>7</sup>Additionally, we inquired about the amounts students would be willing to send to a boy with foreign parents holding a German passport ( $I_5$ ) and a girl with foreign parents who possessed a German passport ( $I_6$ ). This line of questioning was included in the data collection process because our study was initially designed to explore the integration of immigrant children who had acquired German citizenship.

The translated instructions and the decision sheets are available in Appendix C. The students were informed that they would engage in the investment game, initially as the first-mover and subsequently as the second-mover. Students were told they had the opportunity to earn actual money, and their eventual earnings would be contingent on their individual decisions as well as those of another participant. Importantly, we explicitly told students they would be randomly assigned to play with an anonymous student from a completely different school (i.e., not from their own class or school).

To be more concrete, we informed participants that we would determine their ultimate earnings through the following steps: (i) we would randomly pair two participants from different schools in the same federal state; (ii) we would then randomly assign the roles of first-mover and second-mover; (iii) the relevant characteristics of both the first-mover and the second-mover, namely gender and immigration background, would be established using survey information; (iv) we would execute the investment decision made by the first-movers using the characteristics of the second-mover; (v) the back transfer decision of the second-mover would be implemented, taking into account the characteristics of the first-mover and their investment choice from step (iv); (vi) finally, based on the combination of choices from steps (iv) and (v), we would calculate the participants' ultimate earnings. Payments were disbursed within a two-week timeframe and delivered in sealed envelopes, each bearing the student's distinctive identification code, by either the head teacher or the school's secretary.

In applying this method to determine participants' earnings, we categorized mixed-background children (those with one German-born and one foreign-born parent) as having foreign-born parents.<sup>8</sup> This categorization lines up with how we define our measures of classroom diversity in Section 2.2.

**In-Group Out-Group Investment Gap.** To create our main dependent variable, we rely on the investment decisions of native students during the first stage of the investment game. On average, native students invest 2.86 euros, which corresponds to 57 percent of their initial 5 euro endowment. Figure 2 shows a histogram of all possible investment decisions made by the native students in our sample. The two most frequent investment choices are transfers equivalent to either 50% or 100% of their initial endowment. These patterns resemble those found in similar lab-in-the-field experiments.<sup>9</sup>

To measure the in-group bias in cooperation among native students, we take the difference between their average investments in fellow natives averaged over both genders, given by  $\frac{1}{2}(I_1 + I_2)$ , and their average investments in immigrants, given by  $\frac{1}{2}(I_3 + I_4)$ . We denote this measure as the in-group out-group investment gap (IG). The summary statistics for this measure are included in Table 1. Among native students, the average IG is 0.09 euros and the standard deviation is 0.76. Hence, on average there is not much in-group favoritism, but it varies considerably.

---

<sup>8</sup>There were no questions raised by participants about the handling of mixed-background children during the experiment.

<sup>9</sup>See, for example, Bellemare and Kröger (2007) and Falk and Zehnder (2013).

We note that experimenter demand effects and social desirability bias should be constant across classroom types, unless these factors are directly affected by diversity. If they are directly affected by classroom diversity, then this is a mechanism for our findings, rather than a threat to identification.

## 2.4 Estimation Sample

Our main estimation sample consists of 2,163 native German students. Summary statistics in Table 1 indicate that these students were 15.8 years old on average at the time of the study. Approximately half of them are male (54%), and a majority identify as Christians (68%), specifically Catholics or Protestants. In terms of socioeconomic background, most students come from families falling into one of the following categories: (i) two-parent households where at least one parent has a high level of education (20%); (ii) two-parent households with both parents having a low level of education (27%); and (iii) single-parent households with a parent who has a low level of education (27%).

## 3 Empirical Approach

### 3.1 Identification

We are interested in estimating how the ethnic composition of classroom peers affects native students' in-group out-group bias. As pointed out by Manski (1993), two challenges in identifying peer effects are correlated unobservables and endogenous group formation. Translated to our setting, the first challenge is that the ethnic makeup of a classroom is likely correlated with both observable and unobservable characteristics, such as average family income and attitudes towards immigrants. The second challenge is that students (and their families) self-select into the neighborhoods and schools they attend in ways which are likely to create a bias.

To deal with these challenges, we take advantage of the quasi-random assignment of students to classrooms within schools. The idea is that while the school a student attends is unlikely to be random, which classroom they are assigned to within a school is as good as random. We model outcomes (e.g., in-group out-group bias) for native individual  $i$  in classroom  $k$  in school  $s$  as:

$$Y_{i,k,s} = \alpha + f(\text{diversity}_k) + \delta X_i + \gamma Z_k + \theta_s + \epsilon_{i,k,s} \quad (1)$$

where  $f(\cdot)$  is a function of the mix of ethnic peers ( $\text{diversity}_k$ ) a native is exposed to in their classroom. We will model classroom diversity either as a function of the fraction of immigrant peers or as a function of the fraction of culturally distant and the fraction of culturally close immigrant peers. The vector  $X_i$  contains individual and family background characteristics,  $Z_k$  is the number of students in the class, and  $\epsilon_{i,k,s}$  is the error term. Crucially, the estimating equation includes school fixed effects,  $\theta_s$ , to account for the fact that individuals are not randomly assigned to schools.<sup>10</sup>

---

<sup>10</sup>Note that all but one of our outcome variables are calculated as within-person differences. For example, the

Figure 3 illustrates the identifying variation we use when estimating our model. The left panel displays a histogram for the fraction of immigrant peers in classes, where we have first regressed out the school fixed effects. The scatterplot in the right panel likewise shows the mix of Muslim versus non-Muslim immigrant peers in a classroom after netting out school fixed effects. Both graphs reveal substantial residual variation in our diversity measures.

Our identification strategy is related to work which exploits either within school variation across classes or natural variation in cohort composition across time within a given school (Antecol et al., 2015; Hoxby, 2002; Hanushek et al., 2003; Carrell and Hoekstra, 2019; Carrell et al., 2018; Balestra et al., 2022).<sup>11</sup>

### 3.2 Validity

Our identification strategy relies on quasi-random assignment of students to classes within a school. In Germany, all students have the right to receive a non-discriminatory education, and schools often have explicit rules to ensure this applies in the assignment of students to classrooms.

The right to a non-discriminatory education evolved over time. In 1971, the Standing Conference of the Ministers of Education and Cultural Affairs of the German Federal States issued a guideline that immigrants should be treated equally in all matters related to schooling (Puskeppeleit and Krüger-Potratz, 1999). Despite this, up to the 1990’s classrooms were often partially segregated along native versus immigrant lines. A major citizenship reform introducing birthright citizenship in 2000 coincided with a renewed push to integrate classrooms (Nieden and Karakayali, 2016).

In a high profile case in Berlin in 2012, German parents lobbied for and were successful at creating classrooms which were highly segregated. Immigrant parents filed a complaint with the Berlin Senate arguing this was a discriminatory practice and the Senate ruled in their favor, requiring students to be reassigned in a non-segregated fashion (Nieden and Karakayali, 2016). Today, all German federal states have school laws that specify a right to non-discriminatory education (Federal Anti-Discrimination Agency, 2022).

Although these laws do not explicitly mandate random assignment, in practice many schools state that they do not discriminate based on migration background when assigning students to classrooms.<sup>12</sup> In our own conversations with school officials and principals, we were also told that migration background is not used as a criterion for making classroom assignments. Consistent with this, a recent government report highlights that in Germany “the assignment of students to classes

---

investment gap (IG) is the difference in how much a native invests in a native versus how much they invest in an immigrant. This specification nets out the constant effect within an individual, which provides an increase in precision.

<sup>11</sup>Another approach uses the random assignment of peers to social groups. Prominent examples are the random assignment of students to classrooms, such as in project STAR (Chetty et al., 2011), roommates in university dorms (Sacerdote, 2001), freshmen to university sections (Feld and Zölitz, 2017), and cadets to squads in the military (Lyle, 2007; Dahl et al., 2021).

<sup>12</sup>See, for example, <https://www.goetheschule-asperg.de/index.php/eltern/faq> (accessed November 14, 2023).

is difficult to influence” (Federal Government of North Rhine-Westphalia, 2022, p. 92).<sup>13</sup>

Several empirical tests provide strong support that students are quasi-randomly assigned to classrooms. Our first test appears in Table 2. We regress the fraction of immigrant peers in a native’s class on background characteristics of the native student and the number of students in the class. In column 1, we do not include school fixed effects. Several variables have sizable and statistically significant effects. In particular, native students with a higher fraction of immigrant peers are less likely to be Protestant or non-religious, are more likely to come from families with lower socio-economic backgrounds, and are older. An F-test reveals that the variables are also jointly significant (p-value<.001). Column 2 of Table 2 reports what happens when school fixed effects are added to the same regression. All of the coefficient estimates are now close to zero, with none of the 15 variables being individually statistically significant. The joint F-test is not statistically significant either (p-value=.279). The tests in column 2 are similar to the balancing tests performed for actual experiments, to check whether random assignment to treatment has been implemented correctly.

In Appendix Table A1, we perform a similar test, but this time using the fraction of Muslim immigrant peers as the dependent variable. As before, when school fixed effects are not included, the covariates are both individually and jointly statistically significant. But when school fixed effects are added to the regression, only one of the coefficients is significant at the 10% level (roughly what would be expected by chance) and the covariates are not jointly significant.

Related to these balancing tests, for our main regression model in Section 4, we explore what happens when we add additional covariates beyond the school fixed effects. As expected with conditional random assignment, the estimates are virtually identical.

Finally, we conduct a simulation test. Following Carrell and West (2010), we randomly assign students to classrooms within schools, using the actual share of immigrants at the school level and the actual class sizes within a school. We repeat this counterfactual exercise 1,000 times. To test random assignment, for every set of re-sampled classes, we calculate the empirical p-value as the proportion of simulations where exposure to immigrant students is smaller than that observed in the original class. If class composition is random, the distribution of p-values within a school should be approximately uniform, which is testable using a one-sample Kolmogorov-Smirnov test. We reject uniformity for just one out of 57 schools in our sample at the 5% confidence level, and obtain the same result (rejecting uniformity just once) when repeating the simulation exercise for the share of non-Muslim immigrants.

In conclusion, the battery of tests we conduct all provide strong support for the quasi-random assignment of immigrant and native students to classrooms.

---

<sup>13</sup>There is limited scope to influence assignment; for example, while some schools allow parents to name one friend they would like to have their child be in the same class with, others state this type of preferences will not be considered.

### 3.3 Modeling Classroom Diversity

To model classroom diversity,  $f(\text{diversity}_k)$ , we build on the measures of Montalvo and Reynal-Querol (2005) which distinguish between ethnic polarization and ethnic fractionalization:

$$\text{polar}_k = 4 \sum_{j=1}^N \pi_{j,k}^2 (1 - \pi_{j,k}) \quad \text{and} \quad \text{frac}_k = 2 \sum_{j=1}^N \pi_{j,k} (1 - \pi_{j,k}), \quad (2)$$

where  $\pi_{j,k}$  represents the proportion of individuals from ethnic group  $j$  in classroom  $k$ . The polarization index reaches its maximum value when a classroom is characterized by a bipolar composition, meaning it consists of only two large and equally sized groups. Conversely, the fractionalization index attains its maximum value when the proportions of each of the  $N$  groups within the classroom are equal, each accounting for  $1/N$  of the total population.

In the first step of our analysis, we view classrooms as being composed of two groups: natives (with share  $1 - \pi_k$ ) and immigrants (with share  $\pi_k$ ). In the case of only two groups, where  $\text{polar}_k = 4\pi_k(1 - \pi_k)$  and  $\text{frac}_k = 2\pi_k(1 - \pi_k)$ , the indices are equivalent up to a scaling factor, both reaching their maximum when natives and immigrants constitute equally large groups. Therefore, we start by modeling  $f(\text{diversity}_k)$  in equation 1 as a quadratic polynomial,  $\beta_1\pi_k + \beta_2\pi_k^2$ , where  $\pi_k$  denotes the share of immigrant classmates. Our first regression specification is:

$$Y_{i,k,s} = \alpha + \beta_1\pi_k + \beta_2\pi_k^2 + \gamma X_i + \delta Z_k + \theta_s + \epsilon_{i,k,s} \quad (3)$$

In the second step of our analysis, we consider two immigrant subgroups within classrooms: immigrants culturally close to their native peers and culturally distant immigrant peers. As laid out in Section 2.2, our main proxy for cultural distance is based on religion, specifically the distinction between non-Muslim and Muslim immigrants. Thus, a classroom is now defined as being composed of non-Muslim immigrants (share  $\pi_{C,k}$ , with  $C$  for culturally “close”), Muslim immigrants (share  $\pi_{D,k}$ , with  $D$  for culturally “distant”), and native Germans (share  $\pi_{N,k}$ , with  $N$  for “natives”). In this case, using  $\pi_{N,k} = 1 - \pi_{C,k} - \pi_{D,k}$ , we can express  $\text{polar}_k$  and  $\text{frac}_k$  as functions of the share of non-Muslim immigrant peers ( $\pi_{C,k}$ ) and the share of Muslim immigrant peers ( $\pi_{D,k}$ ):

$$f(\text{diversity}_k) = \underbrace{4(\pi_{C,k} + \pi_{D,k} - \pi_{C,k}^2 - \pi_{D,k}^2 - \pi_{C,k}\pi_{D,k})}_{=\text{frac}_k} - \underbrace{12(\pi_{C,k}\pi_{D,k} - \pi_{C,k}^2\pi_{D,k} - \pi_{C,k}\pi_{D,k}^2)}_{=\text{polar}_k} \quad (4)$$

The first term in brackets is the expression for the fractionalization index. The sum of the two terms on the right hand side is the reformulated polarization index. Thus, our second regression specification, which flexibly nests both polarization and fractionalization, is:

$$\begin{aligned}
Y_{i,k,s} = & \alpha + \beta_1\pi_{C,k} + \beta_2\pi_{D,k} + \beta_3\pi_{C,k}^2 + \beta_4\pi_{D,k}^2 + \beta_5\pi_{C,k}\pi_{D,k} + \\
& \beta_6\pi_{C,k}^2\pi_{D,k} + \beta_7\pi_{D,k}^2\pi_{C,k} + \gamma X_i + \delta Z_k + \theta_s + \epsilon_{i,k,s}
\end{aligned}
\tag{5}$$

Note that the specification of classroom diversity in equation 5 is close to a third order expansion of the terms  $\pi_{C,k}$  and  $\pi_{D,k}$ . The only difference is that a third-order expansion would have also included  $\pi_{C,k}^3$  and  $\pi_{D,k}^3$ .

In an extension in Section 4.3, we will also consider more than two immigrant groups based on geographical regions. For this extension, we model diversity using the expressions for polarization and fractionalization in equation 2.

## 4 Results

### 4.1 Share of Immigrant Peers in the Classroom: An Inverted-U Relationship

We start by examining how natives' in-group bias varies with the proportion of immigrant peers in a classroom in Table 3. The dependent variable is the in-group out-group investment gap. For comparison purposes, we start with a linear specification in panel A. In column 1, we include only basic controls: student's gender, age, the class size, and school fixed effects. There is no statistically significant linear relationship between the in-group bias of native students and the proportion of immigrant peers. This null effect continues to hold when we additionally include controls for students' religious background (column 2) and family background (column 3). Based on panel A, one might be tempted to conclude that in-group bias is unaffected by classroom diversity, but this would be incorrect, as panel B demonstrates.

In panel B, we model diversity as a quadratic function of the proportion of immigrant peers, as specified in equation 3. No matter which set of controls are included, the coefficients on the polynomial terms are individually and jointly significant (p-value for joint F-test = .0002). Moreover, we cannot reject that the peak in in-group bias occurs when polarization is highest (50% natives and 50% immigrants).<sup>14</sup>

Figure 4 illustrates this inverted-U shaped relationship. The x-axis is the fraction of immigrant peers in a classroom and the y-axis is the investment gap in standard deviation units. The in-group out-group investment gap initially widens as the proportion of immigrant classmates increases, but then it reverses course and starts to narrow again as the share of immigrant peers continues to rise. At the turning point of a 45% immigrant share, a one-standard deviation change in the immigrant share in either direction (0.21) results in a statistically significant 10% of a standard

---

<sup>14</sup>The test is whether the two quadratic coefficients are equal but opposite in sign; this test has a p-value of .248 in panel B, column 3.

deviation reduction in in-group bias (p-value=0.018). A similar pattern emerges when estimating diversity even more flexibly by using a third-order polynomial in the fraction of immigrant peers (see Appendix Figure A1).

## 4.2 Diversity Among Immigrant Peers and the Role of Cultural Distance

Our first analysis modeled classroom diversity solely as a function of the proportion of immigrants in a classroom, treating immigrants as a homogeneous group. A natural question is whether diversity within the immigrant group matters. Once there are more than two groups (natives versus immigrants), the distinction between polarization and fractionalization also becomes relevant. This is because the expressions for polarization and fractionalization in equation 2 are equivalent up to a scaling factor when there are just two groups, but differ in the case of three groups or more (equation 4).

With this in mind, we now shift our focus to exploring how diversity among immigrants matters, distinguishing between immigrants with different cultural backgrounds. Using the three group case, we ask the following questions. What role does cultural diversity within the immigrant group play in shaping the inverted-U relationship between natives' in-group bias and the immigrant share? Is this relationship due to a high level of polarization, where students are exposed to a two-point symmetric distribution of classmates, with natives being equal in number to one large and culturally distinct group of immigrants? Or is it explained by fractionalization, where natives, culturally close immigrants, and culturally distant immigrants are each minorities?

In Germany, perhaps the most salient measure for whether an immigrant is culturally close or distant to a native German is captured by religious affiliation – specifically, whether an immigrant is Muslim. Our second analysis uses the share of Muslim and non-Muslim peers to flexibly model the effects of classroom diversity using equation 5. The results can be found in Table 4. The estimates describe a three-dimensional surface where the x-axis is the share of culturally distant peers, the y-axis is to the share of culturally close peers, and the z-axis is the predicted in-group bias among natives (as measured by the in-group out-group investment gap).

Instead of plotting three-dimensional figures, we show heatmaps, as they make it easier to visualize our results. A two-dimensional heatmap plots the predicted values of the in-group out-group investment gap using colors, with darker shading representing higher values and lighter shading representing lower values. The heatmap depicted in Figure 5 corresponds to the estimates appearing in Table 4.<sup>15</sup>

The peak in native's in-group bias is observed in classrooms marked by what we define as *cultural polarization (CP)*. This occurs when native German peers constitute a slight majority of 58%, while

---

<sup>15</sup>The range for the heatmap is limited to the areas where we have nontrivial amounts of data: non-Muslim shares between 0% to 50%, Muslim shares between 0% to 70%, and the combined total of these shares not exceeding 70%. All control variables are evaluated at their means.

Muslim immigrant peers make up a large minority of 42%. The table displayed to the right of the heatmap in Figure 5 calculates the difference between predicted in-group bias at its peak (*CP*) and four other classroom scenarios. Each of these scenarios are also labeled on the heatmap.

In the first scenario, referred to as *non-cultural polarization (NCP)*, native Germans continue to make up a slight majority of 58%, but now the non-Muslim immigrants are the large minority group of 42%. In comparison to the peak in-group bias found in culturally polarized classrooms, native students in these classrooms exhibit 35% of a standard deviation lower in-group bias.<sup>16</sup> This can be seen visually in the heatmap by the darker color shading at *CP* versus *NCP*. In the second scenario, which we term *immigrant fractionalization (IF)*, native peers once again constitute a slight majority of 58%, while both Muslim and non-Muslim immigrant peers now form two medium-sized minority groups, each comprising 21% of the class. In comparison to the peak in culturally polarized classrooms, the in-group bias of native students is 32% of a standard deviation lower in this setting. Revisiting the inverted-U relationship between the in-group bias of native students and the fraction of immigrant peers, these findings imply the peak is primarily driven by culturally polarized classrooms. In contrast, fractionalization and non-cultural polarization seem to alleviate in-group biases among native students.

The finding that culturally polarized classrooms drive differences in in-group bias gains further support when we consider two additional classroom scenarios. The first, labeled *overall fractionalization (OF)*, evenly divides the classroom among native peers, Muslim immigrant peers, and non-Muslim immigrant peers, with each group constituting one-third of classmates. The second, denoted as *native supermajority (NSM)*, features native peers as a large majority of 80%, with both Muslim and non-Muslim immigrant peers forming small minority groups of 10% each. In both of these scenarios, the in-group bias of native students is approximately 40% of standard deviation lower compared to culturally polarized classrooms.

### 4.3 Robustness and Extensions

As a first robustness test, we augment regression equation 5. As a reminder, equation 5 is a flexible nesting of both the polarization and fractionalization indices. It is close to a third order expansion of the non-Muslim and Muslim immigrant fractions in a classroom, but does not include the two cubics  $\pi_{C,k}^3$  and  $\pi_{D,k}^3$ . When we add these terms into the regression, the resulting heatmap is similar (see Appendix Figure A2).

As a second robustness check, we use an alternative measure for cultural distance based on linguistic differences between the German language and the country-of-origin language for an immigrant’s parents. Drawing upon language tree data from the Ethnologue database, we calculate the cladistic distance between the language spoken by each immigrant’s origin country and the German language.

---

<sup>16</sup>For context, one standard deviation equals 76 euro cents and the average in-group out-group bias is 9 euro cents, so this represents a 300% increase relative to the mean.

For immigrant children whose parents come from different countries, we determine the weighted cladistic distance based on their mother’s and father’s country-of-origin language. We categorize immigrant peers as culturally distant or close, respectively, based on whether their linguistic distance is above or below the mean.<sup>17</sup> On average, native students have 21% of their classmates categorized as culturally distant (std. dev.=0.17) and 17% as culturally close (std. dev.=0.11). Appendix Figure A3 demonstrates substantial variation in the proportions of linguistically close and distant immigrant peers across different classrooms after netting out school fixed effects. This linguistic measure of cultural diversity has a correlation of 0.94 with the one based on the fraction of Muslims.

In Appendix Figure A4, we re-estimate equation 5 and produce a heatmap using the shares of linguistically close and distant immigrant peers as the main independent variables. This yields qualitatively similar results to the graph using the shares of Muslim and non-Muslim immigrants. Native’s in-group out-group investment gap peaks when the share of native peers is 58% and the share of linguistically distant peers is 42%. Compared to the other classroom scenarios described above, natives’ in-group bias in culturally polarized classrooms is predicted to be 19 to 31% of a standard deviation higher.

As a final robustness check, we expand the number of immigrant groups based on parents’ country of origin.<sup>18</sup> Countries of origin are mapped into 11 regions: Turkey, Balkan States, Eastern Europe, Post Soviet Bloc, Southern Europe, Central and Northern Europe, Middle East, Asia, Africa, other countries, and unidentified (see Appendix Table A2).

Since we now have 11 immigrant groups instead of two, we estimate a more parsimonious model for diversity. Specifically, we use the summary polarization and fractionalization measures in equation 2 as two right-hand side variables. Appendix Figure A5 illustrates that there is independent variation in these two variables in the 11 group case after netting out classroom fixed effects. The figure plots a classroom’s residualized polarization on the x-axis and residualized fractionalization on the y-axis. Holding constant polarization, there can be sizable differences in fractionalization and *vice versa*.

Table 5 contains the in-group out-group investment gap estimates for the 11 immigrant group case. The first column only includes the polarization measure, and finds a large coefficient. The second column only includes the fractionalization measure, and also finds a large coefficient. But when both measures are included in column 3, only the polarization variable matters. The polarization coefficient is large and statistically significant; the estimated effect of 0.66 implies that a one standard deviation increase in polarization (0.153) leads to 10% of a standard deviation increase in in-group bias. In sharp contrast, the fractionalization coefficient is close to zero and insignificant. Of course, this more parsimonious model does not distinguish between different types of polarization.

---

<sup>17</sup>We normalize cladistic distance to be between 0 to 1. The average distance is 0.81. For context, Turkish and Arabic both have values of 1, while French and Polish are 0.69 and 0.56, respectively.

<sup>18</sup>If both parents are immigrants, we use the mother’s origin country; if the mother is an immigrant and the father is not, we use the mother’s country; if the father is an immigrant and the mother is not, we use the father’s country.

## 5 Mechanisms

The prior section found that natives in culturally polarized classrooms have the largest in-group bias, with significantly more bias compared to classrooms characterized by non-cultural polarization, immigrant fractionalization, overall fractionalization, or a native supermajority. In this section, we explore several possible mechanisms which could be driving these results.

As first pointed out by Becker (1957) in the context of labor markets, discrimination can be either profit maximizing or costly, depending on whether it reflects accurate statistical discrimination or taste-based discrimination, respectively. More recently, researchers have argued that inaccurate statistical discrimination (Bohren et al., forthcoming) is also costly, as individuals will make decisions based on incorrect beliefs. In this discussion of mechanisms, we refer to accurate statistical discrimination as “statistical discrimination”, inaccurate statistical discrimination as “negative stereotypes”, and discrimination due to preferences as “taste discrimination”.

In our setting, an important rationale for transferring money in the investment game is *trust*. This refers to the sender’s beliefs about the recipient’s trustworthiness and their expectation of generous reciprocation. Hence, statistical discrimination would be a correct belief among natives that the immigrants they play with are less trustworthy than natives. Negative stereotypes would be an incorrect belief about the trustworthiness of immigrants relative to natives. Taste discrimination would reflect preferences for making transfers to natives versus immigrants, rather than beliefs about trust.

### 5.1 Statistical Discrimination

By design, our experiment rules out the possibility of statistical discrimination playing a role. In our setting, statistical discrimination would be a rational perception among natives that the immigrants they play with will reciprocate less generously than natives. We set up our experiment so that any statistical discrimination would be constant across classroom types. Specifically, students were told they would engage in the investment game with an individual (native or immigrant) chosen randomly from a different school in their federal state of residence, rather than someone from their own class. Since all natives play with the same set of immigrants on average, there is no role for differential statistical discrimination across classrooms with different types of diversity. Thus, it cannot explain the peak in in-group bias observed in culturally polarized classrooms versus other classroom scenarios. Since statistical discrimination is ruled out by design, the remaining explanations are negative stereotypes and taste discrimination, both of which are costly.

### 5.2 Negative Stereotypes

We next turn to negative stereotypes, which in our setting is an inaccurate belief by natives that immigrants are less trustworthy partners in the investment game compared to natives. Our goal is

not to test whether negative stereotypes exist, but rather whether there are differences in negative stereotypes across classrooms which could explain the observed pattern of in-group bias across classrooms.

The idea is that natives in polarized classrooms could rationally develop more mistrust towards the immigrants in their classroom because of limited or adversarial interactions with them. If natives inaccurately generalize their classroom experiences, they could develop negative stereotypes, leading them to place less trust in *all* immigrants compared to natives. Consequently, they would rationally favor fellow natives over immigrants in general, given their biased beliefs.

We start by examining differences using a general measure of trust. Specifically, in our survey, we asked participants how much do you trust people with German nationality and how much do you trust people with a foreign nationality. Students could answer on a scale from 0 to 10, with a 0 indicating “very low degree of trust” and a 10 “very high degree of trust”. These same two questions are asked on the Eurobarometer.

Using the within subject difference in trust (trust in natives - trust in immigrants) as the outcome variable, we estimate equation 5. Figure 6 displays the associated heatmap. The largest in-group out-group trust gap occurs with 54% natives and 46% Muslim immigrant peers. As shown in the graph, this peak closely aligns with the peak in cultural polarization defined in Table 5 and marked in the current figure with the label *CP*. More generally, the two heat plots have very similar patterns.

We find statistically significant increases in the trust gap in culturally polarized classrooms relative to each of the other 4 classroom configurations. The differences are large. For example, the trust gap is 59% of a standard deviation higher in culturally polarized classrooms relative to non-culturally polarized classrooms and 40% higher relative to immigrant fractionalized classrooms.

As a second test for negative stereotypes, we use information from the other side of the experiment, where participants play as the receiver and are asked how much money they would return based on whether the sender was a native or immigrant.

We calculate how natives would have been treated by natives versus immigrants had they been paired with someone from their own classroom. This is a counterfactual exercise, as students were specifically told they would not be paired with a classmate, but rather someone from a different school. But it still serves as a proxy measure for the differential treatment a native is likely to experience in their actual classroom, which they may mistakenly extrapolate to the broader population.

We measure differential treatment by calculating how much, on average, native and immigrant receivers transfer back to native senders. As a reminder, receivers specify how much they would transfer back for each of 11 possible investments (50 cent increments from 0 to 5 euros). Therefore, we average the back transfers for native and immigrant receivers over the 11 possible investments from a native sender. This approach holds constant the amount received from a native sender,

enabling a direct comparison.<sup>19</sup>

We estimate equation 5 using this measure for the native-immigrant gap in reciprocity, and graph the associated heatmap in Figure 7. The largest reciprocity gap occurs with 59% natives and 41% Muslim immigrant peers, which again closely lines up with the cultural polarization peak defined in Table 5. Relative to Figure 5, the heat plot in the current graph has the same general pattern, but with more pronounced differences. We find statistically significant increases in the reciprocity gap in culturally polarized classrooms which are between 123% to 141% of a standard deviation higher compared to non-culturally polarized, immigrant fractionalized, overall fractionalized, and native supermajority classrooms.

### 5.3 Taste Discrimination

A final possible mechanism is taste discrimination, where natives in polarized classrooms simply prefer to make larger transfers to natives relative to immigrants. A taste preference is conceptually different from a belief that natives are more trustworthy and hence will reciprocate more generously. We explore this possibility using two different tests.

We begin by again using information from the other side of the experiment, but this time we focus solely on the decisions natives make when they play as the receiver. Native receivers are asked how much money they would return based on both the amount of the transfer and whether the sender was a native or immigrant. Note that when natives play as the receiver (second mover), they have no financial incentive to return any money to the sender (first mover). This is because they are playing a one-shot game, where the senders and receivers are completely anonymous and come from different schools.

For these reasons, the choices made by a receiver in our game are akin to those made by the second mover in a dictator game. In a one-shot dictator game against an anonymous opponent, a self-interested dictator's optimal strategy is to keep all of the money and return nothing. If they do return money, this can be attributed to other-regarding preferences. The incentives faced by a self-interested receiver in our investment game are similar. And since the receiver conditions their response on whether they are playing against a native or an immigrant, we can measure a native receiver's taste discrimination (i.e., the difference in their other-regarding preferences).

Specifically, we measure taste discrimination by calculating how much, on average, native receivers transfer back to native versus immigrant senders. As a reminder, receivers specify how much they would transfer for each of 11 possible investments (50 cent increments from 0 to 5 euros). Therefore, we average native receivers back transfers over the 11 possible investments, enabling a direct comparison. We then estimate equation 5 using this as the dependent variable.

---

<sup>19</sup>While we also asked senders how much they would expect to get back if they played against an immigrant versus a native, this does not allow for a clean comparison. The reason is that these expectations were asked unconditionally. Hence a native sender could partly expect to get back more from natives versus immigrants because they chose to invest more in natives versus immigrants.

Figure 8 displays the heatmap associated with this measure of taste discrimination. The peak in taste discrimination occurs when natives are 70% of classroom peers and Muslim immigrants are 30%. This maximum is shifted somewhat to the left compared to the peak in cultural polarization for in-group bias observed in Figure 5 and marked by *CP* in the figure. For this reason, while there is still a similar pattern in the heatmap, but shifted somewhat to the left, the differences between culturally polarized classrooms and the four other classroom types are not statistically significant. But we note that the differences are all positive and relatively large (between .20 and .43 of a standard deviation).

As a second test for taste discrimination, we construct a measure of anti-immigrant sentiment among natives. Survey participants were asked whether it is fair that workers of Turkish descent are allowed to work in Germany” and similar questions for immigrants of Polish and French origins. They could answer on a scale of 1 to 4, where 1=strongly agree, 2=agree somewhat, 3=disagree somewhat, and 4=strongly disagree. Using a principal components analysis, we construct an index of anti-immigrant sentiment, normalizing the index to have mean 0 and standard deviation 1. A larger value indicates a native has more anti-immigrant sentiment. Note that this outcome is somewhat different from all of the other outcomes used in the paper, as it is measured in levels rather than as a difference.

The heatmap for anti-immigrant sentiment is shown in Figure 9. The peak occurs close to the peak in cultural polarization for in-group bias, as marked with the label *CP* in the graph. And as we have seen in the many heatmaps preceding this one, the pattern is similar to that observed in Figure 5. The difference in anti-immigrant sentiment for a culturally polarized classroom versus non-culturally polarized classroom is 44% of a standard deviation. There are likewise large and statistically significant differences relative to fractionalized immigrant, overall fractionalized, and native supermajority classrooms. This provides support for taste discrimination contributing to the in-group bias observed in Figure 5.

We recognize that taste discrimination and negative stereotypes could interact with each other. In particular, an increase in negative stereotypes could affect an individual’s preferences, amplifying their taste discrimination. And taste discrimination could contribute to negative stereotypes.

## 5.4 Payoff Losses

In general, discrimination can lead to either higher or lower payoffs, depending on the form it takes. But in our setting, the profit maximizing force of statistical discrimination is ruled out by design. The implication is that natives in culturally polarized classrooms should experience payoff losses due to their in-group bias. The reason is that the remaining explanations of negative stereotypes and taste discrimination, both of which we find empirical evidence for, are costly.

In Figure 10, we test this prediction about payoff losses. For each native, we calculate the expected difference in their payoffs as the sender (first mover) if they are randomly matched with an immigrant

versus a native interaction partner (second mover).<sup>20</sup> Using the expected payoff difference as the outcome variable, we estimate equation 5.

The resulting heatmap is displayed in Figure 10. The payoff loss reaches a peak when there are 54% natives and 46% Muslim immigrants in a classroom. In comparison, the peak in in-group bias shown in Figure 5 occurs at a similar point (58% natives and 42% Muslim peers). More generally, comparing the two heat plots, they have almost identical patterns. To test this more formally, we compare the predicted losses in culturally polarized classrooms versus the four other classroom types. Importantly, classroom types are defined based on the definitions used in Figure 5; we do not redefine what a culturally polarized classroom is to reflect the peak in payoff losses in Figure 10.

We find larger, and statistically significant, payoff losses in culturally polarized classrooms relative to other classroom types. Payoff losses are 0.37 standard deviations higher compared to non-culturally polarized classrooms. Likewise, payoff losses are 0.33 to 0.46 standard deviations higher relative to classrooms with immigrant fractionalization, overall fractionalization, or native supermajorities. These findings confirm the prediction that discrimination is costly to natives in our setting.

## 6 A Framework

In this section, we posit a simple model that provides a basis for rationalizing our empirical findings. Drawing on the idea that in-group bias results from limited meaningful intergroup interaction, it features individuals who deliberately choose who to identify and socially engage with. Limited intergroup interaction would naturally explain the negative stereotypes and taste discrimination we observe. Individuals are of different types—e.g., natives and multiple immigrant groups—and these types differ in terms of their endowed cultural identity. They must decide whether to join an *inclusive* group with a shared identity or preserve their innate cultural identity by joining an *exclusive* group comprised solely of others of their own type. This decision involves striking a balance between the wish to be part of a larger group and the willingness to forfeit one’s own cultural identity in favor of that of the inclusive group, where the latter shifts endogenously with the distribution of types in the population. Individual heterogeneity is introduced through additive extreme value-distributed preferences for inclusive versus exclusive group membership. Given that each individual faces a binary decision and cares about how many others make the same choice—that is, join the same group—our model can be seen as a version of the social interactions model by Brock and Durlauf (2001), but with an in-built Hotelling structure where multiple interacting types are endowed with different cultural identities. In our framework, natives’ attitudes toward immigrants naturally correspond to the proportion of natives who join and identify with the inclusive group.

---

<sup>20</sup>As a reminder, all individuals fill out decisions sheets as the first mover playing with a native and with an immigrant as well decision sheets as the second mover playing with a native and with an immigrant. Participants are told they will be randomly chosen to play either as the first mover or the second mover and that they will be paired randomly with a native or an immigrant.

**Basic Structure.** Consider an economy with a continuum of individuals who are of  $J \geq 2$  types. Let the proportions of types in the population be denoted  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_J)$ . There is a space of cultural identities which we take to be  $\Theta \equiv [-1, 1]^{J-1}$ . Each individual  $i$  of type  $j$  is endowed with some exogenously given identity  $\boldsymbol{\theta}_j \in \Theta$  and the inclusive group has an endogenous identity  $\boldsymbol{\theta} \in \Theta$ , making the type-specific (Euclidean) distance  $d_j \equiv \|\boldsymbol{\theta} - \boldsymbol{\theta}_j\|$ . We assume a type-specific cost of joining the inclusive group that depends on this cultural distance,  $h_j = h_j(d_j)$  where, for each type  $j$ ,  $h_j(\cdot)$  is twice continuously differentiable and satisfies  $h_j(0) = 0$ ,  $h_j'(\cdot) > 0$  and  $h_j''(\cdot) > 0$ . Finally, let  $\beta > 0$  parameterize the strength of preference for group size. Individual  $i$  of type  $j$  chooses between two options: either stay with the type- $j$  exclusive group (option 0) or join the inclusive group (option 1). Let  $\mu_j \in [0, 1]$  denote the proportion of type  $j$  who choose to join the inclusive group. The associated utilities for individual  $i$  are then

$$u_{ij}^0 = \beta\pi_j(1 - \mu_j) + \varepsilon_{ij}^0, \quad u_{ij}^1 = \beta \sum_{j'=1}^J \pi_{j'}\mu_{j'} - h_j + \varepsilon_{ij}^1, \quad j = 1, \dots, J$$

where we used that the size of the type- $j$  exclusive group is  $\pi_j(1 - \mu_j)$  and size of the inclusive group is  $\sum_{j'=1}^J \pi_{j'}\mu_{j'}$ , and where  $\varepsilon_{ij}^0$  and  $\varepsilon_{ij}^1$  are the individual's i.i.d. extreme value distributed choice-specific random preferences. For any given type distribution  $\boldsymbol{\pi}$  and inclusive group identity  $\boldsymbol{\theta}$  (and hence fixed vector of joining costs  $\mathbf{h} = (h_1, \dots, h_J)$ ), a “joining equilibrium” is a vector of type-specific inclusive-group joining rates  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_J)$  that satisfies

$$\log\left(\frac{\mu_j}{1 - \mu_j}\right) = \beta \left( \sum_{j'=1}^J \pi_{j'}\mu_{j'} - \pi_j(1 - \mu_j) \right) - h_j, \quad j = 1, \dots, J$$

for all  $J$  types simultaneously. In Appendix B, we prove that such a joining equilibrium exists and is guaranteed to be unique for  $\beta \in (0, 2)$  for any finite number of types  $J$ . The upper limit,  $\beta \leq 2$ , ensures the absence of multiple coordination equilibria.

Endogeneity of the inclusive group's identity  $\boldsymbol{\theta}$ , and hence of the type-specific distance costs, plays a central role. The interactions within the inclusive group occurs with a shared identity, for instance by mixing preferred activities or by adopting a common clothing style or slang. To close the model we need to specify how  $\boldsymbol{\theta}$  is determined. In line with the assumption of positive preferences for group size, we assume that  $\boldsymbol{\theta}$  is chosen by the inclusive group so as to maximize the group's size (or “popularity”) in the joining equilibrium that ensues,

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{j=1}^J \pi_j \mu_j.$$

Two factors will make the inclusive group's identity  $\boldsymbol{\theta}$  align relatively closely with the endowed identity  $\boldsymbol{\theta}_j$  of type  $j$ . First, if type  $j$ 's has a *relatively large population share*  $\pi_j$  (as this makes type  $j$

an important potential source of members for the inclusive group). Second, if type  $j$  has a *relatively strong aversion to cultural distance* (as placing a  $\theta$  close to  $\theta_j$  will then encourage participation by type  $j$  individuals).

**A Three-Type Case.** A three-type version of the model—i.e., one featuring natives and two immigrant groups—can be used to revisit our main empirical findings. Thus, consider the  $J = 3$  case where  $\Theta = [-1, 1]^2$ . Let the three types be denoted by  $N$ ,  $C$ , and  $D$ , for “natives”, culturally “close” immigrants, and culturally “distant” immigrants, respectively. We assume that their endowed identities are represented by points on the unit circle:  $\theta_N = (-1, 0)$ ,  $\theta_C = (-1/2, \sqrt{3}/2)$  and  $\theta_D = (1, 0)$ . The left panel of Figure 11 illustrates  $\Theta$  and the endowed identities. As is evident, the cultural identity of type- $C$  immigrants is close to that of natives, whereas type- $D$  immigrants exhibit a large cultural distance from the two other groups.

For comparative statics purposes, we use two parameters to represent the population type distribution:  $\pi \in [0, 1]$  denotes the overall immigrant share and  $s \in [0, 1]$  the share of culturally distant immigrants among the immigrants. Hence  $\pi_N = 1 - \pi$ ,  $\pi_C = \pi(1 - s)$  and  $\pi_D = \pi s$ . Our interest lies in how natives’ inclusive behavior—i.e., their propensity to engage with immigrants in the inclusive group ( $\mu_N$ )—varies with changes in the population type distribution ( $\pi, s$ ). To relate the model to our heatmaps discussed earlier, we specifically focus on  $(\pi, s)$  combinations that represent cultural polarization ( $CP$ ), non-cultural polarization ( $NCP$ ), fractionalization ( $IF$  or  $OF$ ), and a native supermajority ( $NSM$ ).

Figure 11 illustrates our two main comparative statics. They are derived from a numerical example constructed to capture the qualitative features of our empirical findings using minimal parameterization. In the example, we set  $\beta = 1.75$ ,  $h_j(d) = \gamma_j d^\sigma$  with a common elasticity  $\sigma = 1.1$  but with type-specific constant terms:  $\gamma_N = 0.35$  and  $\gamma_C = \gamma_D = 0.6$ . Hence, in the example, natives have a lower aversion to cultural distance than culturally close and distant immigrants.

The first comparative static question we ask is: How does altering the mix of immigrants ( $s$ ) influence the inclusive behavior of natives ( $\mu_n$ ), all while keeping the fraction of natives ( $\pi$ ) constant at 50 percent? In context of our heatmaps, this exercise can be seen as capturing the curvature of the line connecting the points ( $NCP, IF, CP$ ); that is, it enables a comparison between cultural polarization ( $\pi = \frac{1}{2}, s = 1$ ) and non-cultural polarization ( $\pi = \frac{1}{2}, s = 0$ ), as well as immigrant fractionalization ( $\pi = \frac{1}{2}, s = \frac{1}{2}$ ). The solid green lines in Figure 11 illustrate the results. Panel (a) shows the cultural identity of the inclusive group ( $\theta$ ), while panel (b) depicts natives’ propensity to mix with immigrants ( $\mu_N$ ). When  $s = 0$ , the only two types in the population are natives and culturally close immigrants ( $NCP$ ). The inclusive group’s identity  $\theta$  will be a point on the chord between  $\theta_C$  and  $\theta_N$ , in this case closer to  $\theta_C$  than to  $\theta_N$  reflecting that  $\gamma_C > \gamma_N$ . Conversely, when  $s = 1$ , only natives and culturally distant immigrants exist ( $CP$ ) and  $\theta$  will be a point on the chord between  $\theta_D$  and  $\theta_N$  but closer to the former reflecting that  $\gamma_D > \gamma_N$ .

An increase in  $s$  will have both a direct and an indirect effect on the natives’ inclusive behavior.

A *direct* effect of  $s$  obtains via the impact on the overall immigrant joining rate,  $(1 - s)\mu_C + s\mu_D$ , holding  $\theta$  constant: if  $s$  increases the overall immigrant joining rate, this will encourage joining also by natives. However, this direct effect will be muted, as it involves two offsetting forces. On the one hand, an increase in  $s$  tends to increase  $\mu_C$  by reducing the own-type exclusive group for type-C immigrants. On the other hand, it decreases  $\mu_D$  by increasing the own-type exclusive group for type-D immigrants. An *indirect* effect of  $s$  obtains via a shifting of the inclusive group's identity  $\theta$ . As  $s$  increases away from zero towards 0.5, immigrants become increasingly fractionalized, leaving the natives as the single largest population type. To maximize the inclusive group's size,  $\theta$  will move closer to the native culture, explaining why the green solid line in the left panel is curved towards  $\theta_N$ . As a result, natives' propensity to mix with immigrants in the inclusive group,  $\mu_N$ , initially increases in  $s$ . However, as the proportion of distant immigrants continues to rise, the identity of the inclusive groups begins to shift away from that of natives and leans towards that of distant immigrants. Consequently, a larger number of natives start to exclusively engage with in-group members. This effect is amplified by a social multiplier, so that even small increases in  $s$  result in large drops in  $\mu_n$ . At the endpoint, where the fraction of distant immigrants is unity ( $s = 1$ ), natives' propensity to mix with immigrants reaches its lowest level. Overall, our first comparative static result rationalizes the lower inclusive behavior of natives under cultural polarization compared to scenarios involving non-cultural polarization and immigrant fractionalization. More specifically, it mirrors the curvature of the line that links the points ( $NCP, IF, CP$ ) in Figure 5.

We now turn to our second comparative static result, namely how an increase in the share of immigrants ( $\pi$ ), with a constant immigrant mix ( $s$ ) of 50 percent, affects natives' inclusive behavior ( $\mu_n$ ). This analysis can be seen as corresponding to the curvature of the line connecting the points ( $NSM, IF, OF$ ) in our heatmaps; that is, it allows for a comparison between scenarios involving a native supermajority ( $\pi = \frac{1}{3}, s = \frac{1}{2}$ ), immigrant fractionalization ( $\pi = \frac{1}{2}, s = \frac{1}{2}$ ), and overall fractionalization ( $\pi = \frac{1}{3}, s = \frac{1}{2}$ ).

As a recap, the heatmap depicted in Figure 5 revealed minimal variation in in-group bias among natives across the scenarios of ( $NSM, IF, OF$ ). The hatched blue lines in Figure 11 demonstrate that our model predicts this pattern.<sup>21</sup> At low values of  $\pi$ , natives constitute the dominant majority ( $NSM$ ), and the inclusive group's identity,  $\theta$ , naturally aligns closely with the natives' endowed identity,  $\theta_N$ . As a result, the cultural distance cost is low for natives, causing many of them to identify with immigrants in the inclusive group (reflected by high  $\mu_N$ ). As  $\pi$  increases, two counteracting forces influence the inclusive behavior of natives. On the one hand, the identity of the inclusive group ( $\theta$ ) shifts away from that of natives ( $\theta_N$ ) as they begin to lose their majority status. This discourages natives from engaging with immigrants. On the other hand, the native own-type group shrinks as  $\pi$  increases, encouraging natives to identify with immigrants. The latter effect weakly dominates for relatively low values of  $\pi$ , but as  $\pi$  continues to increase, the two effects almost completely offset each other. On the whole, and consistent with the heatmap depicted in

<sup>21</sup>We restrict attention in the example to  $\pi \geq 0.15$  as the binary choice model with logistic errors is not a suitable representation in an environment where one type forms an overwhelming majority in the population.

Figure 5, the model predicts that natives exhibit similar behavioral patterns in scenarios involving a native supermajority, immigrant fractionalization, and overall fractionalization.

When considering our two comparative static results in tandem, cultural polarization emerges as the environment least favorable for fostering inclusive behavior among natives. Thus, our framework can explain why we observe a peak in natives' in-group bias in culturally polarized classroom.<sup>22</sup> It is, however, worth noting that alternative economic models could be adapted to account for our findings. Indeed, our approach is closely linked to existing literature on (i) cultural identity and assimilation and (ii) the endogenous formation of friendships and homophily. Classic contributions to the former literature include Lazear (1999), Bisin and Verdier (2000), and Carvalho (2013). In the latter, Currarini et al. (2009) stands out for constructing a random matching model to investigate homophily in the creation of friendship networks, highlighting that the most significant in-group bias emerges from middle-sized groups.

## 7 Conclusion

Our lab-in-the field experiment reveals that in-group bias among natives peaks in culturally polarized classrooms, where German natives form a slim majority and Muslim immigrants a large minority. In contrast, in classrooms characterized by non-cultural polarization, overall fractionalization, immigrant fractionalization, or a native majority, there are significantly lower levels of in-group bias. The estimated gaps for payoff losses, negative stereotypes, and taste discrimination follow a remarkably similar pattern, supplying evidence on mechanisms.

These findings provide a fresh perspective on the impact of migration induced-diversity, reconciling evidence from prior studies. It has long been argued that diversity in schools, workplaces, and neighborhoods can serve as a catalyst for close intergroup interaction, thereby helping to dismantle negative stereotypes and encourage deeper cross-cultural understanding (Allport, 1954; Pettigrew and Tropp, 2008). However, recent empirical studies have found mixed results, with migration-induced diversity either improving cross-cultural relations and understanding (Calderon et al., 2023) or failing to foster meaningful contact and resulting in hostility instead (Algan et al., 2016). What makes diversity unifying in some settings but divisive in others? Our research points to two key factors: (i) whether diversity takes the form of polarization or fractionalization and (ii) the cultural distance between groups.

From a policy perspective, our results suggest that extra efforts are needed to counteract low levels of trust in culturally polarized environments. To address this in schools, potential solutions could involve modifying the curricula to incorporate lessons focused on inclusion (Alan et al., 2021)

---

<sup>22</sup>In Appendix B, we present comparative static results when there are natives and just one immigrant type. Specifically, we investigate the two limiting cases where we manipulate the immigrant share ( $\pi$ ) with  $s = 0$  and  $s = 1$ , respectively. In the context of our heatmaps, they correspond to traversing along a vertical ray ( $s = 0$ ) and a horizontal ray ( $s = 1$ ) from the origin, respectively. These comparative statics also provide an explanation for the inverted U-shape observed in-group bias when varying the immigrant share in Figure 4.

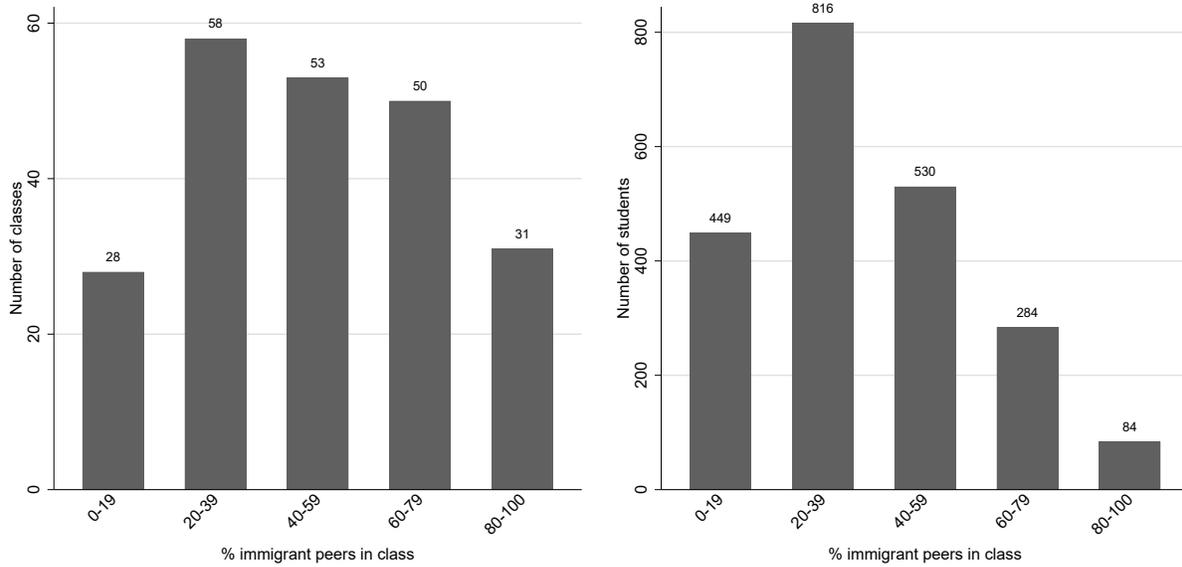
or implementing randomized seating assignments to disrupt own-group attachment and promote cross-group friendships (Faur and Laursen, 2022). Our findings also have wider implications. For example, governments use a variety of assignment rules to allocate refugees and immigrants across regions, and researchers have recently developed algorithms to improve these practices (Bansak et al., 2018). Our findings suggest it could be important to factor in whether the destination communities will become culturally polarized.

## References

- Alan, Sule, Ceren Baysan, Mert Gumren, and Elif Kubilay (2021) “Building Social Cohesion in Ethnically Mixed Schools: An Intervention on Perspective Taking,” *The Quarterly Journal of Economics*, 136 (4), 2147–2194.
- Alesina, Alberto and Eliana La Ferrara (2000) “Participation in Heterogeneous Communities,” *Quarterly Journal of Economics*, 115 (3), 847–904.
- (2002) “Who Trusts Others?,” *Journal of Public Economics*, 85 (2), 207–234.
- Alesina, Alberto and Marco Tabellini (forthcoming) “The Political Effects of Immigration: Culture or Economics?” *Journal of Economic Literature*.
- Algan, Yann, Camille Hémet, and David D Laitin (2016) “The Social Effects of Ethnic Diversity at the Local Level: A Natural Experiment with Exogenous Residential Allocation,” *Journal of Political Economy*, 124 (3), 696–733.
- Allport, Gordon (1954) *The Nature of Prejudice*: Cambridge, MA: Addison-Wesley Publishing.
- Antecol, Heather, Ozkan Eren, and Serkan Ozbeklik (2015) “The Effect of Teacher Gender on Student Achievement in Primary School,” *Journal of Labor Economics*, 33 (1), 63–89.
- Balestra, Simone, Beatrix Eugster, and Helge Liebert (2022) “Peers with Special Needs: Effects and Policies,” *Review of Economics and Statistics*, 104 (3), 602–618.
- Bansak, Kirk, Jeremy Ferwerda, Jens Hainmueller, Andrea Dillon, Dominik Hangartner, Duncan Lawrence, and Jeremy Weinstein (2018) “Improving Refugee Integration through Data-Driven Algorithmic Assignment,” *Science*, 359 (6373), 325–329.
- Bazzi, Samuel, Arya Gaduh, Alexander D. Rothenberg, and Maisy Wong (2019) “Unity in Diversity? How Intergroup Contact Can Foster Nation Building,” *American Economic Review*, 109 (11), 3978–4025.
- Becker, Gary (1957) *The Economics of Discrimination*: University of Chicago Press, Chicago.
- Bellemare, Charles and Sabine Kröger (2007) “On Representative Social Capital,” *European Economic Review*, 51 (1), 183–202.
- Bifulco, Robert, Jason M. Fletcher, and Stephen L. Ross (2011) “The Effect of Classmate Characteristics on Post-secondary Outcomes: Evidence from the Add Health,” *American Economic Journal: Economic Policy*, 3 (1), 25–53.
- Bisin, Alberto and Thierry Verdier (2000) ““Beyond the Melting Pot”: Cultural Transmission, Marriage, and the Evolution of Ethnic and Religious Traits,” *Quarterly Journal of Economics*, 115 (3), 955–988.
- Bohren, J. Aislinn, Kareem Haggag, Alex Imas, and Devin G. Pope (forthcoming) “Inaccurate Statistical Discrimination: An Identification Problem,” *Review of Economics and Statistics*.
- Boucher, Vincent, Semih Tumen, Michael Vlassopoulos, Jackline Wahba, and Yves Zenou (2022) “Ethnic Mixing in Early Childhood: Evidence from a Randomized Field Experiment and a Structural Model,” *Unpublished Manuscript*.
- Brenøe, Anne Ardila and Ulf Zölitz (2020) “Exposure to More Female Peers Widens the Gender Gap in STEM Participation,” *Journal of Labor Economics*, 38 (4), 1009–1054.
- Brock, William A. and Steven N. Durlauf (2001) “Discrete Choice with Social Interactions,” *Review of Economic Studies*, 68 (2), 235–260.
- Brown, Bradford B. (2011) “Peer Groups and Peer Cultures,” in Feldman, S. S. and G. R. Elliott eds. *At the Threshold: The Developing Adolescent*, 171–196: Harvard University Press.
- Calderon, Alvaro, Vasiliki Fouka, and Marco Tabellini (2023) “Racial Diversity and Racial Policy Preferences: The Great Migration and Civil Rights,” *Review of Economic Studies*, 90 (1), 165–200.
- Carrell, Scott E, Mark Hoekstra, and Elira Kuka (2018) “The Long-Run Effects of Disruptive Peers,” *American Economic Review*, 108 (11), 3377–3415.
- Carrell, Scott E, Mark Hoekstra, and James E West (2019) “The Impact of College Diversity on Behavior Toward Minorities,” *American Economic Journal: Economic Policy*, 11 (4), 159–182.
- Carrell, Scott E and James E West (2010) “Does Professor Quality Matter? Evidence from Random

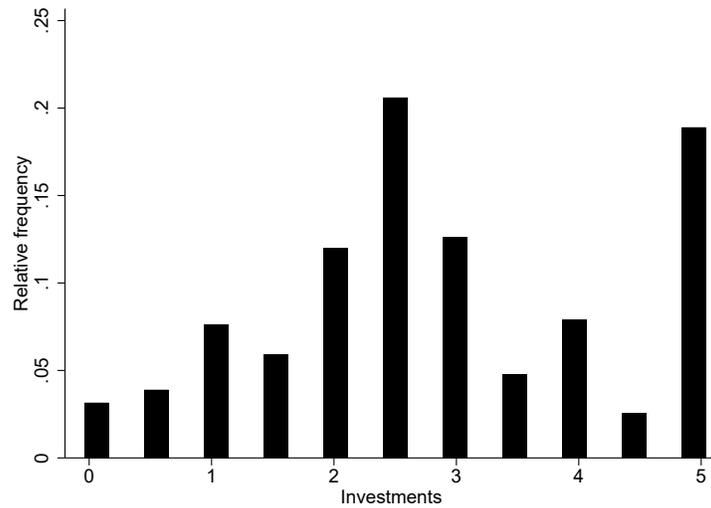
- Assignment of Students to Professors,” *Journal of Political Economy*, 118 (3), 409–432.
- Carrell, Scott and Mark Hoekstra, “Externalities in the Classroom: How Children Exposed to Domestic Violence Affect Everyone’s Kids,” *American Economic Journal: Applied Economics*, 2 (1), 2010.
- Carvalho, Jean-Paul (2013) “Veiling,” *Quarterly Journal of Economics*, 128 (1), 337–370.
- Chetty, Raj, John Friedman, Nathaniel Hilger, Emmanuel Saez, Diane Whitmore Schanzenbach, and Danny Yagan (2011) “How Does Your Kindergarten Classroom Affect Your Earnings? Evidence from Project Star,” *Quarterly Journal of Economics*, 126 (4), 1593–1660.
- Corno, Lucia, Eliana La Ferrara, and Justine Burns (2022) “Interaction, Stereotypes, and Performance: Evidence from South Africa,” *American Economic Review*, 112 (12), 3848–75.
- Currarini, Sergio, Matthew O Jackson, and Paolo Pin (2009) “An Economic Model of Friendship: Homophily, Minorities, and Segregation,” *Econometrica*, 77 (4), 1003–1045.
- Dahl, Gordon B, Christina Felfe, Paul Frijters, and Helmut Rainer (2022) “Caught between Cultures: Unintended Consequences of Improving Opportunity for Immigrant Girls,” *Review of Economic Studies*, 89 (5), 2491–2528.
- Dahl, Gordon, Andreas Kotsadam, and Dan-Olof Rooth (2021) “Does Integration Change Gender Attitudes? The Effect of Randomly Assigning Women to Traditionally Male Teams,” *Quarterly Journal of Economics*, 136 (2), 987–1030.
- Dee, Thomas S. (2004) “Are There Civic Returns to Education?” *Journal of Public Economics*, 88 (9), 1697–1720.
- Falk, Armin and Christian Zehnder (2013) “A City-Wide Experiment on Trust Discrimination,” *Journal of Public Economics*, 100, 15–27.
- Faur, Sharon and Brett Laursen (2022) “Classroom Seat Proximity Predicts Friendship Formation,” *Frontiers in Psychology*, 13, 796002.
- Federal Anti-Discrimination Agency (2022) *Diskriminierung an Schulen erkennen und vermeiden: Praxiseitfaden zum Abbau von Diskriminierung an Schulen*: Antidiskriminierungsstelle des Bundes.
- Federal Government of North Rhine-Westphalia (2022) *11. Kinder- und Jugendbericht der Landesregierung Nordrhein-Westfalen*: Ministerium für Kinder, Familie, Flüchtlinge und Integration des Landes Nordrhein-Westfalen.
- Federal Statistical Office of Germany (2009) *Bevölkerung und Erwerbstätigkeit. Bevölkerung mit Migrationshintergrund - Ergebnisse des Mikrozensus 2018*: Statistische Bundesamt.
- Feld, Jan and Ulf Zölitz (2017) “Understanding Peer Effects: On the Nature, Estimation, and Channels of Peer Effects,” *Journal of Labor Economics*, 35 (2), 387–428.
- Felfe, Christina, Martin G Kocher, Helmut Rainer, Judith Saurer, and Thomas Siedler (2021) “More Opportunity, More Cooperation? The Behavioral Effects of Birthright Citizenship on Immigrant Youth,” *Journal of Public Economics*, 200, 104448.
- Fershtman, Chaim and Uri Gneezy (2001) “Discrimination in a Segmented Society: An Experimental Approach,” *Quarterly Journal of Economics*, 116 (1), 351–377.
- Figlio, David N, Paola Giuliano, Riccardo Marchingiglio, Umut Åzek, and Paola Sapienza (forthcoming) “Diversity in Schools: Immigrants and the Educational Performance of U.S. Born Students,” *Review of Economic Studies*.
- Hanson, Gordon and Craig McIntosh (2016) “Is the Mediterranean the New Rio Grande? US and EU Immigration Pressures in the Long Run,” *Journal of Economic Perspectives*, 30 (4), 57–82.
- Hanushek, Eric, John Kain, Jacob Markman, and Steven Rivkin (2003) “Does Peer Ability Affect Student Achievement?,” *Journal of Applied Econometrics*, 18 (5), 527–544.
- Horowitz, Donald L (1985) *Ethnic Groups in Conflict*: University of California Press.
- Hotelling, Harold (1929) “Stability in Competition,” *Economic Journal*, 39 (153), 41–57.
- Hoxby, Caroline (2002) “The Power of Peers: How Does the Makeup of a Classroom Influence Achievement,” *Education Next*, 20 (2), 56–63.
- Kremer, Michael and Dan Levy (2008) “Peer Effects and Alcohol Use among College Students,” *Journal of Economic Perspectives*, 22 (3), 189–206.

- Lavy, Victor, M Daniele Paserman, and Analia Schlosser (2012) “Inside the Black Box of Ability Peer Effects: Evidence from Variation in the Proportion of Low Achievers in the Classroom,” *Economic Journal*, 122 (559), 208–237.
- Lazear, Edward P (1999) “Culture and Language,” *Journal of Political Economy*, 107 (S6), S95–S126.
- Lowe, Matt (2021) “Types of Contact: A Field Experiment on Collaborative and Adversarial Caste Integration,” *American Economic Review*, 111 (6), 1807–1844.
- Luttmer, Erzo FP (2001) “Group Loyalty and the Taste for Redistribution,” *Journal of Political Economy*, 109 (3), 500–528.
- Lyle, David (2007) “Estimating and Interpreting Peer and Role Model Effects from Randomly Assigned Social Groups at West Point,” *Review of Economics and Statistics*, 89 (2), 289–299.
- Manski, Charles (1993) “Identification of Endogenous Social Effects: The Reflection Problem,” *Review of Economic Studies*, 60 (31), 531–542.
- McAuliffe, M. and A. Triandafyllidou (2021) *World Migration Report 2022: United Nations International Organization for Migration*, Geneva.
- Milligan, Kevin, Enrico Moretti, and Philip Oreopoulos (2004) “Does Education Improve Citizenship? Evidence from the United States and the United Kingdom,” *Journal of Public Economics*, 88 (9), 1667–1695.
- Montalvo, José G and Marta Reynal-Querol (2005) “Ethnic Polarization, Potential Conflict, and Civil Wars,” *American Economic Review*, 95 (3), 796–816.
- Mousa, Salma (2020) “Building Social Cohesion between Christians and Muslims through Soccer in Post-ISIS Iraq,” *Science*, 369 (6505), 866–870.
- Nieden, Birgit zur and Juliane Karakayali (2016) “Harte Tür: Schulische Segregation nach Herkunft in der postmigrantischen Gesellschaft,” 81–96.
- Ohinata, Asako and Jan C Van Ours (2013) “How Immigrant Children Affect the Academic Achievement of Native Dutch Children,” *Economic Journal*, 123 (570), F308–F331.
- Pettigrew, Thomas F and Linda R Tropp (2008) “How Does Intergroup Contact Reduce Prejudice? Meta-Analytic Tests of Three Mediators,” *European Journal of Social Psychology*, 38 (6), 922–934.
- Puskeppeleit, Jürgen and Marianne Krüger-Potratz (1999) *Bildungspolitik und Migration: Texte und Dokumente zur Beschulung ausländischer und ausgesiedelter Kinder und Jugendlicher: 1950-1999: Bd 1. Arbeitsstelle Interkulturelle Pädagogik*.
- Rao, Gautam (2019) “Familiarity Does Not Breed Contempt: Generosity, Discrimination, and Diversity in Delhi Schools,” *American Economic Review*, 109 (3), 774–809.
- Sacerdote, Bruce (2001) “Peer Effects with Random Assignment: Results for Dartmouth Roommates,” *Quarterly Journal of Economics*, 116 (2), 681–704.
- (2011) “Peer Effects in Education: How Might They Work, How Big Are They and How Much Do We Know Thus Far?” in Hanushek, Eric A., Stephen Machin, and Ludger Woessmann eds. *Handbook of the Economics of Education*, 3, 249–277: Elsevier.
- Tropp, Linda R and Suchi Saxena (2018) “Re-Weaving the Social Fabric Through Integrated Schools: How Intergroup Contact Prepares Youth to Thrive in a Multiracial Society,” *Research Brief No. 13.*, *National Coalition on School Diversity*.
- Wells, Amy Stuart, Lauren Fox, and Diana Cordova-Cobo (2016) “How Racially Diverse Schools and Classrooms Can Benefit all Students,” *The Education Digest*, 82 (1), 17.



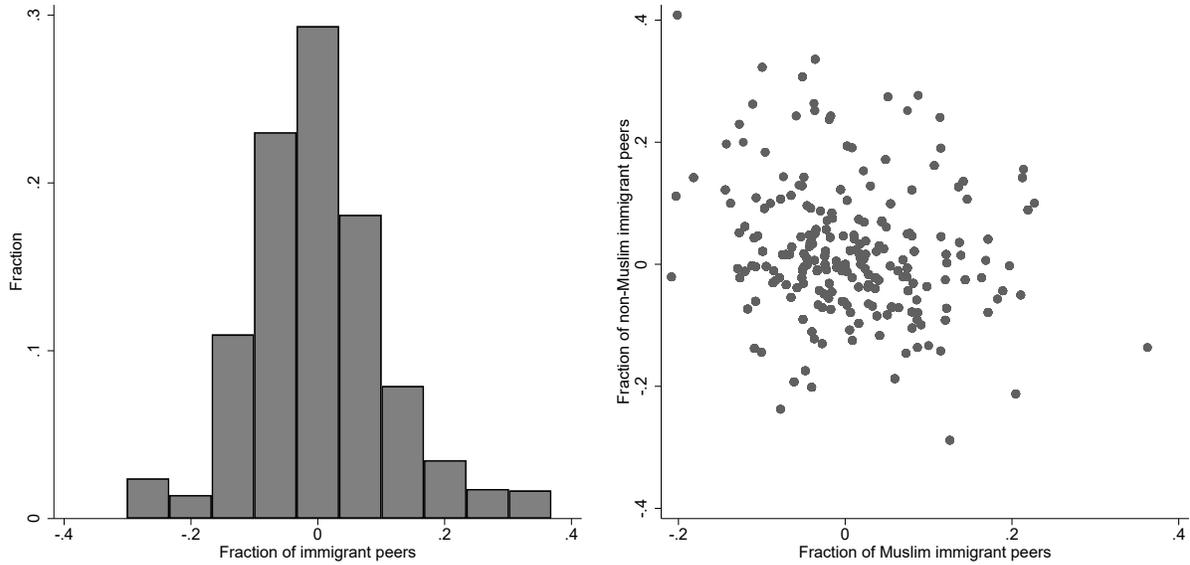
**Figure 1.** Share of Immigrant Peers at the Classroom and Student Level

*Notes:* There are 220 classes in the left panel and 2,163 native German students in the right panel.



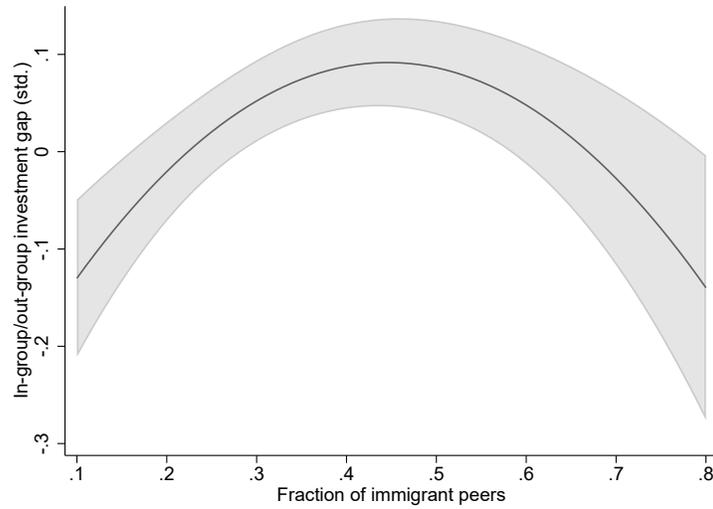
**Figure 2.** Investment Decisions Made by Native Senders

*Notes:* Histogram of investment decisions made by native senders. Senders could invest between 0 and 5 euros, in 50 cent increments.



**Figure 3.** Identifying Variation after Removing School Fixed Effects

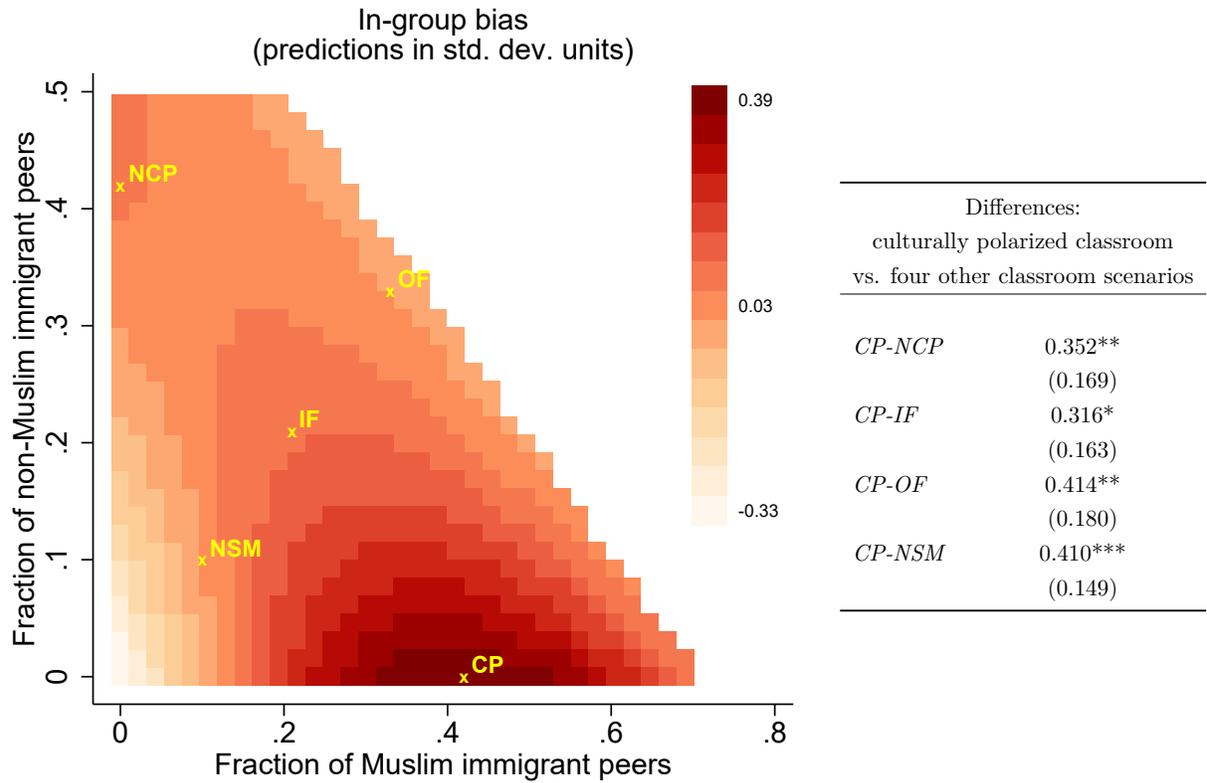
*Notes:* The left panel is a histogram of the residualized immigrant share, net of school fixed effects. The right panel is a scatter plot of the residualized non-Muslim immigrant share against the residualized Muslim immigrant share, both net of school fixed effects



**Figure 4.** In-Group Bias: An Inverted-U in the Share of Immigrant Peers

*Notes:* The figure uses estimates from column 3 of Table 3 to predict the in-group/out-group investment gap for immigrant shares, evaluating control variables at their means. The grey shaded area denotes pointwise 95% confidence intervals.

N=2,163. 1 standard deviation = 76 euro-cents.



**Figure 5.** In-Group Bias: How Type of Diversity Matters

*Notes:* In-group bias is measured by the in-group/out-group investment gap in the first stage of the investment game. The heatmap shows predicted values for in-group bias based on estimates for equation (5). The peak in in-group bias occurs in a classroom exhibiting **cultural polarization (CP)**, with native German peers forming a slight majority group [58%] and Muslim immigrant peers a large minority group [42%]. The table shows the difference between predicted in-group bias at this peak and four other classroom scenarios:

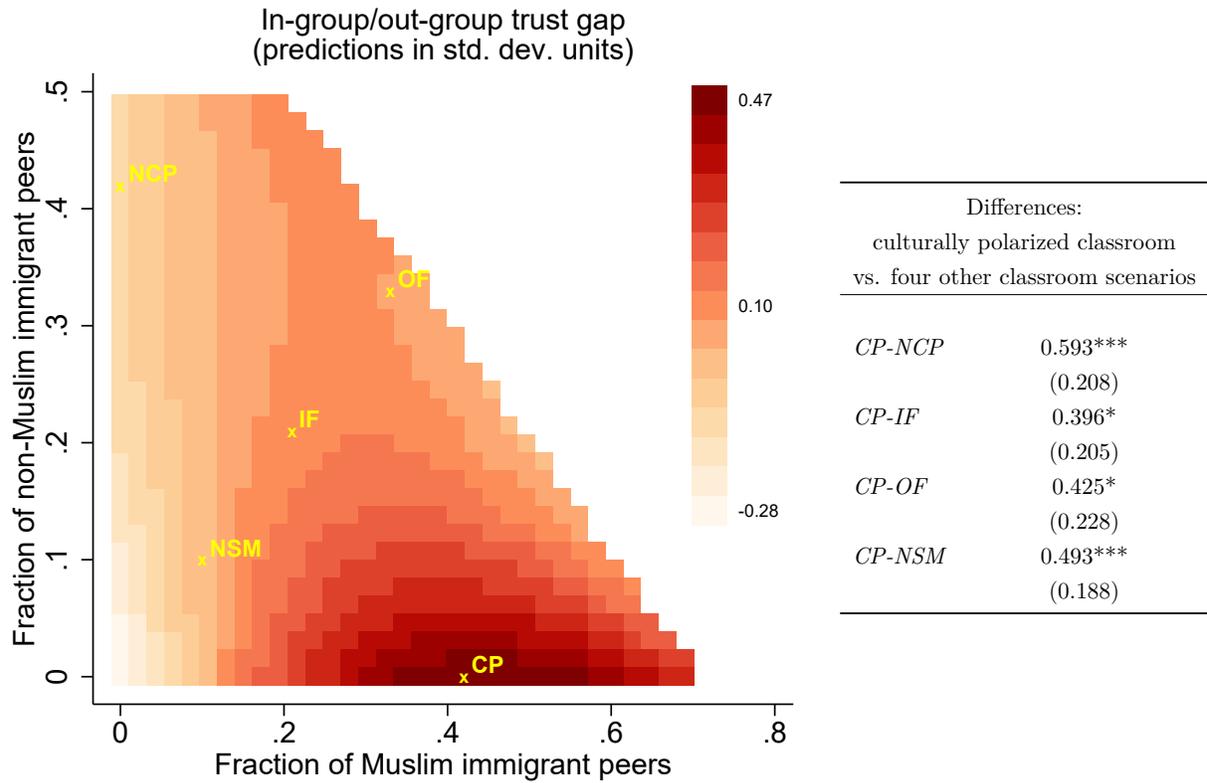
**Non-cultural polarization (NCP):** native German peers form a slight majority group [58%] and non-Muslim immigrant peers a large minority group [42%]

**Immigrant fractionalization (IF):** native peers form a slight majority group [58%] and both Muslim and non-Muslim immigrant peers equally large, medium-sized minority groups [21% each]

**Overall fractionalization (OF):** native peers, Muslim immigrant peers, and non-Muslim immigrant peers form equally large groups [33% each]

**Native supermajority (NSM):** native peers form a large majority group [80%] and both Muslim and non-Muslim immigrant peers small minority groups [10% each]

N=2,163. 1 standard deviation = 76 euro-cents. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .



**Figure 6.** Negative Stereotypes: Biased Beliefs about Immigrants' Trustworthiness

*Notes:* The in-group/out-group trust gap is measured through survey questions asking students how much they trust people with German and foreign nationality, respectively. The heatmap shows predicted values for the trust gap based on estimates for equation (5). The table shows the difference between predicted in-group/out-group trust gap in a classroom exhibiting **cultural polarization (CP)** [58% natives, 42% Muslim immigrants] and four other classroom scenarios:

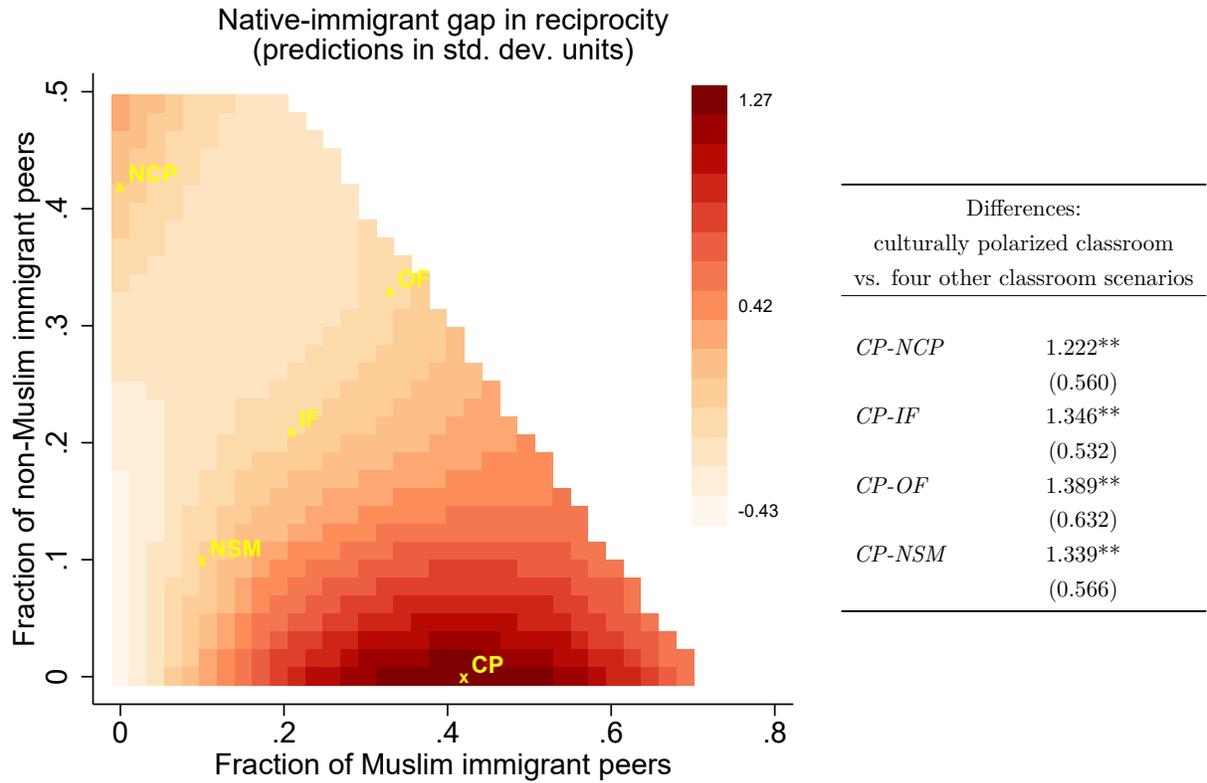
**Non-cultural polarization (NCP):** [58% natives, 42% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [58% natives, 21% Muslim immigrants, 21% non-Muslim immigrants];

**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

**Native supermajority (NSM):** [80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,108. 1 standard deviation = 1.59 points on a 0-10 Lickert scale. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.



**Figure 7.** Negative Stereotypes: Role of Peers' Behavior

*Notes:* For each individual  $i$ , the native-immigrant gap in peers' reciprocity is the difference between how much, on average, their native and immigrant classmates transfer back to natives in the second stage of the investment game (with the back transfers of each classmate averaged over the 11 possible investments from a native sender). The heatmap shows predicted values for this gap based on estimates for equation (5). The table shows the difference between the predicted gap in native vs. immigrant classmates' return behavior in a classroom exhibiting **cultural polarization (CP)** [58% natives, 42% Muslim immigrants] and four other classroom scenarios:

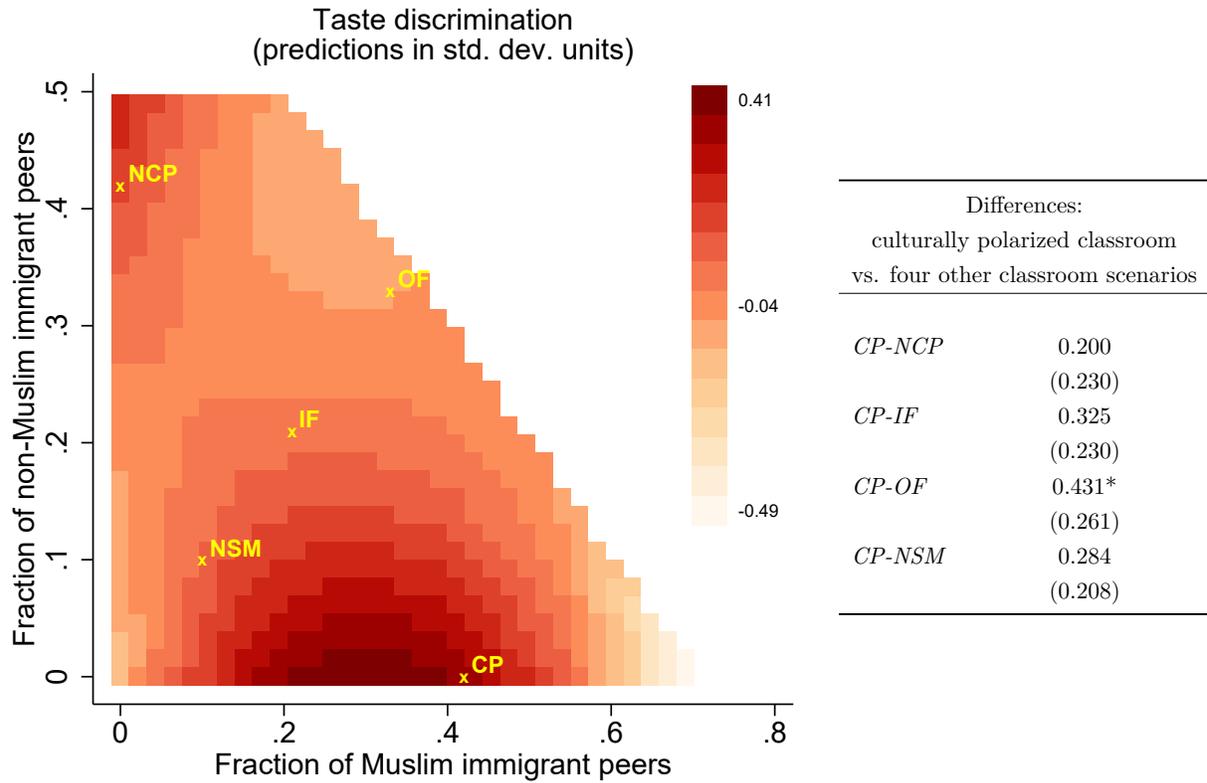
**Non-cultural polarization (NCP):** [58% natives, 42% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [58% natives, 21% Muslim immigrants, 21% non-Muslim immigrants];

**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

**Native supermajority (NSM):** [80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,101. 1 standard deviation = 1.3 euros. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ .



**Figure 8.** Taste Discrimination: Natives' Behavior when Playing as the Receiver

*Notes:* Taste discrimination is measured by the in-group/out-group return gap in the second stage of the investment game for natives. The heatmap shows predicted values for taste discrimination based on estimates for equation (5). The table shows the difference between predicted taste discrimination in a classroom exhibiting **cultural polarization (CP)** [58% natives, 42% Muslim immigrants] and four other classroom scenarios:

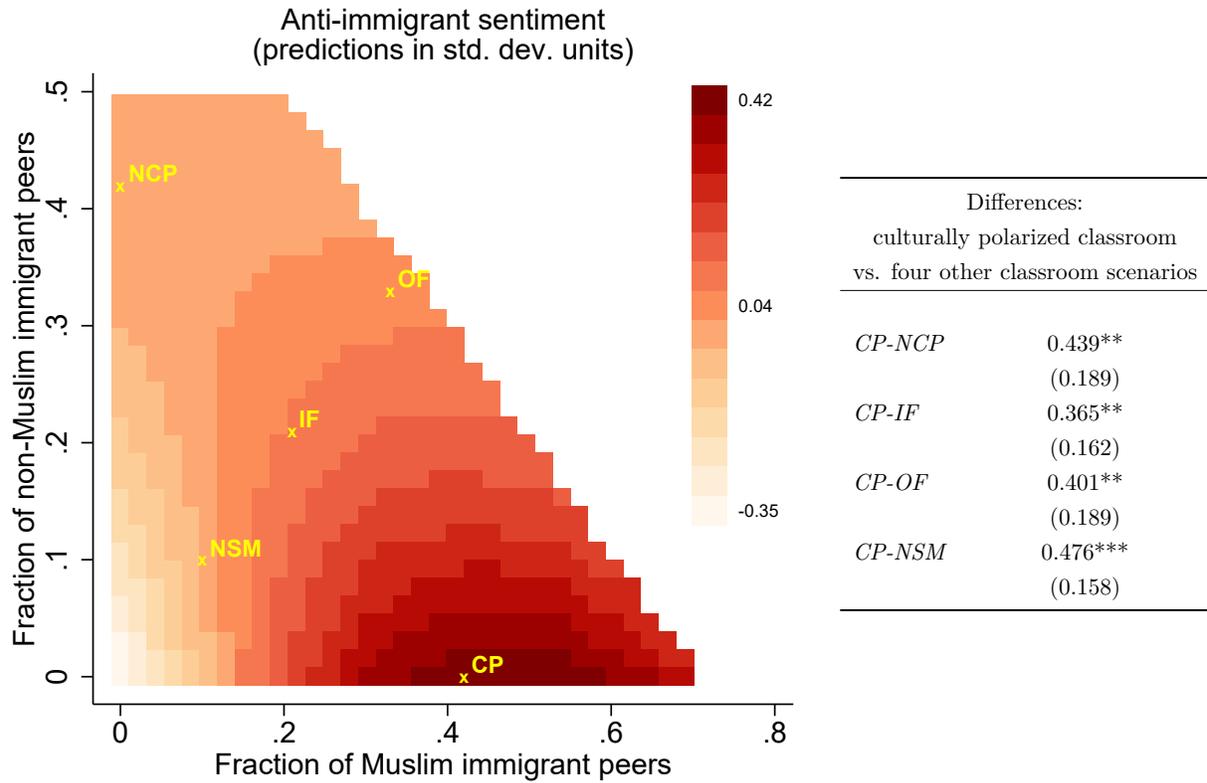
**Non-cultural polarization (NCP):** [58% natives, 42% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [58% natives, 21% Muslim immigrants, 21% non-Muslim immigrants];

**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

**Native supermajority (NSM):** [80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,088. 1 standard deviation = 1.05 euros. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.



**Figure 9.** Taste Discrimination: Anti-Immigrant Sentiment

*Notes:* Anti-immigration sentiment is measured by students' disagreement with three survey questions asking whether it is fair that workers of Turkish, Polish, and French descent are allowed to work in Germany (4=strongly disagree, 3=disagree somewhat, 2=agree somewhat, 1=strongly agree). Using a principal component analysis, we create an index of anti-immigration sentiment based on these questions. We normalize the index to have mean 0 and standard deviation 1. The heatmap shows predicted values for the index of anti-immigration sentiment based on estimates for equation (5). The table shows the difference between predicted values of anti-immigration sentiment in a classroom exhibiting **cultural polarization (P1)** [58% natives, 42% Muslim immigrants] and four other classroom scenarios:

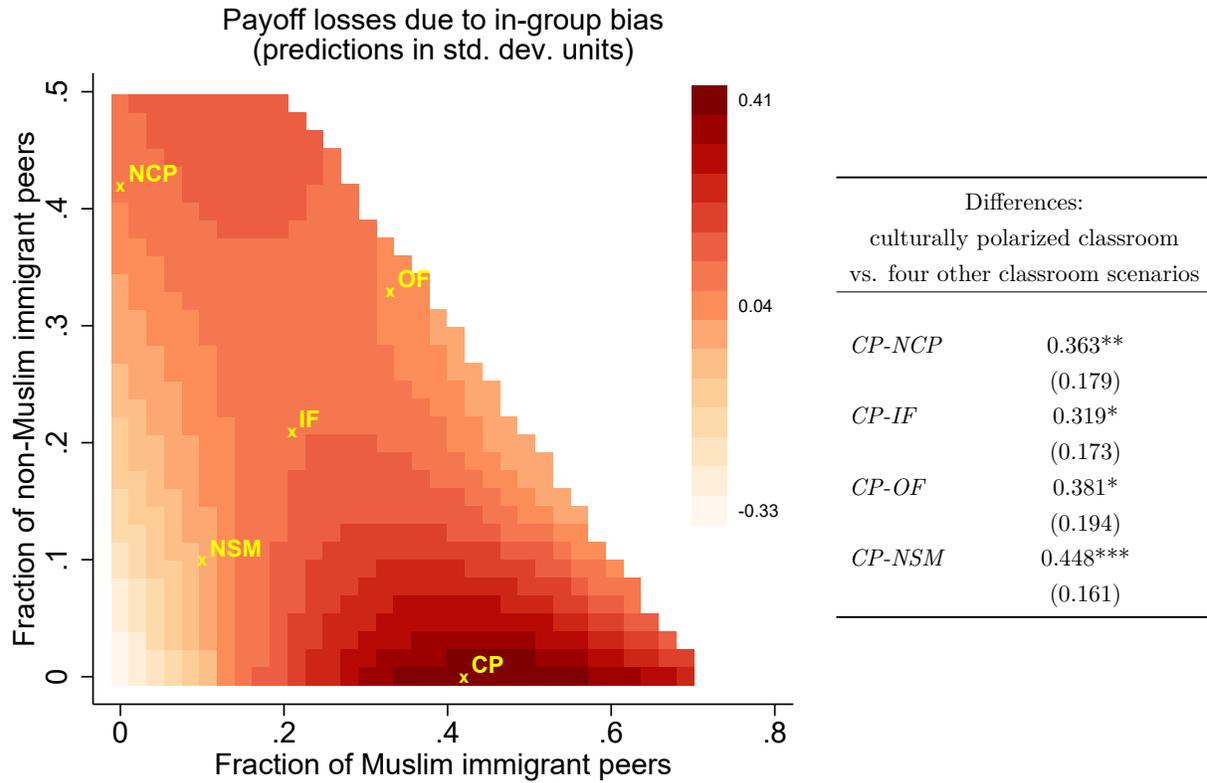
**Non-cultural polarization (NCP):** [58% natives, 42% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [58% natives, 21% Muslim immigrants, 21% non-Muslim immigrants];

**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

**Native supermajority (NSM):** [80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,088. 1 standard deviation = 0.85 points on 1-4 disagreement scale. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.



**Figure 10.** Payoff Losses Due to In-Group Bias

*Notes:* Payoff losses due to in-group bias are the differences in expected payoffs of senders when randomly matched with an immigrant versus a native interaction partner. The heatmap shows predicted values for payoff losses based on estimates for equation (5). The table shows the difference between predicted payoff losses in a classroom exhibiting **cultural polarization (CP)** [58% natives, 42% Muslim immigrants] and four other classroom scenarios:

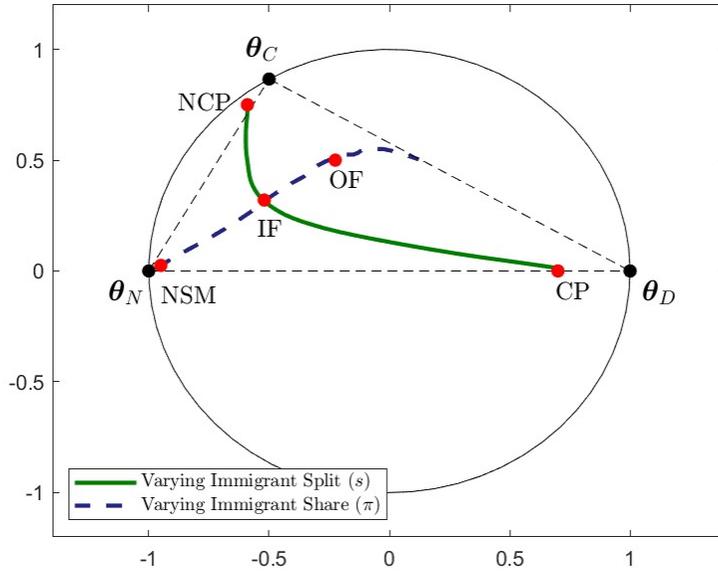
**Non-cultural polarization (NCP):** [58% natives, 42% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [58% natives, 21% Muslim immigrants, 21% non-Muslim immigrants];

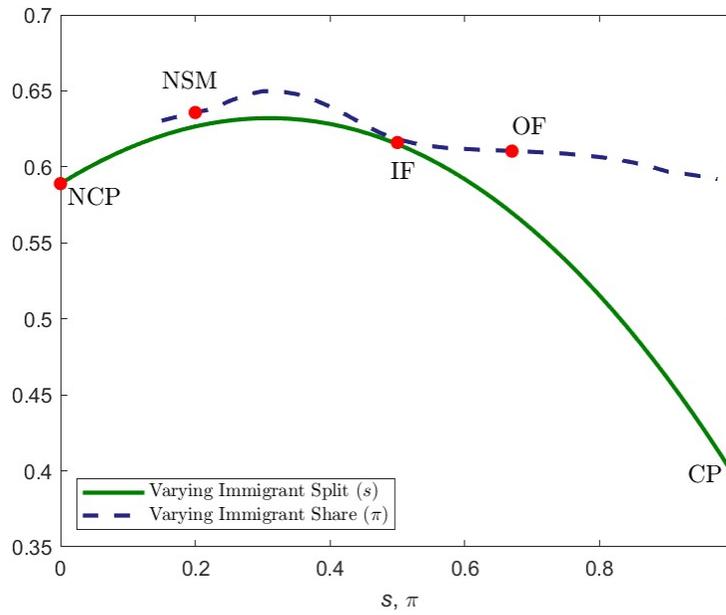
**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

**Native supermajority (NSM):** [80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,163. 1 standard deviation = 36 euro-cents. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.



(a) Cultural Identity  $\theta$  of the Inclusive Group



(b) Native Inclusive Group Joining Rate  $\mu_N$

**Figure 11.** A Three-Type Example with Natives and Two Immigrant Types

**Table 1.** Summary Statistics

	Mean	Std. Dev.
<i>Main dependent variable</i>		
In-group/out-group investment gap	0.094	0.758
<i>Main independent variables</i>		
Fraction immigrant peers	0.381	0.213
Fraction Muslim immigrant peers	0.188	0.178
Fraction non-Muslim immigrant peers	0.194	0.109
<i>Background variables</i>		
Male	0.539	0.499
Age	15.828	0.620
Age missing	0.017	0.130
Catholic	0.143	0.350
Protestant	0.534	0.499
Other religion or not religious	0.323	0.468
SES: Two-parent hh; high education	0.197	0.398
SES: Two-parent hh; low education	0.265	0.442
SES: Single-parent hh; high education	0.070	0.255
SES: Single-parent hh; low education	0.268	0.443
SES: missing	0.200	0.400
Age mother	46.215	8.455
Age mother: missing	0.051	0.221
Age father	50.356	11.077
Age father: missing	0.087	0.282
Observations	2,163	

*Notes:* Summary statistics for the main estimation sample. When defining Socioeconomic status (SES), a high education household (hh) is one where at least one parent has either a high school or university degree.

**Table 2.** Balancing Tests with and without School Fixed Effects

Dependent variable:	Fraction of immigrant peers	
	(1)	(2)
Male	-0.007 (0.010)	-0.005 (0.005)
Age	0.060*** (0.012)	-0.004 (0.005)
Age missing	-0.039 (0.034)	-0.031 (0.019)
Protestant	-0.093*** (0.019)	0.005 (0.008)
Other religion or not religious	-0.112*** (0.021)	-0.004 (0.008)
SES: Two-parent hh; low education	0.029** (0.014)	-0.009 (0.007)
SES: Single-parent hh; high education	0.016 (0.018)	-0.000 (0.009)
SES: Single-parent hh; low education	0.051*** (0.014)	-0.007 (0.006)
SES: missing	0.076*** (0.015)	0.000 (0.007)
Age mother	-0.002* (0.001)	0.001 (0.001)
Age mother: missing	0.074* (0.044)	-0.015 (0.023)
Age father	-0.001 (0.001)	-0.000 (0.000)
Age father: missing	0.060 (0.044)	0.027 (0.019)
Class size	-0.004 (0.003)	0.001 (0.003)
Observations	2,163	2,163
R-squared	0.108	0.744
F-statistic	6.26	1.20
p-value	0.000	0.279
School FE		✓

*Notes:* OLS regressions with the fraction of immigrant peers as the dependent variable. Standard errors are reported in parentheses and are clustered at the classroom level. \*\*\*, \*\*, \* indicate significance at the 1%, 5%, and 10% level, respectively.

**Table 3.** In-Group Bias: An Inverted-U in the Share of Immigrant Peers

Dependent variable:	In-group/out-group investment gap		
	(1)	(2)	(3)
Panel A: Linear Specification			
Fraction immigrant peers	0.035 (0.164)	0.030 (0.166)	0.054 (0.164)
Observations	2,163	2,163	2,163
R-squared	0.051	0.052	0.057
Panel B: Quadratic Specification			
Fraction immigrant peers	1.662*** (0.394)	1.668*** (0.397)	1.651*** (0.396)
Fraction immigrant peers squared	-1.881*** (0.455)	-1.896*** (0.459)	-1.850*** (0.457)
p-value: coeffs. jointly equal to zero	0.0002	0.0002	0.0002
p-value: coeffs. equal but opposite in sign	0.204	0.193	0.248
Observations	2,163	2,163	2,163
R-squared	0.055	0.057	0.062
Basic controls	✓	✓	✓
Religious background		✓	✓
Family background			✓

*Notes:* OLS estimates of equation (3). The dependent variable is the in-group/out-group investment gap, normalized to be mean 0 and standard deviation 1. Basic controls include school fixed effects, a student's gender and age, and class size. Religious background includes three dummy variables for a student's religious affiliation (Catholic, Protestant, other religion or not religious). Family background includes dummy variables for student's SES as listed in Table 1. Standard errors are reported in parentheses and are clustered at the classroom level. \*\*\*, \*\*, \* indicate significance at the 1%, 5%, and 10% level, respectively.

**Table 4.** In-Group Bias: How Type of Diversity Matters

Dependent variable:	In-group/out-group investment gap		
	(1)	(2)	(3)
Fraction Muslim immigrant peers ( $\pi_{D,k}$ )	3.910*** (1.011)	3.892*** (1.006)	3.759*** (1.024)
Fraction non-Muslim immigrant peers ( $\pi_{C,k}$ )	1.896*** (0.693)	1.885*** (0.698)	1.825*** (0.702)
$\pi_{D,k}^2$	-4.787*** (1.243)	-4.733*** (1.240)	-4.489*** (1.264)
$\pi_{C,k}^2$	-1.903 (1.280)	-1.944 (1.289)	-1.881 (1.290)
$\pi_{D,k} \times \pi_{C,k}$	-14.884** (6.487)	-14.640** (6.492)	-13.764** (6.631)
$\pi_{D,k}^2 \times \pi_{C,k}$	11.339* (6.332)	10.645* (6.387)	9.469 (6.487)
$\pi_{D,k} \times \pi_{C,k}^2$	10.965 (8.441)	11.132 (8.425)	10.635 (8.540)
p-value: coeffs. jointly equal to zero	0.0005	0.0005	0.0008
Observations	2,163	2,163	2,163
R-squared	0.057	0.059	0.064
Basic controls	✓	✓	✓
Religious background		✓	✓
Family background			✓

OLS estimates of equation (5). The dependent variable is the in-group/out-group investment gap, normalized to be mean 0 and standard deviation 1. Basic controls include school fixed effects, a student's gender and age, and class size. Religious background includes three dummy variables for a student's religious affiliation (Catholic, Protestant, other religion or not religious). Family background includes dummy variables for a student's SES as listed in Table 1. Standard errors are reported in parentheses and are clustered at the classroom level. \*\*\*, \*\*, \* indicate significance at the 1%, 5%, and 10% level, respectively.

**Table 5.** In-Group Bias: Classroom Polarization versus Classroom Fractionalization

Dependent variable:	In-group/out-group investment gap		
	(1)	(2)	(3)
Classroom polarization	0.643*** (0.142)		0.661*** (0.198)
Classroom fractionalization		0.368** (0.142)	-0.028 (0.192)
Observations	2,163	2,163	2,163
R-squared	0.062	0.059	0.062
Basic controls	✓	✓	✓
Religious background	✓	✓	✓
Family background	✓	✓	✓

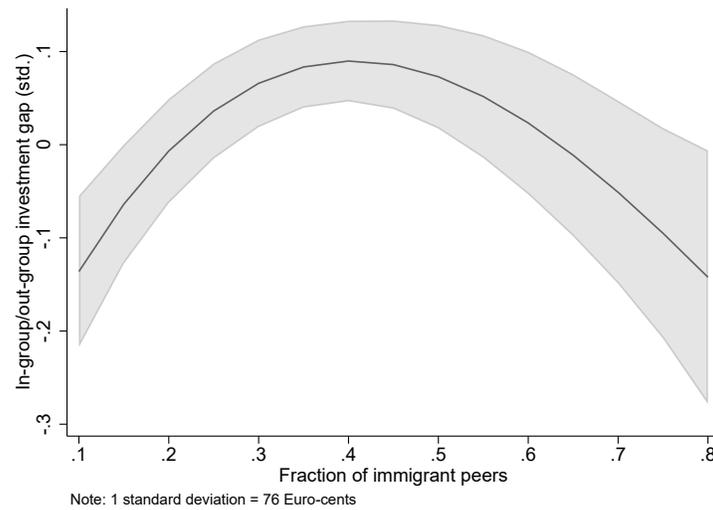
OLS estimates. The dependent variable is the in-group/out-group investment gap, normalized to be mean 0 and standard deviation 1. Classroom polarization and fractionalization are defined using equation 2 and the 11 origin countries given in Table A2. All regressions include the full set of control variables used in column (3) of Table 3. Standard errors are reported in parentheses and are clustered at the classroom level. \*\*\*, \*\*, \* indicate significance at the 1%, 5%, and 10% level, respectively.

# Appendix for Online Publication

## “Diversity and Discrimination in the Classroom”

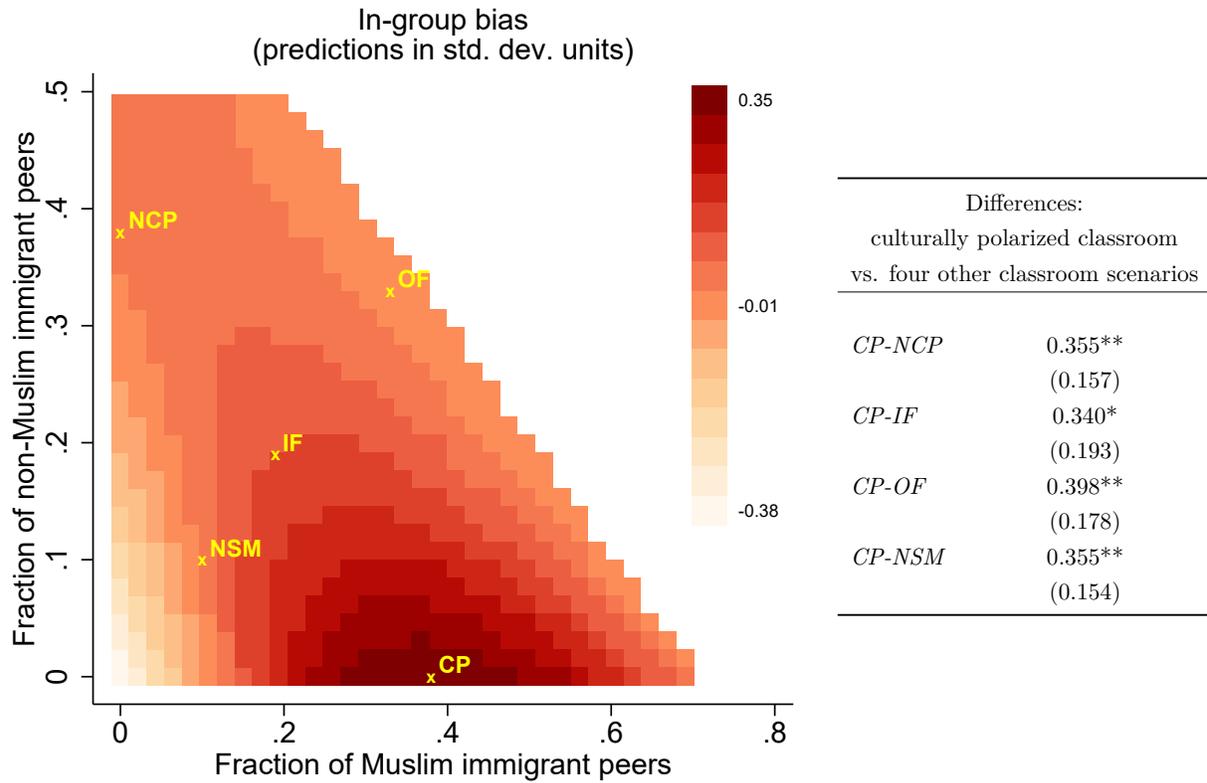
by Dan Anderberg, Gordon B. Dahl, Christina Felfe, Helmut Rainer, and Thomas Siedler

### Appendix A: Additional Figures and Tables



**Figure A1.** In-Group Bias Using a Cubic Specification

*Notes:* The figure shows predictions similar to those in Figure 4, but with diversity modeled as a cubic polynomial in the the proportion of immigrant peers. The grey shaded area denotes pointwise 95% confidence intervals.



**Figure A2.** Robustness: Third-Order Expansion Estimates for In-Group Bias

*Notes:* In-group bias is measured by the in-group/out-group investment gap in the first stage of the investment game for natives. The heatmap shows predicted values for in-group bias using a complete third-order expansion of the non-Muslim and Muslim immigrant fractions in a classroom. The peak in in-group bias occurs in a classroom exhibiting **cultural polarization (CP)**, with native German peers forming a slight majority group [62%] and Muslim immigrant peers a large minority group [38%]. The table shows the difference between predicted in-group bias at this peak and four other classroom scenarios:

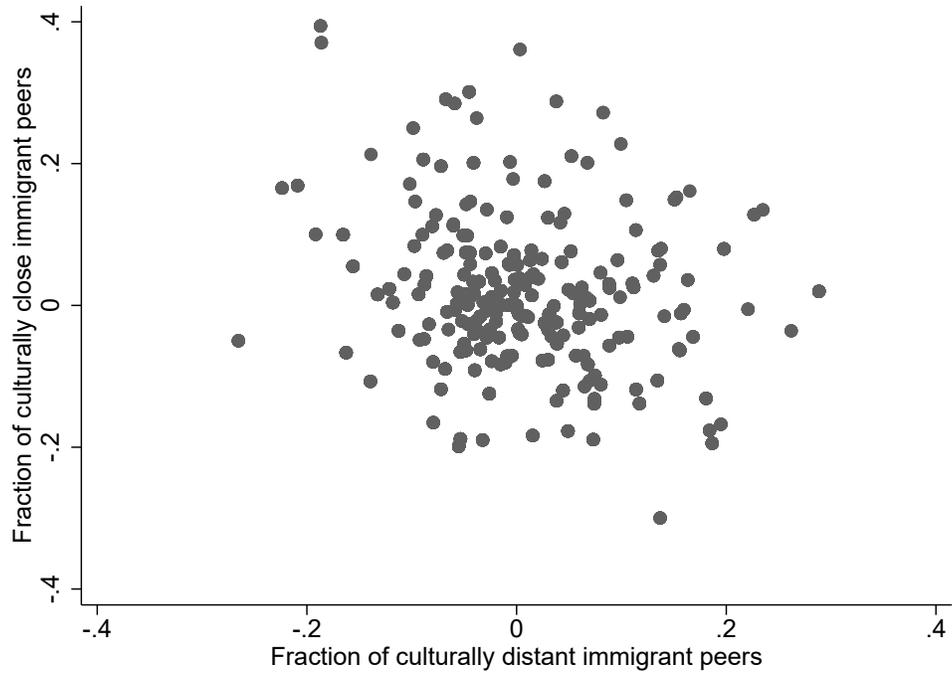
**Non-cultural polarization (NCP):** [62% natives, 38% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [62% natives, 19% Muslim immigrants, 19% non-Muslim immigrants];

**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

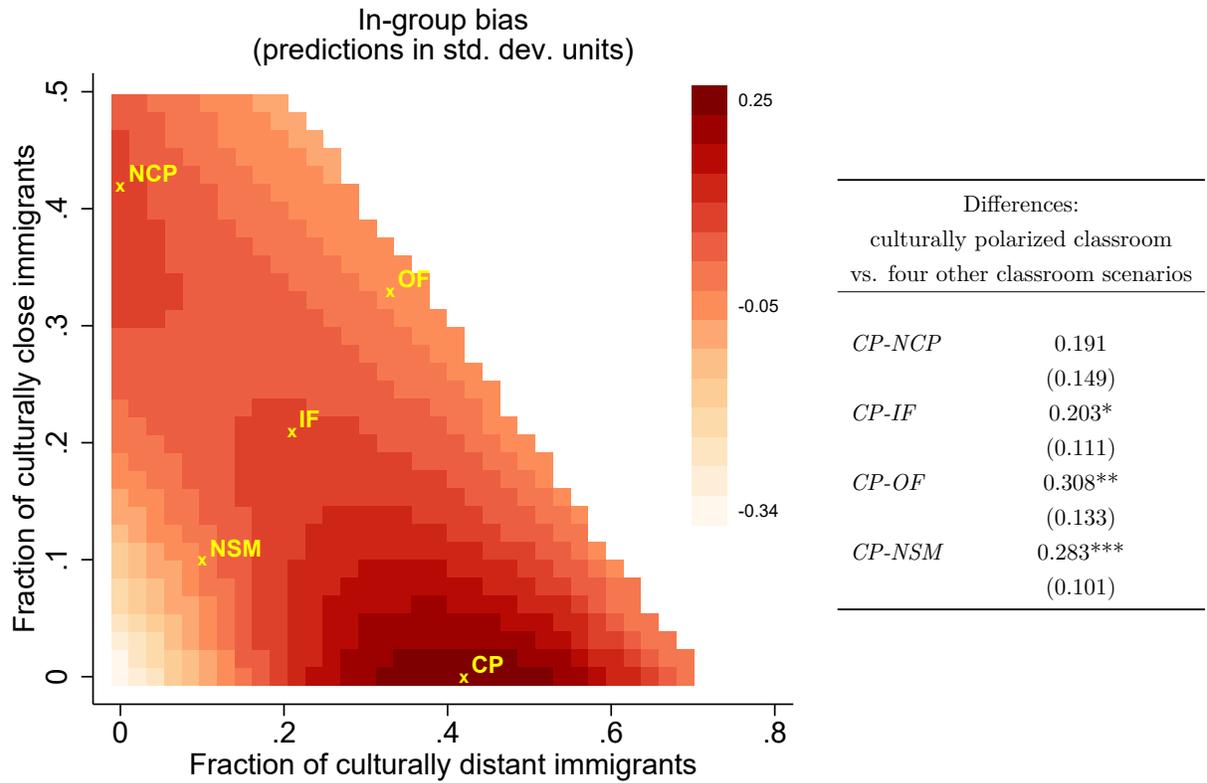
**Native supermajority (NSM):** [80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,163. 1 standard deviation = 76 euro-cents. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.



**Figure A3.** Linguistic Distance: Identifying Variation after Removing School Fixed Effects

*Notes:* This figure parallels the right panel of Figure 3, but using linguistic distance to define culturally close and distant immigrants.



**Figure A4.** Robustness: Using Linguistic Distance to Define Culturally Close and Distant Immigrant Peers

*Notes:* In-group bias is measured by the in-group/out-group investment gap in the first stage of the investment game. The heatmap shows predicted values for in-group bias based on estimates for equation 5. The peak in in-group bias occurs in a classroom exhibiting cultural polarization (CP), with native German peers forming a slight majority group [58%] and Muslim immigrant peers a large minority group [42%]. The table shows the difference between predicted in-group bias at this peak and four other classroom scenarios:

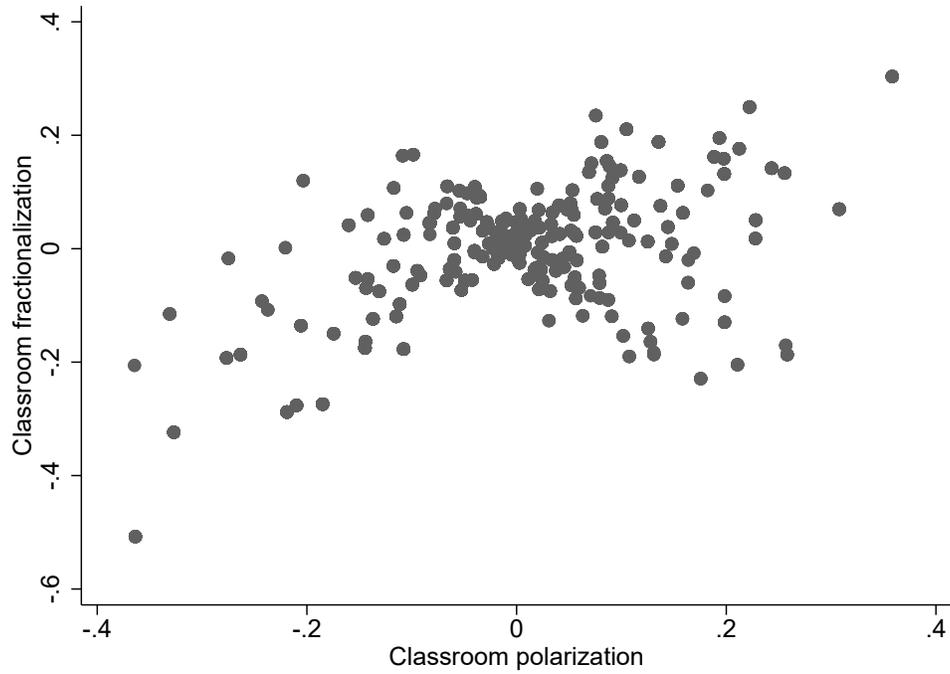
**Non-cultural polarization (NCP):** [58% natives, 42% non-Muslim immigrants];

**Immigrant fractionalization (IF):** [58% natives, 21% Muslim immigrants, 21% non-Muslim immigrants];

**Overall fractionalization (OF):** [34% natives, 33% Muslim immigrants, 33% non-Muslim immigrants];

**Native supermajority (NSM):**[80% natives, 10% Muslim immigrants, 10% non-Muslim immigrants].

N=2,163. 1 standard deviation = 76 euro-cents. \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.



**Figure A5.** Polarization and Fractionalization: Identifying Variation after Removing School Fixed Effects

*Notes:* Classroom polarization and fractionalization are defined using equation 2 and the 11 origin countries given in Table A2. The figure shows a scatter plot of the residualized fractionalization index against the residualized polarization index, both net of school fixed effects.

**Table A1.** Additional Balancing Tests with and without School Fixed Effects

Dependent variable:	Fraction of Muslim immigrant peers	
	(1)	(2)
Male	-0.008 (0.008)	-0.005* (0.003)
Age	0.050*** (0.010)	0.001 (0.004)
Age missing	-0.028 (0.029)	-0.020 (0.015)
Protestant	-0.095*** (0.016)	-0.001 (0.007)
Other religion or not religious	-0.116*** (0.018)	-0.009 (0.007)
SES: Two-parent hh; low education	0.026** (0.010)	-0.003 (0.004)
SES: Single-parent hh; high education	0.016 (0.013)	0.005 (0.006)
SES: Single-parent hh; low education	0.043*** (0.011)	0.002 (0.005)
SES: missing	0.070*** (0.013)	0.009 (0.007)
Age mother	-0.002** (0.001)	0.000 (0.001)
Age mother: missing	0.085** (0.041)	0.008 (0.022)
Age father	-0.002* (0.001)	-0.000 (0.000)
Age father: missing	0.067* (0.038)	0.030 (0.018)
Class size	-0.006** (0.002)	-0.003 (0.002)
Observations	2,163	2,163
R-squared	0.149	0.773
F-statistic	8.24	1.37
p-value	0.000	0.171
School FE		✓

*Notes:* OLS regressions with the fraction of Muslim immigrant peers as the dependent variable. Standard errors are reported in parentheses and are clustered at the classroom level. \*\*\*, \*\*, \* indicate significance at the 1%, 5%, and 10% level, respectively.

**Table A2.** Origin Countries of Peers, Mapped into 11 Regions

	Mean	St. Dev.
Overall fraction immigrant peers	0.381	0.213
<i>Fraction of immigrant peers from:</i>		
Turkey	0.143	0.149
Balkan countries	0.037	0.050
Eastern European countries	0.045	0.053
Post Soviet countries	0.040	0.050
Southern European countries	0.020	0.040
Central and Northern European countries	0.013	0.030
Middle Eastern countries	0.024	0.043
Asian countries	0.023	0.038
African countries	0.023	0.044
Other countries	0.010	0.026
Unidentified countries	0.003	0.015
Observations	2,163	

*Notes:* Origin country of peers defined by parent's country of origin. See text for details.

## Appendix B: A Framework

In this appendix we present further details of the framework described in Section 6.

### 1.1 Setup and Equilibrium Group Joining Rates

The setup is as outlined in Section 6 and is repeated here only for completeness. Consider an economy with a continuum of individuals who are of  $J \geq 2$  types. Let the proportions of types in the population be denoted  $\boldsymbol{\pi} = (\pi_1, \dots, \pi_J)$ .  $J + 1$  groups form endogenously. For each type  $j$  there is an *exclusive* group comprising only  $j$ -type individuals. The final group is an *inclusive* group with members of all types. Individual  $i$  of type  $j$ , chooses between two options: join the own-type exclusive group (option 0) or join the inclusive group (option 1). Joining the inclusive group involves a type-specific joining cost  $h_j$  which below will be related to cultural distance. Finally, let  $\beta > 0$  parameterize the (common) strength of preference for group size. The individual's utilities of the two options are,

$$u_{ij}^0 = \beta\pi_j(1 - \mu_j) + \varepsilon_{ij}^0, \quad u_{ij}^1 = \beta \sum_{j'=1}^J \pi_{j'}\mu_{j'} - h_j + \varepsilon_{ij}^1, \quad (\text{B.1})$$

where  $\mu_j$  is the proportion of type- $j$  individuals who join the inclusive group, and where we used that the size of the type- $j$  exclusive group is  $\pi_j(1 - \mu_j)$  while the size of the inclusive group is  $\sum_{j'=1}^J \pi_{j'}\mu_{j'}$ .  $\varepsilon_{ij}^0$  and  $\varepsilon_{ij}^1$  are the individual's choice-specific random preferences. To rule out multiple coordination equilibria we place an upper limit on  $\beta$ , and we also make a specific, but standard, distributional assumption on the random preferences.

**Assumption 1.** (*Preferences*) *The preference for group-size satisfies  $\beta \in (0, 2)$ . Any individual  $i$ 's random preferences  $\varepsilon_{ij}^0$  and  $\varepsilon_{ij}^1$  are i.i.d. extreme value distributed.*

With the individual random preferences being extreme value distributed, the proportion  $\mu_j$  of type  $j$  individuals who join the inclusive group satisfies

$$\rho_j(\boldsymbol{\mu}) \equiv \log\left(\frac{\mu_j}{1 - \mu_j}\right) - \beta \left( \sum_{j'=1}^J \pi_{j'}\mu_{j'} - \pi_j(1 - \mu_j) \right) + h_j = 0, \quad j = 1, \dots, J, \quad (\text{B.2})$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_J)$ . Similarly, let  $\mathbf{h} = (h_1, \dots, h_J)$  which, for the moment, we take as given.

The inclusive group joining rate of any one type affects the joining incentives of all other types. Hence an equilibrium in joining rates requires both “within-type” consistency (for each type  $j$ ,  $\mu_j$  should satisfy (B.2) for that type given the joining rates of all other types) and “across-type” consistency (within-type consistency should hold for all  $J$  types simultaneously).

**Definition 1.** *A “joining equilibrium”, given the type distribution  $\boldsymbol{\pi}$  and the type-specific joining costs  $\mathbf{h}$ , is a vector  $\boldsymbol{\mu}$  that satisfies (B.2) for all  $J$  types simultaneously.*

The following notes that a unique joining equilibrium exists.

**Proposition 1.** (*Existence of joining equilibrium*). *For any given type distribution  $\boldsymbol{\pi}$  and type-specific joining costs  $\mathbf{h}$ , a unique joining equilibrium  $\boldsymbol{\mu}$  exists.*

**Proof.** We start by showing existence. For any type  $j$  and any given  $\boldsymbol{\mu}_{-j}$  (the vector  $\boldsymbol{\mu}$  without the  $j$ 'th component), the equation  $\rho_j(\boldsymbol{\mu}) = 0$  has a unique solution  $\mu_j \in (0, 1)$ . To see this, note

that, for any  $\boldsymbol{\mu}_{-j}$ ,  $\rho_j(\boldsymbol{\mu})$  is continuously differentiable in  $\mu_j$  and goes to  $-\infty$  as  $\mu_j \rightarrow 0$  and to  $+\infty$  as  $\mu_j \rightarrow 1$ . Moreover, differentiating (B.2) gives that,

$$\frac{\partial \rho_j}{\partial \mu_j} = S_j \equiv \frac{1}{\mu_j(1-\mu_j)} - 2\beta\pi_j > 0, \quad j = 1, \dots, J, \quad (\text{B.3})$$

where the sign follows from the fact that  $\mu_j(1-\mu_j) \leq 1/4$  and  $\beta \in (0, 2)$  (Assumption 1). This establishes that a unique solution  $\mu_j \in (0, 1)$  to  $\rho_j(\boldsymbol{\mu}) = 0$  given  $\boldsymbol{\mu}_{-j}$  exists. Denote this unique solution by  $\hat{\mu}_j(\boldsymbol{\mu}_{-j})$ . By the implicit function theorem  $\hat{\mu}_j(\cdot)$  is also a continuous function. Extending the argument of  $\hat{\mu}_j(\cdot)$  to  $\boldsymbol{\mu}$  (where it is understood that  $\hat{\mu}_j(\cdot)$  only depends on  $\boldsymbol{\mu}_{-j}$ , not on  $\mu_j$ ) and forming  $\hat{\boldsymbol{\mu}}(\boldsymbol{\mu}) = (\hat{\mu}_1(\boldsymbol{\mu}), \dots, \hat{\mu}_J(\boldsymbol{\mu}))$  we thus have a continuous function that maps  $[0, 1]^J$  into itself. A fixed point of this mapping is a joining equilibrium and by Brouwer's theorem, such a fixed point is guaranteed to exist.

Consider next uniqueness. Let  $\boldsymbol{\rho}(\boldsymbol{\mu}) = (\rho_1(\boldsymbol{\mu}), \dots, \rho_J(\boldsymbol{\mu}))$  and note that a joining equilibrium is a  $\boldsymbol{\mu}$  such that  $\boldsymbol{\rho}(\boldsymbol{\mu}) = \mathbf{0}$ , the null vector of length  $J$ . If the mapping  $\boldsymbol{\rho}(\boldsymbol{\mu})$  can be shown to be univalent, then uniqueness follows since the null vector can then have at most one pre-image (and given existence we know that it has at least one pre-image). Gale and Nikaido (1965) showed that if the Jacobian matrix of a map  $\boldsymbol{\rho}(\cdot)$  is a P-matrix, then  $\boldsymbol{\rho}(\cdot)$  is univalent. A matrix  $\mathbf{A} \in \mathbb{R}^{J \times J}$  is a P-matrix if all its principal minors are positive, and a sufficient condition for this is the dominant diagonal condition:  $|a_{jj}| \geq \sum_{j' \neq j} |a_{jj'}|$  for all  $j$ . That is, for each row in  $\mathbf{A}$ , the absolute value of the diagonal term  $a_{jj}$  is no less than the sum of the absolute values of the off-diagonal terms in that row. Using that the diagonal terms of the Jacobian of  $\boldsymbol{\rho}(\boldsymbol{\mu})$  take the form (B.3) and the off-diagonal terms are  $\partial \rho_j / \partial \mu_{j'} = -\beta\pi_{j'}$ ,  $j' \neq j$ , the Jacobian of  $\boldsymbol{\rho}(\cdot)$  has a dominant diagonal if, for all  $j$ ,  $S_j \geq \beta \sum_{j' \neq j} \pi_{j'}$ . But, since  $\sum_{j' \neq j} \pi_{j'} = 1 - \pi_j$ , this is equivalent to  $[\mu_j(1-\mu_j)]^{-1} \geq \beta(1+\pi_j)$  which holds since the left hand side is not less than four and the right hand side is, using Assumption 1, strictly less than four. #

## 1.2 Cultural Identities and Inclusive-Group Joining Costs

There is a space of cultural identities which we take to be  $\Theta \equiv [-1, +1]^{J-1}$ . Each type  $j$  is endowed with some *exogenously* given identity  $\boldsymbol{\theta}_j \in \Theta$ . The inclusive group has an *endogenous* identity  $\boldsymbol{\theta} \in \Theta$ , making the type-specific (Euclidean) distance  $d_j \equiv \|\boldsymbol{\theta} - \boldsymbol{\theta}_j\|$ . We assume a type-specific cost of joining the inclusive group that depends on this cultural distance:  $h_j = h_j(d_j)$  where, for each type  $j$ ,  $h_j(\cdot)$  is twice continuously differentiable and satisfies  $h_j(0) = 0$ ,  $h'_j(\cdot) > 0$  and  $h''_j(\cdot) > 0$ . To close the model we need to specify how  $\boldsymbol{\theta}$  is determined. In line with the assumption of positive preferences for group size, we assume that  $\boldsymbol{\theta}$  is chosen so as to maximize the inclusive group's size (or "popularity") in the joining equilibrium that ensues.

**Definition 2.** (*Inclusive group identity*) The cultural identity of the inclusive group  $\boldsymbol{\theta} \in \Theta$  maximizes the size of the inclusive group,  $\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta} \in \Theta} \Omega(\boldsymbol{\theta})$ , where  $\Omega(\boldsymbol{\theta}) \equiv \sum_{j=1}^J \pi_j \mu_j$ .

If the types have different population frequencies, then it is natural that  $\boldsymbol{\theta}$  is set close to the endowed identity of a relatively more frequent type. However, the type specific aversions to distance also matter. The case of just two types allows us to illustrate this.

### 1.2.1 Comparative Statics with Two-Types

Consider the  $J = 2$  case: natives (type 1) and immigrants (type 2). With  $J = 2$ , the identity space is  $\Theta = [-1, +1]$  and we can assume that the endowed cultural identities of natives and immigrants

are  $\theta_1 = -1$  and  $\theta_2 = +1$  respectively. The identity of the inclusive group is a scalar  $\theta \in \Theta$  whereby  $d_1 = 1 + \theta$  and  $d_2 = 1 - \theta$ , with corresponding type-specific joining costs  $h_1 = h_1(1 + \theta)$  and  $h_2 = h_2(1 - \theta)$ . The inclusive group identity  $\theta$  maximizes  $\Omega(\theta) \equiv \sum_{j=1}^2 \pi_j \mu_j$  where the joining rates  $\boldsymbol{\mu} = (\mu_1, \mu_2)$  simultaneously satisfy (B.2) for  $j = 1, 2$ .

We will consider deviations from “symmetry” where the two types are of equal proportion in the population,  $\pi_1 = \pi_2 = 1/2$ , and have the same distance cost function,  $h_j(\cdot) = h(\cdot)$  for some common  $h(\cdot)$ . Symmetric population proportions and costs also imply that  $\theta = 0$  is a stationary point of  $\Omega(\theta)$  (see below) at which both types have the same inclusive group joining rate,  $\mu^*$ , satisfying

$$\log\left(\frac{\mu^*}{1 - \mu^*}\right) - \frac{\beta}{2}(3\mu^* - 1) + h(1) = 0. \quad (\text{B.4})$$

Consider first how  $\mu_1$  and  $\mu_2$  are affected by  $\theta$  at symmetry? Simple comparative statics show that

$$\left.\frac{\partial \mu_2}{\partial \theta}\right|_{sym} = -\left.\frac{\partial \mu_1}{\partial \theta}\right|_{sym} = \frac{h'(1)}{S^* + C^*} > 0, \quad (\text{B.5})$$

where  $S^* = [\mu^*(1 - \mu^*)]^{-1} - \beta > 0$  and  $C^* = \beta/2 > 0$ , whereby  $S^* + C^* = [\mu^*(1 - \mu^*)]^{-1} - \beta/2 > 0$ . This trivially confirms that  $\theta = 0$  is indeed a stationary point of  $\Omega(\theta)$  under symmetric population proportions and distance costs<sup>23</sup>

$$\left.\frac{\partial \Omega(\theta)}{\partial \theta}\right|_{sym} = \frac{1}{2} \sum_{j=1}^2 \left.\frac{\partial \mu_j}{\partial \theta}\right|_{sym} = 0. \quad (\text{B.6})$$

Similarly, at symmetry, and using  $\pi_1 = 1 - \pi_2$ ,

$$\left.\frac{\partial \mu_1}{\partial \pi_2}\right|_{sym} = -\left.\frac{\partial \mu_2}{\partial \pi_2}\right|_{sym} = \frac{\beta(1 - \mu^*)}{S^* + C^*} > 0. \quad (\text{B.7})$$

Turning to the cross-partials, first for the equilibrium joining rates, it can be shown that

$$\left.\frac{\partial^2 \mu_j}{\partial \theta \partial \pi_2}\right|_{sym} = \frac{h'(1)}{D^*} \left( \frac{(2\mu^* - 1)\beta(1 - \mu^*)}{(\mu^*)^2(1 - \mu^*)^2(S^* + C^*)} + 3\beta \right), \quad (\text{B.8})$$

where  $D^* = (S^* + C^*)(S^* - C^*) > 0$ . More importantly for our purposes however, the cross partial of  $\Omega(\theta)$  at symmetry (after substituting using (B.5), and (B.8) and also using the expressions for  $S^*$ ,  $C^*$  and  $D^*$ ) can be shown to be

$$\left.\frac{\partial \Omega(\theta)}{\partial \theta \partial \pi_2}\right|_{sym} = \frac{8\mu^*(1 - \mu^*)((2 - \beta) + \beta\mu^*(2 - \mu^*))h'(1)}{(2 - \beta\mu^*(1 - \mu^*))^2(2 - 3\beta\mu^*(1 - \mu^*))} > 0, \quad (\text{B.9})$$

where the sign follows from  $\beta \in (0, 2)$  (Assumption 1). As a result, it follows that  $\theta$  is increasing in  $\pi_2$  at symmetry. That is, when type 2 (type 1) becomes a majority (minority) type, the inclusive group’s identity moves closer to type 2 (away from type 1).

<sup>23</sup>Different sufficient conditions can be given for  $\theta = 0$  to be a (local) maximum of  $\Omega(\theta)$  under symmetric type proportions and distance costs. Either sufficient convexity of the common cost function,  $h''(1)/(h'(1))^2 > 1$  (for any  $\mu^*$ ) or any convexity  $h''(1) > 0$  and  $\mu^* > 1/2$  (where the latter would be implied by  $h(1) < \beta/4$ ).

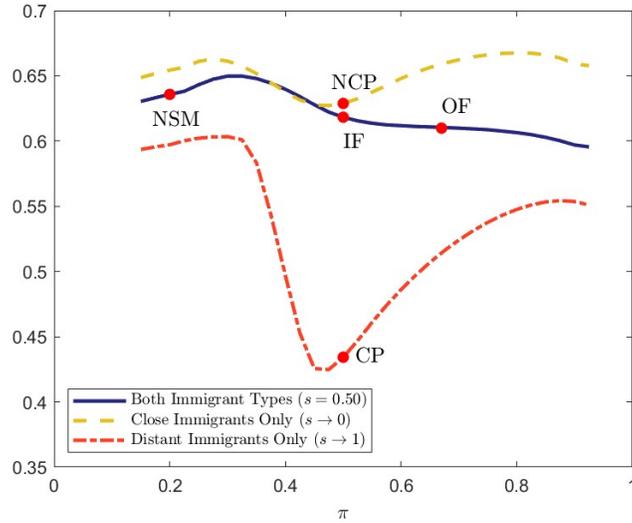
Consider next the effect of introducing type-differences in the distance cost. Specifically, let  $h_1(d) = (1 - \epsilon) h(d)$  and  $h_2(d) = (1 + \epsilon) h(d)$ , where we then consider the introduction of  $\epsilon > 0$  from symmetry ( $\epsilon = 0$ ). It can then be shown that,

$$\left. \frac{\partial^2(\theta)}{\partial\theta\partial\epsilon} \right|_{sym} = \frac{h'(1)}{D^*} \left[ (S^* + C^*) + \frac{(2\mu^* - 1)}{(\mu^*)^2 (1 - \mu^*)^2} \frac{h(1)}{(S^* + C^*)} \right]. \quad (\text{B.10})$$

The main effect here – captured by the first component – is that  $\epsilon > 0$  increases the cost for type 2 relative to type 1, whereby  $\theta$  is moved closer to type 2 to retain participation from this group. The second component in (B.10) (which is also positive whenever  $\mu^* > 1/2$ ) reflects an impact on the relative size of the social multipliers,  $S_1$  and  $S_2$ .

### 1.2.2 Illustrative Example: Limiting Cases

The example provided in Section 6 had three types: natives ( $N$ ), culturally close immigrants ( $C$ ), culturally distant immigrants ( $D$ ). The overall immigrant share was denoted by  $\pi$  and the split between the immigrant types was denoted  $s$ , whereby  $\pi_N = 1 - \pi$ ,  $\pi_C = \pi(1 - s)$  and  $\pi_D = \pi s$ . When either  $s \rightarrow 0$  or  $s \rightarrow 1$  the example limits to just natives and one immigrant type. In the former limiting case, the only immigrant type to exist is the culturally close type. In the latter limiting case, the only immigrant type to exist is the culturally distant type. Here we showcase how the example given also provides a good correspondence to our empirical findings in these limiting cases. For instance, increasing the immigrant share  $\pi$  with  $s = 0$ , is equivalent to moving on a vertical ray out from the origin in our heatmaps, that is, along the  $y$ -axis. Similarly, increasing  $\pi$  whilst for  $s = 1$ , is equivalent to moving on a horizon ray out from the origin in our heatmaps, that is, along the  $x$ -axis.



**Figure B.1.** Native Inclusive Group Joining Rate  $\mu_N$  in the Limiting Cases  $s = 0$  and  $s = 1$

Figure B.1 illustrates the equilibrium native inclusive group joining rate  $\mu_n$  in these two limiting cases for the particular parameterization used in Section 6 ( $\gamma_N = 0.35$ ,  $\gamma_C = 0.6$ ,  $\gamma_D = 0.6$ ,  $\beta = 1.75$  and  $\sigma = 1.1$ ). The blue solid line is for  $s = 1/2$  as already shown in Figure 11. The yellow dashed line is for the limiting case where  $s \rightarrow 0$  and thus shows how the native inclusive group joining rate

varies with the immigrant share when the only immigrant type is the culturally close type. In line with the empirical findings, the native joining rate drops when  $\pi$  approaches a half, giving rise to non-cultural polarization (*NCP*). However, the drop in the native joining rate in this case is modest relative to the drop that occurs under cultural polarization (*CP*). Indeed, in the limiting case with only culturally distant immigrants,  $s \rightarrow 1$ , the native joining rate exhibits a strong U-shape as the immigrant share  $\pi$  is varied, directly matching the inverted U-shape in in-group bias when moving along the  $x$ -axis in Figure 4.

## Appendix C: Translation of Experimental Instructions and Decision Sheets

### *Experimental Instructions*

Participation in this game is voluntary!

Thank you very much for participating. From now on, please do not speak with anyone else apart from us about the game. Unfortunately, if you break this rule, we will have to exclude you from the game.

The objective of this game is to examine how people make decisions. There are no “right” or “wrong” decisions in the game and our aim is not to test your knowledge. Make your decisions exactly as you wish. During this game, you will be earning real money. We guarantee that you will receive a cash payout within two weeks. You will receive your money in an envelope marked with your ID number, so please make sure you keep your ID number in a safe place! These envelopes will be passed out by one of your teachers or can be collected from the secretary’s office.

The amount of money you earn depends on your decisions and the decisions of the other participants. We will now describe the rules in detail. It is therefore especially important that you listen very carefully.

There are no “right” or “wrong” decisions in this game. You should make your decisions based on your own personal deliberations. Your decisions will remain anonymous, which means that no one else will know what you decide.

If you have any questions after reading these instructions, please raise your hand. Someone will then come over to you and answer your questions in private (i.e., quietly).

### **Process:**

There are two roles in this game: **sender** and **responder**.

The game starts as follows: Each sender and each responder receives 5 euros. The sender must decide how much of the 5 euros he/she wishes to give to the responder.

The amount the sender gives to “his/her” responder will then be tripled. In other words, the responder receives precisely three times the amount the sender has given him/her.

Next, it is the responder’s turn. He/she now has three times the amount the sender has given him/her plus his/her own 5 euros. The responder must now decide how much of this money he/she would like to return to “his/her” sender. Please note: The sum the responder returns to the sender is not tripled.

### Payment:

At the end of the game, the sender receives the sum that he/she kept plus the sum that the responder returned to him/her.

**Payment to sender = 5 euros - sum sent + sum returned (by responder)**

The responder receives the sum he/she was given by the sender (times 3), minus the sum he/she returned to the sender.

**Payment to responder = 5 euros + 3 x sum sent (by sender) - sum returned**

### Decisions:

You will be required to make one decision in the role of sender and one in the role of responder. You can also choose between different “categories” of senders and responders; you obviously do not have to treat these groups differently, however. These categories are described on the decision sheet. You can, for instance, choose whether you send or return money to a boy or a girl. It is your decision, there is no “right” or “wrong”.

### Calculating your payment:

Some of the following points will be easier to understand once you have seen the decision sheets. We will now go through the points and then look at the decision sheets together. If you still have questions after that, we will be happy come back to these points.

Once the game has been carried out in several schools, the following will happen:

1. Two students from different schools will be randomly paired; you will therefore not know “your” sender or “your” responder personally; however, he or she will be around the same age as you and will also go to school in North Rhine-Westphalia.
2. Who is to play the role of the sender and who the role of the responder will also be randomly decided.
3. Next, we identify the category (see decision sheet) that the sender and responder are each from. This information is extracted from the questionnaire you completed. The sender can be a girl and the responder a boy, for instance.
4. Next, the sender’s decision is implemented based on the actual category of the responder.
5. Finally, the responder’s decision is implemented based on the actual category of the sender and the actual amount received from “their” sender.

6. We now know how much the sender has sent and how much the responder has returned. Based on this, we can calculate the payment to both the sender and the responder. This money is then placed in the appropriate envelopes marked with the corresponding sender and responder ID numbers and taken to the schools.
7. At the end, you will be able to collect the envelope containing your payment at your school.

Now look at the decision sheets. This will help you to better understand some of the points described above. Think carefully about the decisions you wish to make. You have plenty of time! If you have any questions, please raise your hand. Someone will then come over to you and answer your questions in private (i.e., quietly).

**ID:**



**Please KEEP your ID!!!!**

**ID:**





# Receiver 1

You are the receiver. The sender is a boy with German parents. How much do you want to send back to him? Please fill in an amount for each possible case (at most one decimal place = 10-cent steps.)

<i>Assume the sender has sent you the following amount:</i>	<i>The sender still has:</i>	<i>You have:</i>	<i>Which amount do you want to send back:</i>	<i>Potential amount to send back:</i>
0 EURO	5 EURO	5 EURO	_____ EURO	(0 to 5 EURO)
0.5 EURO	4.5 EURO	6.5 EURO	_____ EURO	(0 to 6.5 EURO)
1 EURO	4 EURO	8 EURO	_____ EURO	(0 to 8 EURO)
1.5 EURO	3.5 EURO	9.5 EURO	_____ EURO	(0 to 9.5 EURO)
2 EURO	3 EURO	11 EURO	_____ EURO	(0 to 11 EURO)
2.5 EURO	2.5 EURO	12.5 EURO	_____ EURO	(0 to 12.5 EURO)
3 EURO	2 EURO	14 EURO	_____ EURO	(0 to 14 EURO)
3.5 EURO	1.5 EURO	15.5 EURO	_____ EURO	(0 to 15.5 EURO)
4 EURO	1 EURO	17 EURO	_____ EURO	(0 to 17 EURO)
4.5 EURO	0.5 EURO	18.5 EURO	_____ EURO	(0 to 18.5 EURO)
5 EURO	0 EURO	20 EURO	_____ EURO	(0 to 20 EURO)

## Receiver 2

You are the receiver. The sender is a girl with German parents. How much do you want to send back to her? Please fill in an amount for each possible case (at most one decimal place = 10-cent steps.)

<i>Assume the sender has sent you the following amount:</i>	<i>The sender still has:</i>	<i>You have:</i>	<i>Which amount do you want to send back:</i>	<i>Potential amount to send back:</i>
0 EURO	5 EURO	5 EURO	_____ EURO	(0 to 5 EURO)
0.5 EURO	4.5 EURO	6.5 EURO	_____ EURO	(0 to 6.5 EURO)
1 EURO	4 EURO	8 EURO	_____ EURO	(0 to 8 EURO)
1.5 EURO	3.5 EURO	9.5 EURO	_____ EURO	(0 to 9.5 EURO)
2 EURO	3 EURO	11 EURO	_____ EURO	(0 to 11 EURO)
2.5 EURO	2.5 EURO	12.5 EURO	_____ EURO	(0 to 12.5 EURO)
3 EURO	2 EURO	14 EURO	_____ EURO	(0 to 14 EURO)
3.5 EURO	1.5 EURO	15.5 EURO	_____ EURO	(0 to 15.5 EURO)
4 EURO	1 EURO	17 EURO	_____ EURO	(0 to 17 EURO)
4.5 EURO	0.5 EURO	18.5 EURO	_____ EURO	(0 to 18.5 EURO)
5 EURO	0 EURO	20 EURO	_____ EURO	(0 to 20 EURO)

## Receiver 3

You are the receiver. The sender is a boy with foreign parents. How much do you want to send back to him? Please fill in an amount for each possible case (at most one decimal place = 10-cent steps.)

<i>Assume the sender has sent you the following amount:</i>	<i>The sender still has:</i>	<i>You have:</i>	<i>Which amount do you want to send back:</i>	<i>Potential amount to send back:</i>
0 EURO	5 EURO	5 EURO	_____ EURO	(0 to 5 EURO)
0.5 EURO	4.5 EURO	6.5 EURO	_____ EURO	(0 to 6.5 EURO)
1 EURO	4 EURO	8 EURO	_____ EURO	(0 to 8 EURO)
1.5 EURO	3.5 EURO	9.5 EURO	_____ EURO	(0 to 9.5 EURO)
2 EURO	3 EURO	11 EURO	_____ EURO	(0 to 11 EURO)
2.5 EURO	2.5 EURO	12.5 EURO	_____ EURO	(0 to 12.5 EURO)
3 EURO	2 EURO	14 EURO	_____ EURO	(0 to 14 EURO)
3.5 EURO	1.5 EURO	15.5 EURO	_____ EURO	(0 to 15.5 EURO)
4 EURO	1 EURO	17 EURO	_____ EURO	(0 to 17 EURO)
4.5 EURO	0.5 EURO	18.5 EURO	_____ EURO	(0 to 18.5 EURO)
5 EURO	0 EURO	20 EURO	_____ EURO	(0 to 20 EURO)

## Receiver 4

You are the receiver. The sender is a girl with foreign parents. How much do you want to send back to her? Please fill in an amount for each possible case (at most one decimal place = 10-cent steps.)

<i>Assume the sender has sent you the following amount:</i>	<i>The sender still has:</i>	<i>You have:</i>	<i>Which amount do you want to send back:</i>	<i>Potential amount to send back:</i>
0 EURO	5 EURO	5 EURO	_____ EURO	(0 to 5 EURO)
0.5 EURO	4.5 EURO	6.5 EURO	_____ EURO	(0 to 6.5 EURO)
1 EURO	4 EURO	8 EURO	_____ EURO	(0 to 8 EURO)
1.5 EURO	3.5 EURO	9.5 EURO	_____ EURO	(0 to 9.5 EURO)
2 EURO	3 EURO	11 EURO	_____ EURO	(0 to 11 EURO)
2.5 EURO	2.5 EURO	12.5 EURO	_____ EURO	(0 to 12.5 EURO)
3 EURO	2 EURO	14 EURO	_____ EURO	(0 to 14 EURO)
3.5 EURO	1.5 EURO	15.5 EURO	_____ EURO	(0 to 15.5 EURO)
4 EURO	1 EURO	17 EURO	_____ EURO	(0 to 17 EURO)
4.5 EURO	0.5 EURO	18.5 EURO	_____ EURO	(0 to 18.5 EURO)
5 EURO	0 EURO	20 EURO	_____ EURO	(0 to 20 EURO)

## Receiver 5

You are the receiver. The sender is a boy with foreign parents who possesses German citizenship. How much do you want to send back to him? Please fill in an amount for each possible case (at most one decimal place = 10-cent steps.)

<i>Assume the sender has sent you the following amount:</i>	<i>The sender still has:</i>	<i>You have:</i>	<i>Which amount do you want to send back:</i>	<i>Potential amount to send back:</i>
0 EURO	5 EURO	5 EURO	_____ EURO	(0 to 5 EURO)
0.5 EURO	4.5 EURO	6.5 EURO	_____ EURO	(0 to 6.5 EURO)
1 EURO	4 EURO	8 EURO	_____ EURO	(0 to 8 EURO)
1.5 EURO	3.5 EURO	9.5 EURO	_____ EURO	(0 to 9.5 EURO)
2 EURO	3 EURO	11 EURO	_____ EURO	(0 to 11 EURO)
2.5 EURO	2.5 EURO	12.5 EURO	_____ EURO	(0 to 12.5 EURO)
3 EURO	2 EURO	14 EURO	_____ EURO	(0 to 14 EURO)
3.5 EURO	1.5 EURO	15.5 EURO	_____ EURO	(0 to 15.5 EURO)
4 EURO	1 EURO	17 EURO	_____ EURO	(0 to 17 EURO)
4.5 EURO	0.5 EURO	18.5 EURO	_____ EURO	(0 to 18.5 EURO)
5 EURO	0 EURO	20 EURO	_____ EURO	(0 to 20 EURO)

## Receiver 6

You are the receiver. The sender is a girl with foreign parents who possesses German citizenship. How much do you want to send back to her? Please fill in an amount for each possible case (at most one decimal place = 10-cent steps.)

<i>Assume the sender has sent you the following amount:</i>	<i>The sender still has:</i>	<i>You have:</i>	<i>Which amount do you want to send back:</i>	<i>Potential amount to send back:</i>
0 EURO	5 EURO	5 EURO	_____ EURO	(0 to 5 EURO)
0.5 EURO	4.5 EURO	6.5 EURO	_____ EURO	(0 to 6.5 EURO)
1 EURO	4 EURO	8 EURO	_____ EURO	(0 to 8 EURO)
1.5 EURO	3.5 EURO	9.5 EURO	_____ EURO	(0 to 9.5 EURO)
2 EURO	3 EURO	11 EURO	_____ EURO	(0 to 11 EURO)
2.5 EURO	2.5 EURO	12.5 EURO	_____ EURO	(0 to 12.5 EURO)
3 EURO	2 EURO	14 EURO	_____ EURO	(0 to 14 EURO)
3.5 EURO	1.5 EURO	15.5 EURO	_____ EURO	(0 to 15.5 EURO)
4 EURO	1 EURO	17 EURO	_____ EURO	(0 to 17 EURO)
4.5 EURO	0.5 EURO	18.5 EURO	_____ EURO	(0 to 18.5 EURO)
5 EURO	0 EURO	20 EURO	_____ EURO	(0 to 20 EURO)