

Collaboration and Homophily in Global Teams

Gábor Békés^{a,*}, Gianmarco I.P. Ottaviano^b

^a*Central European University, KRTK and CEPR*

^b*Bocconi University, Baffi-CAREFIN, IGIER, CEP and CEPR*

Abstract

How do barriers related to nationality and language affect collaboration in multinational teams? We address this question by exploiting a newly assembled exhaustive dataset recording all 10.7 million passes by professional European football players from 132 countries fielded by all 154 teams competing in the top five men leagues over eight sporting seasons, together with full information on players' and teams' characteristics. We measure collaboration as the average number of passes per minute between a pair of players in a half season. We use a discrete choice model of players' passing behavior as a baseline to separately identify excess collaboration within nationality or language due to preferences ('choice homophily') from collaboration due to opportunities ('induced homophily'). Our dataset allows us to estimate the model using a rich set of player and play characteristics as well as player fixed effects. We find strong evidence of homophily: conditioning on players' and teams' characteristics, player pairs of same nationality exhibit an average number of passes per minute that is 2.8 percent higher than player pairs of different nationality. Same nationality is about as likely to lead to more passes as doubling the player pair's valuation, which is a consensus measure of players' skills. Shared language has about half the impact of same nationality. Pairs of same nationality are also more likely to engage in deeper collaboration, disproportionately participating in more complex pass sequences. These findings show that homophily based on nationality and language is pervasive even in teams of very high skill individuals with clear common objectives and aligned incentives, and involved in interactive tasks that are well-defined and not particularly language-intensive.

Keywords: Homophily, organizations, teams, diversity, big data

Preliminary and incomplete: 12 October 2021. See latest version [HERE](#).

*Corresponding author, mail: bekesg@ceu.edu, Central European University, Quellenstrasse 51, Vienna, Austria.

**Bekes thanks the support of the 'Firms, Strategy and Performance' Lendület grant of the Hungarian Academy of Sciences. We thank Gabor Kezdi, Balazs Murakozy, Miklós Koren and seminar participants at CERS-HAS, University of Reading, IMT Lucca, and CEU for useful comments and

1. Introduction

To compete in the global economy, companies are increasingly relying on a multinational workforce (Neeley, 2015). This allows them to build teams that offer the best expertise from around the world, and draw on the benefits of international diversity by bringing together people from many cultures with varied work experiences and perspectives. Teams like these, however, also face several hurdles. When team members come from different countries and cultural backgrounds, communication can rapidly deteriorate, misunderstanding can ensue, and cooperation can degenerate into distrust.

In this paper we systematically investigate how barriers related to nationality and language diversity affect collaboration in multinational teams (Lazear, 1999*b,a*). We define ‘collaboration’ as the situation of two or more people working together to create or achieve the same thing (Cambridge Dictionary), and study teams that are not geographically dispersed, as dispersion may *per se* inhibit collaboration (Joshi et al., 2002). We show that a ‘border effect’ between team members of different nationality or language may indeed hamper collaboration, pretty much like it hampers international trade in goods and services between the regions of different countries (Head and Mayer, 2014). Team members of the same nationality or language will collaborate more than team members of different nationality or language.

We base our investigation on the unique features of a newly assembled dataset recording all passes made by professional European football players in the top five men leagues (England, France, Germany, Italy, Spain) over eight sporting seasons (2011-12 to 2018-19) together with full information on players’ and teams’ characteristics as well as their performances.¹

We build a unique dataset recording all the 10.7 million passes in 14,608 games by players from 132 countries fielded by 154 teams. To measure same nationality difference, we aggregate our data to time periods with stable squads of players - half of seasons. Estimations are carried out for 355 thousand player pairs * half seasons.

We measure collaboration as the average number of passes per minute between a pair of players in a given time period (which we call their ‘pass rate’), and study how it is affected by the pairs’ nationality and language. Passes are the essential building blocks of football. They represent how players work together for the common objective of scoring or preventing the opponent from scoring a goal. Importantly, passes are positively correlated with winning.²

This type of sports data has several advantages. First, the European football industry is very globalized: fans are spread around the world, and multinational teams

suggestions. We are grateful to Endre Borza, Bence Szabo for outstanding and extensive research assistance.

¹With ‘European football’, or simply ‘football’ henceforth, we refer to ‘association football’, which is commonly known as ‘football’ in Europe and ‘soccer’ in the United States (Tovar, 2020) .

²See Appendix D

are the rule in the top five leagues³. Second, football players are very mobile internationally, and their mobility decisions are typically made for work-related reasons, with pay being the most prominent of them.

Third, in the top five leagues players are very diverse in terms of origin as they come from over a hundred countries. At the same time, they are all very high skilled workers hardly facing obstacles with integration outside the workplace. Moreover, while language matters for collaboration, football tasks are not particularly language-intensive (Nüesch and Haas, 2013). Fourth, all sorts of player as well as team characteristics and performance indicators are precisely measured, and fastidiously recorded. Moreover, extensive media coverage can be readily used to shed light on any odd data patterns. Fifth, while team composition is exogenous to players' decisions, collaboration with other team members is mostly up to their individual choices. Sixth, the 'rules of the game' are codified, and crystal clear to players and teams⁴

All these features allow us to investigate collaboration in competitive global teams of high skilled workers with precise common objectives, leveraging a big dataset on interactions in an actual workplace rather than in an artificial experimental lab (Jackson et al., 2003), while exploiting an extremely rich set of team and worker controls.

We are not the first to exploit team sports data to analyze the potential gains and losses from employing culturally diverse work teams. In the case of the National Hockey League in North America, (Kahane et al., 2013) find that the presence of European players (with Europe being the typical origin of foreign players) does increase firm-level performance: teams that employ a higher proportion of European players perform better. However, their results also indicate that teams perform better when their European players come from the same country rather than being spread across many European countries. When teams have players from a wide array of European countries, integration costs associated with language and cultural differences may start to override any gains from diversity. Parallel evidence based on European football leads to mixed conclusions. In the top German league multinational teams have been found to perform worse than teams with less national diversity (Nüesch and Haas, 2013), whereas the opposite has been found in the top continental tournament (Ingersoll et al., 2017). Studying the top leagues of England and Spain, Tovar (2020) suggests that conflicting results may derive from a hump-shaped relation between team performance and predominant nationality. This echoes (Kahane et al., 2013) in that an optimal degree of diversity may exist. What distinguishes our analysis from these and related works is that we zoom in on collaboration and we can measure it accurately through the pass data.

³On average teams in our sample have a squad of players from 13 countries and field a starting XI with players from 6 countries.

⁴This last point is relevant because it is not that some nationality has an advantage in understand and exploiting unwritten/unclear rules, and thus it passes more independently of collaboration propensity.

The key methodological challenge that our investigation faces has been highlighted in the studies on homophily, defined as the tendency to associate with similar others (Lawrence and Shah, 2020). That team members of the same nationality or language collaborate more than team members of different nationality or language is a statement about homophily. It highlights common nationality or language as the antecedents of homophily, that is, the specific attributes that serve as its basis, while singling out collaboration as the targeted consequence of homophily (Ertug et al., 2021). In this respect, in studying homophilous behavior an important distinction has been made between two underlying mechanisms: opportunities and preferences ((McPherson and Smith-Lovin, 1987; McPherson et al., 2001).

According to the former mechanism, individuals’ distributions across categories within a social context define the probability they choose similar others (Lawrence and Shah, 2020). This may mechanically ‘induce’ homophily, irrespective of whether players have any actual preference for similar others, and thus it may not tell much about their real tendency to associate with similar others. Lawrence and Shah (2020) offer the following simple example. Consider a group of 100 geoscientists who associate with one another during a conference workshop. If 40 percent are geochemists and 60 percent are hydrologists, the expected rate for geochemists associating with other geochemists is 0.40. Only when the proportion of geochemists’ associations with other geochemists exceeds this baseline, it demonstrates a preference for geochemists to associate with other geochemists. It is this preference that distinguishes ‘choice homophily’ from ‘induced homophily’. Hence, to be of any interest at all, the statement that team members of the same nationality or language collaborate more has to be based on choice homophily after controlling for induced homophily.

Defining the baseline is quite straightforward in the previous example. It is much less so when individuals may or may not differ along several potential attributes that could confound the roles of the targeted antecedents of homophily, making it harder to ascertain whether individuals are mechanically induced to choose similar others. We address these identification issues by designing the baseline in terms of a discrete choice model of players’ passing behavior. The model determines how the pass rate for a pair of players is pinned down by their characteristics and opportunities during the matches they play together in a given time period. It is implemented empirically by a Poisson regression with a variety of player characteristics as controls. Results are then corroborated by a rich set of robustness checks.

We find that player pairs of same nationality do pass more between them. Conditioning on observable player characteristics (such as team, position, valuation, citizenship), pass features (such as average distance), player pairs of same nationality tend to have a passing rate (pass per minute) of 2.8 percent higher than player pairs of different nationality. Same nationality is about as likely to lead to more passes as doubling the player pair’s valuation (a consensus measure of their skills). Same language has about half the impact of same nationality, suggesting that same nationality has multiple components, with language being only one of them. Pairs of same nationality also engage in

deeper collaboration, participating in complex pass sequences. These include two-pass sequences (ABA or BAB for players A and B), or longer ones. For such sequences, player pairs of same nationality tend to have a passing rate of 5 percent higher than player pairs of different nationality. Shared experience does not affect these results. Once individual experience is controlled for, shared experience has no relevant effect. This suggests that nationality is not a proxy for knowing each other.

These findings show that homophily based on nationality and language is pervasive even in teams of very high skill individuals with clear common objectives and aligned incentives, and involved in interactive tasks that are well-defined and not particularly language-intensive.

The rest of the paper is organized as follows. Section 2 offers a selective overview of the related literature beyond works already referenced in this introduction. Section 3 describes data collection and our dataset. Section 4 introduces the discrete choice model of passing behavior. Section 5 presents the model estimation, whose results are then discussed in Section 6. Section 7 concludes.

2. Related literature

This paper is related to various research streams of the vast literature on diversity and performance in teams, which spans from management (see e.g. Earley and Mosakowski (2000) and (Jackson et al., 2003)) to education studies (see e.g. Terenzini et al. (2001)).

Three streams are particularly relevant to what we do. The first is concerned with ‘diversity spillovers’, which improve team performance in a diverse environment, but not necessarily in a team that is itself diverse (Ottaviano and Peri, 2006, 2005). This stream highlights four main mechanisms (Buchholz, 2021). Diversity increases productivity: (i) when people from different countries work on problems together, in turn identifying better solutions by combining their knowledge (‘interactive problem-solving’), (ii) through increasing the specialization, variety of skills and approaches to tasks within an occupation, though without necessarily requiring interaction between people from different countries of birth (‘complementary task specialization’), (iii) when people from the same country of birth cluster in particular occupations and this clustering facilitates stronger knowledge exchanges (‘niching effects’), (iv) when simply through exposure to a diverse range of knowledge and approaches to problems workers learn and become more productive (‘exposure effects’). The evidence on US Metropolitan Statistical Area reported in Buchholz (2021) supports exposure effects as the main mechanism, but also interactive problem-solving and complementary task specialization seem to play an important role.

This first stream does not leverage information on diversity and collaboration within teams, which is what we do. In this respect, our investigation is more closely related to a second research stream that studies how individuals of different ethnicities may complement each other in production, but workers of the same ethnic background

may collaborate more effectively (Lazear, 1999*b,a*; Lang, 1986). Specifically related to our investigation are works highlighting how distortions due to ethnic diversity and discriminatory worker attitudes affect firms and their organization of production. These studies face stiff data challenges. To systematically examine the effects of culture and language within a firm, one needs a host of detailed data: the nationalities of all workers must be identifiable, each worker’s skills and output, as well as the collective output of the firm, must be measurable, and all other factors of production should be held constant (Kahane et al., 2013). That is why works on firms are typically based on field experiments (Bertrand and Duflo, 2017). For instance, Hjort (2014) studies team production at a plant in Kenya, where an upstream worker supplies and distributes flowers to two downstream workers, who assemble them into bunches. He finds that upstream workers undersupply non-coethnic downstream workers (vertical discrimination) and shift flowers from non-coethnic to coethnic downstream workers (horizontal discrimination), at the cost of lower own pay and total output. Team pay, whereby the two downstream workers are remunerated for their combined output, is shown to mitigate discrimination and its allocative distortions.⁵

In Hjort (2014), the upstream worker’s decision on distributing flowers to the downstream workers resembles the choice a football player faces on passing the ball to his teammates. The context is, however, quite different. Whereas a Kenyan plant is a low skilled, highly charged context in a developing country with ethnic conflicts, a European football team is a high skilled, lowly charged context in a developed area with no real conflicts. Moreover, the flower plant and the football team setups have different pros and cons. The former can exploit an essentially random rotation process to assign workers to positions for identification, but its external validity may be limited. In the latter setup rotation is arguably not random as it depends on the manager’s choices, but the richness of information from which to obtain all sorts of individual and team controls makes the case for external validity stronger. Be it as it may, non-random rotation due to endogenous team formation leads to known biases. Calder-Wang et al. (2021) exploit a dataset of MBA students who participated in a required course to propose and start a real micro-business that allows them to examine horizontal diversity (i.e., within the team) as well as vertical diversity (i.e., team to faculty advisor) and their effect on performance. The course was run in multiple cohorts in otherwise identical formats except for the team formation mechanism used. In several cohorts, students were allowed to choose their teams from among students in their section. In other cohorts, students were randomly assigned to teams based upon a computer algorithm. In the cohorts that were allowed to choose, Calder-Wang et al. (2021) find strong selection

⁵Conflicts exacerbate discrimination. Hjort (2014) finds that a period of ethnic conflict following Kenya’s 2007 election led to a sharp increase in discrimination at the flower plant. Using microdata from GitHub, the world’s largest hosting platform for software projects, Laurensyeva (2019) finds that political conflict that burst out between Russia and Ukraine reduced online cooperation between Russian and Ukrainian programmers.

based upon shared attributes.⁶ Among the randomly-assigned teams, greater diversity along the intersection of gender and race/ethnicity significantly reduced performance. However, the negative effect of this diversity is alleviated in cohorts in which teams are endogenously formed. In this respect, as long as the manager of a football team acts as mediator allowing the team to internalize the effects of diversity, the negative impact of diversity on collaboration we find can be seen as a lower bound estimate.

The third research stream analyzes homophily in scientific publications. Looking into scientific papers written by US-based authors from 1985 to 2008, Freeman and Huang (2015) find evidence of choice homophily as persons of similar ethnicity co-author together more frequently than predicted by their proportion among authors; and that greater homophily is associated with publication in lower impact journals and with fewer citations, even holding fixed the authors' previous publishing performance. By contrast, diversity in inputs by author ethnicity, location, and references leads to greater contributions to science as measured by impact factors and citations. In the same vein, AlShebli and Woon (2018) study the relationship between research impact and five classes of diversity: ethnicity, discipline, gender, affiliation, and academic age. Using randomized baseline models, they establish the presence of homophily in ethnicity, gender and affiliation. However, ethnic diversity has the strongest correlation with scientific impact. To further isolate the effects of ethnic diversity, they use randomized baseline models and again find a clear link between diversity and impact. Differently from these studies, we use a discrete choice model rather than randomized models to separate choice homophily from induced homophily.

3. Data

Our dataset consists of all passes made by professional European football players in the top five men leagues over eight sporting seasons, together with full information on players' and teams' characteristics as well as their performances.⁷ The top five leagues are the German Bundesliga, the French Ligue 1, the Spanish La Liga, the English Premier League, and the Italian Serie A. We have selected these leagues because of their undisputed preeminence as the pinnacle of national football competitions. Moreover, for these leagues data availability is the most comprehensive.

⁶Currarini et al. (2009) study friendship formation in US schools when students have types and may see type-dependent benefits from friendships. They show that any matching process such that types are matched in frequencies in proportion to their relative stocks cannot replicate the generalized inbreeding they observe in Add Health 1994 Data, regardless of type sensitivity of preferences. On the contrary, a model with both type-sensitive preferences and a matching bias generates the observed patterns of inbreeding. See also Currarini et al. (2010).

⁷Data come from private sources and cannot be made publicly available. We cannot print or share any individual piece of information on games. We cannot share the data publicly. However, all codes will be shared together with a smaller dataset to allow checking them and even running queries. Furthermore, we will offer access to run scripts on the full dataset on a secure server. We are currently building the data infrastructure.

The dataset covers all games played in sporting seasons 2011-12 to 2018-19, for which data quality is the highest. A season is the time period between mid-August to mid-May, during which each team plays twice (home and away) with every other team. The Premier League, La Liga, Serie A and Ligue 1 are all composed of 20 teams (playing $20 \times 19 = 380$ games) while there are 18 teams ($18 \times 17 = 306$ games) in the Bundesliga. In any given season, there are 98 teams in our sample. Due to relegation and promotion, we have a total of 154 teams in the sample. Overall, our dataset covers a total of $8 \times (380 \times 4 + 306) = 14,608$ games.

3.1. Players

We have 7,048 players in our sample, for whom we can fully map their entire career, with a typical team relying on a squad of about 30 players. Player information and characteristics are compiled by Transfermarkt (<https://www.transfermarkt.com/>).⁸ They include country of birth, citizenship or citizenships if multiple, date of birth, height, and participation in a national team, which are time-invariant in our dataset. They also include a player’s estimated transfer value, that is, the ”expected value of a player in a free market” as determined by a group of experts. This estimate is based on how much a player may contribute to the team’s success, how well he plays, how useful he may be to another team that would purchase the rights, and so on. As such, a player’s transfer value is considered a consensus measure of the quality of his footballing skills. Transfer values are estimated twice a year in correspondence of the transfer windows.

European football is truly globalized and there are players from 132 countries of citizenship in our sample. French, Spanish and Italian players are the largest citizenship groups, followed by Germans, English, Brazilians and Argentinians. Other countries of citizenship with several players include the Netherlands, Serbia, Senegal, and Uruguay. Table 1 reports the shares of the first nationality for countries with at least a 1% share.

To determine common nationality between players in the presence of multiple citizenships, we define two players as co-national if they share at least one citizenship, or have the same country of birth.

3.2. Passes

Information on passes comes from OPTA and is available via third party sites.⁹ A pass is an event defined as ”any pass attempted from one player to another”, including free kicks, corners, throw-ins, goal assists. The feed data include a timestamp, team id, x and y coordinates of where the pass took place in the pitch, and its outcome (e.g., a flag for a successful pass). Using the next event, the destination player and his coordinates can be recovered.¹⁰

⁸Several small data cleaning and measurement issues are discussed in Appendix B.

⁹OPTA’s data are generated by machines and humans coding events during games. The data has been scraped from a third party website. [MORE]

¹⁰A player can *de facto* also pass to himself. We have dropped these self-passes from the data.

Table 1: Most frequent player nationalities

Country	share (%)
Spain	14.1
France	12.6
Italy	10.5
Germany	8.6
England	6.8
Brazil	4.4
Argentina	3.9
Portugal	1.8
Senegal	1.6
Netherlands	1.6
Belgium	1.6
Serbia	1.4
Switzerland	1.3
Cote d'Ivoire	1.3
Croatia	1.1
Uruguay	1.1

There is 17 passes per minute, on average, and there are about 800 successful passes on average per game, so we have 10.73 million passes.

We have first aggregated the pass data to single games to generate variables such as the sum of passes between any two players in a game. This has given us 1.857 million observations at *player – pair × game* level. We have then aggregated these observations by half-season after dividing seasons into two halves, before and after January 1st. This division in half-seasons is suggested by the timing of the transfer windows, which are located between seasons (summer transfer window) and at the beginning of the calendar year (winter transfer window). It also splits the number of games during a season into two approximately equal halves: the number of games per team in a half-season ranges between 16 and 20 compared with the exact equal split of 17 for the German Bundesliga and 19 for the other top five leagues. This gives us 355,610 observations at *player – pair × half – season* level, which is the level of aggregation we will use for estimation (see Table 2).

Table 2: Number of observations at different aggregation levels

source	level of aggregation	N obs (million)
raw data	player-pair * pass events	10.73
game level	player-pair * game	1.86
regressions	player-pair * half seasons	0.36

Half-seasons have several advantages with respect to alternatives such as seasons, games, or half-games. In a game or half-game two players may not play together for various reasons as squads are large and only eleven players can be fielded at the same time. This may raise a selection issue, which can be tackled by considering all the games in a season instead. If two players never pass to each other during an entire season, either they are never fielded together or one can safely assume they are never fielded in positions that interact. Moreover, considering a season allows one to investigate the role of common experience as players who spend more time together on the pitch may learn to pass more.

On the other hand, the presence of the winter transfer window implies that during a season a team’s squad may change composition. Our assumption of unchanged player quality makes more sense in a half-season than in a season, especially younger players may evolve. In this respect, a half-season strikes a balance between mitigating the selection issue and keeping squad composition fixed. Finally, the fact that half-seasons are separated by transfer windows allows us to cleanly map players’ careers as they change teams, thereby combining player and pass information in a consistent way.

Under the reasonable assumption that, if two players never pass to each other during a half-season, it must be that it is impossible for them to do so due to fielding or positioning reasons, we drop the corresponding player pairs. There are thus no zero passes for any player pairs in our dataset.

To match player and pass data, entity resolution for players has been an essential and difficult task. We used a matching algorithm based on player names and additional information. The procedure is detailed in Appendix B. We have also made a few decisions regarding data cleaning, such as dropping players who only had a single passing partner. All results are robust to them.

4. A Discrete Choice Model of Passing Behavior

A crucial challenge in assessing how common nationality and common language affect collaboration through passes arises from the conflation of choice and opportunity. As discussed in the introduction, individuals may collaborate more with similar others because they choose to do so (‘choice homophily’), or because collaboration with similar others is forced on them by unrelated circumstances (‘induced homophily’). In this section we develop a discrete choice model to help us disentangle choice from opportunity in an internally consistent way by controlling for observable player characteristics (such as team, position, valuation, citizenship) and pass features (such as average distance).

Consider a football team of $N = 11$ players, indexed from 1 to N , engaged in a half-season consisting of P passing episodes.¹¹ During the half-season each player is assigned to a particular position on the pitch, which implies that a player’s index

¹¹The model can be extended to allow for $N > 11$ and different selections of players being fielded along the half-season.

identifies both his name and his position. Let us focus on two players, labeled o and d , and on the subset of passing episodes $T^{o,d}$ in which both players are on the pitch with player o having ball possession. A ‘pass’ from o to d is defined as a movement of the ball determined by a decision made by player o (‘passer’) to kick or throw the ball to teammate d (‘receiver’). For $d = o$ the passer keeps possession of the ball. We are interested in characterizing the probability that player o passes to player d rather than to any of the other nine teammates.

A passing episode consists of two periods: when the pass is initiated by o (t) and when the pass is received by d ($t + 1$). The passer wants to maximize team payoff and understands that the benefit for the team of one of its players controlling the ball is determined by the ability and position of that player, and by some randomness due to the vagaries of the game. These may include, for instance, the performance of the opposing team, the referee’s decisions or the weather conditions. We use $\ln u_t^d$ to denote the deterministic part of the team’s benefit as determined by player d ’s characteristics, and z_t^d to denote its random part (‘shock’) due to match contingencies. For each receiver this shock is the realization of a random variable Z with continuous differentiable c.d.f. $\Pr[Z \leq z] = G(z)$ over the support $(-\infty, +\infty)$. Any difference in outcomes across the $T^{o,d}$ episodes depends on this shock’s realizations only.

Passer o also understands the challenges he faces in passing the ball to receiver d . We call $\tilde{c}^{o,d}$ the associated ‘passing cost’ capturing such challenges. In particular, this cost may be high if o and d find it hard to collaborate due to different nationality (or language) or because the pass is difficult due to distance or the position of players.

Moreover, passer o realizes the difficulty receiver d may face in taking control of the ball, which depends on the receiver’s characteristics. We use φ^d to denote the probability that receiver d takes control of the ball. We call this the probability of a successful pass.

The passer’s decision can be characterized as the problem of passing the ball to the receiver who generates the highest expected benefit for the team. The value function of this problem can be written recursively as

$$U_t^o = \ln u_t^o + \beta \max_{\{d\}_{d=0}^N} \{ \varphi^d E[U_{t+1}^d] - \tilde{c}^{o,d} + z_t^d \}. \quad (1)$$

According to (1) the team’s benefit U_t^o of controlling the ball in period t is split into two components: the benefit of player o currently controlling the ball (e.g. the player could try to score a goal; or, with the player in control of the ball, the opposing team cannot score) and the option value of player o passing (or keeping control of) the ball at the beginning of the future period. These two components are captured by $\ln u_t^o$ and $\beta \max_{\{d\}_{d=0}^N} \{ \varphi^d E[U_{t+1}^d] - \tilde{c}^{o,d} + z_t^d \}$ respectively, with $\beta \in [0, 1]$ measuring the importance attributed to the latter option and expectation $E[U_{t+1}^d]$ being taken over future realizations of the shock.

Assuming that Z follows the Gumbel distribution (Type-I Extreme Value distribu-

tion)

$$G(z) = \exp(-\exp(-\kappa z))$$

leads to a simple expression for the probability of player o passing to teammate d in period t . Specifically, after taking expectations on both sides of (1), defining $\tilde{V}_t^o \equiv \varphi^o E[U_t^o]$, $V_t^o \equiv \exp \tilde{V}_t^o$ and $c^{o,d} \equiv \exp \tilde{c}^{o,d}$ allows us to obtain

$$V_t^o = u_t^o (\Lambda_t^o)^\beta \text{ with } \Lambda_t^o = \left[\sum_{d=1}^N (V_{t+1}^d)^{\varphi^{dk}} (c^{o,d})^{-k} \right]^{\frac{1}{\kappa}} \quad (2)$$

which expresses the value of ball possession by player o as the (exponential of the) value of him controlling the ball in period t (u_t^o) adjusted for the option value of passing the ball at the beginning of period $t+1$ (Λ_t^o) weighted by the relative importance attributed to this option (β). Based on (2), *ex ante* the probability that player o in possession of the ball in period t successfully passes to teammate d at the beginning of period $t+1$ is given by

$$\pi_t^{o,d} = \frac{(V_{t+1}^d)^{\varphi^{dk}} (c^{o,d})^{-k}}{\sum_{d=1}^N (V_{t+1}^d)^{\varphi^{dk}} (c^{o,d})^{-k}} = (\Lambda_t^o)^{-k} (V_{t+1}^d)^{\varphi^{dk}} (c^{o,d})^{-k} \quad (3)$$

which *ex post* becomes approximately the average share of successful passes that player o makes to player d per episode over a half-season in the subset of passing episodes $T^{o,d}$ when both o and d are fielded and player o has ball possession.¹² The probability that player o successfully passes the ball to player d in period t is an increasing function of the expected discounted payoff for the team if the ball is passed to player d (V_{t+1}^d), adjusted for the probability this player takes control (φ^d), and a decreasing function of the opportunity cost for the team if player o passes the ball rather than keeping it under control. This opportunity cost itself has two components: the ‘passing cost’ from o to d ($c^{o,d}$) and the option value of o keeping the ball (Λ_t^o).

In the data we observe the total number of team passing episodes (P), the number of passing episodes involving a pass from o to d ($P^{o,d}$), and the number of passing episodes when both o and d are fielded and player o has ball possession ($T^{o,d}$) over a half-season. If we define the half-season ‘pass rate’ as $p^{o,d} = P^{o,d}/P$, the model then implies $p^{o,d} = T^{o,d} \pi_t^{o,d}/P$, and thus

$$\log p^{o,d} = \log T^{o,d} + \log (\Lambda_t^o)^{-k} + \log (V_{t+1}^d)^k + \log (c^{o,d})^{-k} - \log P \quad (4)$$

This equation can be estimated using our data by specifying the bilateral passing cost multiplicatively as

$$c^{o,d} = (g^{o,d})^\gamma (l^{o,d})^\lambda \quad (5)$$

¹²The fact that also o is included in the sum $\sum_{d=1}^N (V_{t+1}^d)^k (c^{o,d})^{-k}$ implies $\sum_{d=1}^N \pi_t^{o,d} = 1$.

In (5) $g^{o,d}$ is the physical distance between the two players' positions so that $(g^{o,d})^\gamma$ captures all distance-related frictions that make it difficult to pass the ball independently of the identities of the passer and the receiver. The term $(l^{o,d})^\lambda$ captures, instead, all non-distance-related frictions that make it difficult for player o to pass to receiver d , independently from the positions they are assigned. These non-distance-related frictions may include, for instance, different cultural traits or limited experience in playing together.

In particular, we will measure $g(o, d)$ with *PassDist*, the average distance between passing players. $l^{o,d}$ will allow us measure homophily with *SameNatIndicator* measuring shared cultural background.

With double player fixed effects, this leads to the following regression

$$\Pi_h^{o,d} = I_h^o + I_h^d + \alpha \text{SameNatIndicator}^{o,d} + \gamma \text{PassDist}^{o,d} + I_h^{team} + \Upsilon_h^{o,d} + \varepsilon_h^{o,d} \quad (6)$$

where h is the half-season label, $\Pi_h^{o,d}$ is log pass rate, I_h^o and I_h^d are passer and receiver quality measures or fixed effects, I_h^{team} is a team fixed effect (soaking up coach, style, pitch features), $\Upsilon_h^{o,d}$ is log $T^{o,d}$, and $\varepsilon_h^{o,d}$ is an error term. We have $\gamma \text{PassDist}^{o,d}$, a measure of pass difficulty to capture $g(o, d)$.

As for $\text{SameNatIndicator}^{o,d}$, this is a time-invariant indicator variable equal to 1 if players o and d have the same nationality, and 0 otherwise. In this case, an estimated $\alpha_l > 0$ would measure a 'homophily premium' in passing between the two players. In this respect, the foregone homophily premium in passing between players of different nationalities would be analogous to the 'border effect' found in gravity regression for trade flows between regions belonging to different countries.

5. Model estimation

We implement our discrete choice model empirically by running a regression of the 'pass rate' for a player-pair on the same nationality indicator variable (*SameNatIndicator* in the model).

For an easier interpretation, our baseline estimate makes two assumptions to have a simpler setting. First, we assume that players have a stable pass per minute rate in a season-half. This means that we can replace the number of passes to any player when two players are on the pitch with minutes spent together. Second, we assume that that quality coefficients are similar, and can treat a pass from A to B together with passes with B to A. In extensions, both assumptions are relaxed with very little impact on results.

Thus, the pass rate is defined as the number of passes ('count') divided by the number of minutes played together for a player-pair ('exposure') and the total number of passes in their teams during a half-season. Specifically, we use a Poisson model that converts the pass rate into the pass count by multiplying both sides of the estimating

equation by exposure, taking logs and thus featuring $\log(\text{exposure})$ as a term added to the regression coefficients (offset). We then use $\text{team} \times \text{half_season}$ fixed effects to absorb the total number of passes per team in a half-season.

Our theoretical model needs a generalized linear estimator with a log link function. In such setup, there is a large literature on the benefits of preferring a Poisson model to $\log(\text{count})$ model with a large number of fixed effects.¹³ In particular, a drawback of fixed effect models in general is the incidental parameter bias: having several nuisance parameters to estimate, the estimated coefficient of the variable of interest may be biased.¹⁴ While the Poisson model seems the best choice for us, robustness checks will be provided with the $\log(\text{count})$ as outcome.

In detail, let us index player 1, player 2, and time period (half-season) by $p1$, $p2$, and t respectively. The pass count between player 1 and 2 in period t is then given by $\text{pass_count}_{p1,p2,t}$, while our variable of interest is denoted by $x_{p1,p2,t}$ with potential confounding variables $z_{p1,p2,t}$. Minutes shared by the players on the pitch is denoted by $\ln \text{minutes_shared}$ and, as an exposure variable, the corresponding coefficient is set to unity. Accordingly, the log pass rate is defined as:

$$\text{pass_rate}_{p1,p2,t} = \left(\frac{\text{pass_count}_{p1,p2,t}}{\text{minutes_shared}_{p1,p2,t}} \right)$$

We estimate two versions of our model. First, we consider a version of the Poisson model with a variety of player characteristics as controls:

$$\begin{aligned} \mu_{p1,p2,t} := E(\text{pass_count}_{p1,p2,t} | \dots) = \exp(\gamma g_{p1,p2,t} + \lambda l_{p1,p2,t} + 1 \ln \text{minutes_shared} + \\ + \sum_{j=1}^2 (\eta_j \text{value}_{p(j),t} + \theta_j \text{playerchar}_{p(j),t})) \end{aligned} \quad (7)$$

where, for both players, we include player valuation and total passes in each half-season, as well as $\text{position} \times \text{half_season}$, $\text{nationality} \times \text{half_season}$ and $\text{team} \times \text{half_season}$ dummies. Second, we estimate a version of the Poisson model with $\text{player}_1 \times \text{half_season}$ and $\text{player}_2 \times \text{half_season}$ fixed effects:

¹³For a discussion with regards to gravity models in international trade literature and the use of fixed effects Poisson Pseudo Maximum Likelihood estimation (FE-PPML) methods, see Fally (2015) and Santos-Silva and Tenreyro (2021). The procedure we use is described in (Berge, 2018).

¹⁴See (Hinz et al., 2021). It turns out that FE-PPML estimates can manage this type of bias better than non-linear OLS - one more reason for our preference of the model (Santos-Silva and Tenreyro, 2021). Weidner and Zylkin (2021) shows that the Poisson model still leave some room for potential bias, but with no double player fixed effects and a large number of observations, the bias should be small in our case. In a robustness check, we used the algorithm developed by (Weidner and Zylkin, 2021) and indeed found only a small bias.

$$\mu_{p1,p2,t} := E(\text{pass_count}_{p1,p2,t}|\dots) = \exp(\gamma g_{p1,p2,t} + \lambda l_{p1,p2,t} + 1 \ln \text{minutes_shared} + \gamma_{p1,t}^1 + \gamma_{p1,t}^2) \quad (8)$$

where $\gamma_{p1,t}^1$ and $\gamma_{p1,t}^2$ are $\text{player}_1 \times \text{half_season}$ and $\text{player}_2 \times \text{half_season}$ fixed effects ¹⁵.

To check robustness, we also run OLS regressions with $\log(\text{count})$ as dependent variable. With $\text{player}_1 \times \text{half_season}$ and $\text{player}_2 \times \text{half_season}$ fixed effects, we have:

$$E(\ln \text{pass_rate}_{p1,p2,t}|\dots) = \beta x_{p1,p2,t} + \delta z_{p1,p2,t} + \gamma_{p1,t}^1 + \gamma_{p1,t}^2 \quad (9)$$

In all estimated models, standard errors are clustered at the player 1 level. Standard errors clustered at player 1 * player 2 are somewhat smaller.

6. Homophily in collaboration

In this section, we present our core results along with extensions that help us better understand those results. We also discuss the robustness test based on OLS regressions with $\log(\text{count})$ as dependent variable and player fixed effects. Our dataset described in Section 3 is at the level of player-pairs and season-halves (N=355,610), and has 1568 team*season-half units in our dataset.

6.1. Core results

Let us start with the estimation results of the same-nationality indicator coefficient in our Poisson model. When the estimate is greater than zero, we will refer to it as the 'homophily premium'.

Column 1 in 3 looks at the unconditional difference: when we look at teams (team*half-season) in our data, we find evidence of an homophily premium: players of the same nationality tend to pass 6.8% more to each other than to player of different nationality. With an average of 31.7 passes between players in a season-half, this an unconditional mean corresponds to an average difference of 2.13 passes.

The estimated unconditional difference corresponds to overall homophily, and includes both induced and choice homophily. Our model allows us separating these two, and calculating choice homophily only.

Our model is estimated in the rest of the table with slightly different specifications. Columns 2 and 3 partial out several player and pair level characteristics. Column 4 replaces player characteristics with double player-period fixed effects. Column 5 allows the coefficient of minutes spent together to vary (it remains close to 1). All model

¹⁵We cannot have player1 - player 2 fixed effect as in some gravity models because our variable of interest is time-invariant (at least for the overwhelming majority of observations.)

Table 3: Baseline results

	pass_count				
	(1)	(2)	(3)	(4)	(5)
Shared nationality (0/1)	0.0681*** (0.0111)	0.0262*** (0.0064)	0.0240*** (0.0062)	0.0288*** (0.0067)	0.0282*** (0.0067)
Pass pair player values (ln)		0.0154*** (0.0043)	-0.0279*** (0.0043)		
Average length of passes (ln)		-1.004*** (0.0121)	-1.033*** (0.0125)	-1.159*** (0.0138)	-1.150*** (0.0137)
P1 total passes			0.1036*** (0.0038)		
P2 total passes			0.0696*** (0.0026)		
Total mins together (ln)					1.077*** (0.0048)
Observations	335,610	335,610	335,610	335,610	335,610
Pseudo R ²	0.05837	0.75416	0.75632	0.80188	0.80213
teamid-time fixed effects	✓	✓	✓	✓	✓
p1_position-league_time fixed effects		✓	✓		
p2_position-league_time fixed effects		✓	✓		
p1_citizenship-league_time fixed effects		✓	✓		
p2_citizenship-league_time fixed effects		✓	✓		
wh_player_id1-time fixed effects				✓	✓
wh_player_id2-time fixed effects				✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period. Match pass count is captured via team *time period fixed effects.

estimations in Columns 2 to 5 have team-time period (i.e. half-season) fixed effects absorbing the number of total passes in games from the theoretical model.

In the table, the homophily premium is fairly stable across specifications, with the same nationality indicator estimate ranging between 2.4 and 2.9%. Conditioning on observable player characteristics (such as team, position, valuation, citizenship), pass features (such as average distance) in Column 3, player-pairs of the same nationality tend to have a pass rate of 2.3% higher than player-pairs of different nationality. Note that all player characteristic variables vary over time, valuations change every half-season, and team, position, citizenship fixed effects are interacted with time period fixed effects. Columns 4 and 5 replace observable player characteristics with time varying player-period fixed effects. This allows player attributes to change over time. It also implies that the estimated same nationality coefficient is close to what would be the average of coefficients, estimated one by one for teams and time periods.

To interpret the coefficient of interest, we may think of a player who compares two potential passing partners who are identical (in terms of all features that are fixed in the given period) except the fact that only one of them has his same nationality. The player will pass 2.8% more per minute to his same nationality teammate.

Columns 2 to 4 are estimated for the pass rate, i.e. controlling for exposure to minutes by forcing the coefficient of the log of minutes shared to be 1. Column 5 relaxes this restriction by allowing for some non-linearity with virtually no effect on our point estimate.

To summarize, we estimated an overall homophily coefficient of 6.8% and choice homophily of 2.8%. This implies a positive induced homophily: players of the same nationality tend to cluster in a way that is conducive to more collaboration, too.

6.2. Extending costs

We modeled passing cost as the sum of two (log) additive parts. Variables that would be in our model (section 4 as part of the friction term $(l^{o,d})^\lambda$ capture non-distance-related frictions that would make it difficult for player o to pass to receiver d .

We have tried three such variables. First, a binary variable to capture when players' valuations are close to each other, signaling similar player quality. This is motivated by the possibility of positive assortative matching with higher quality players interacting more with each other. We indeed see that players of similar valuations tend to pass more to each other.

Second, there may be some physical attribute that is somehow typical of players of a certain country but not of others. We have data on player height and thus created a player-height-difference measure - it is indeed correlated with pass count but changes our coefficient only very marginally.

Third, beyond distance, we may we also measure the direction of play with an average forwardness index - capturing how forward a typical pass is (such as midfield

towards opposition goal). As Table 4 suggests, similar player valuation and the player height difference are both correlated with pass count, but do not affect our results.

6.3. *Endogeneity of time spent together*

A special feature of our setup is that teams change over time, players are replaced within games, and they are selected to play for some but not all games. In half-seasons, player-pairs spend on average 337 minutes or 20% of the maximum feasible amount of time together on the pitch.

We exploit this feature to estimate a complementary model, keeping all independent variables, but replacing the pass rate with minutes shared on the pitch as dependent variable. While the pass rate is eventually determined by players' decisions, minutes played together on the pitch are not: they are decided by the manager (coach).

Of course the manager's decisions are not random, and strongly correlate with the expected joint performance of the fielded players. This, in turn, will be based also on both the coach's evaluation of the same nationality premium and his observations of how actual players play together in training. The coach may even teach selected players to play together, thus improving their future passing activity. In this scenario, minutes shared is caused by same nationality via the same mechanisms we cared about.

In Tables 5 and 6 we repeat our core results for the pass rate along with estimating a similar model but now also for minutes-shared and the pass count - without conditioning for minutes-shared.¹⁶ The first table includes player characteristics, and the second table includes the player-level fixed effects.

We find that minutes shared on the pitch are also determined by same nationality, and so the coefficient estimate in the model without minutes-shared is higher. While the baseline model with the pass rate as the dependent variable is the one that corresponds to the theory and the object we care about, the modified model with minutes-shared as the dependent variable may offer a better estimate of what a coach can expect when considering player-pairs. Indeed if we consider minutes shared together to be endogenous, and view it as a homophily channel, the true parameter estimate would be closer to 3.7%.

6.4. *Benchmarking the premium*

Table 5 also allows for benchmarking the estimated coefficient by comparing it to player valuations.¹⁷

When we estimate our baseline model (Column 1), we find that the coefficient of the same nationality indicator is 2.6% while that of the log average player value is 1.5%. This suggests similar magnitudes: the difference in pass frequency between same and

¹⁶We exclude total player passes as they make it hard to use player valuation (Alternative table is .10 and .11 in Appendix.)

¹⁷As total passes and player valuations are very correlated, we excluded them here.

Table 4: Extended results

	pass_count			
	(1)	(2)	(3)	(4)
Shared nationality (0/1)	0.0240*** (0.0062)	0.0223*** (0.0061)	0.0288*** (0.0067)	0.0261*** (0.0066)
Pass pair player valuation (ln)	-0.0279*** (0.0043)	-0.0266*** (0.0043)		
P1 total passes	0.1036*** (0.0038)	0.1030*** (0.0038)		
P2 total passes	0.0696*** (0.0026)	0.0694*** (0.0027)		
Average length of passes (ln)	-1.033*** (0.0125)	-1.007*** (0.0126)	-1.159*** (0.0138)	-1.133*** (0.0137)
P1 and P2 valued similarly (0/1)		-0.0039 (0.0040)		0.0195*** (0.0045)
height_difference		-0.0085*** (0.0005)		-0.0155*** (0.0006)
Average forwardness Ind (0-1)		-0.1258*** (0.0142)		-0.0146 (0.0146)
Observations	335,554	335,554	335,610	335,610
Pseudo R ²	0.75632	0.75766	0.80188	0.80410
teamid-time fixed effects	✓	✓	✓	✓
p1_position-league_time fixed effects	✓	✓		
p2_position-league_time fixed effects	✓	✓		
p1_citizenship-league_time fixed effects	✓	✓		
p2_citizenship-league_time fixed effects	✓	✓		
wh_player_id1-time fixed effects			✓	✓
wh_player_id2-time fixed effects			✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period. Match pass count is captured via team *time period fixed effects.

Table 5: Looking at the role of minutes - Panel A

	pass_count (1)	minutes_shared (2)	pass_count (3)
Shared nationality (0/1)	0.0262*** (0.0064)	0.0384*** (0.0069)	0.0612*** (0.0100)
Pass pair player valuation (ln)	0.0154*** (0.0043)	0.6165*** (0.0080)	0.6478*** (0.0100)
Average length of passes (ln)	-1.004*** (0.0121)		-0.6085*** (0.0118)
Observations	335,554	335,554	335,554
Pseudo R ²	0.75416	0.29533	0.32610
teamid-time fixed effects	✓	✓	✓
p1_position-league_time fixed effects	✓	✓	✓
p2_position-league_time fixed effects	✓	✓	✓
p1_citizenship-league_time fixed effects	✓	✓	✓
p2_citizenship-league_time fixed effects	✓	✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period.

Table 6: Looking at the role of minutes - Panel B

	pass_count (1)	minutes_shared (2)	pass_count (3)
Shared nationality (0/1)	0.0288*** (0.0067)	0.0089*** (0.0030)	0.0370*** (0.0076)
Average length of passes (ln)	-1.159*** (0.0138)		-1.261*** (0.0146)
Observations	335,610	335,610	335,610
Pseudo R ²	0.80188	0.88506	0.74159
wh_player_id1-time fixed effects	✓	✓	✓
wh_player_id2-time fixed effects	✓	✓	✓
teamid-time fixed effects	✓	✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period. Match pass count is captured via team *time period fixed effects.

different nationality players is similar (somewhat greater) than between players of twice the valuations (conditioning on other player characteristics).

This result suggests that homophily has a very large effect on collaboration. The effect may be large partially due to the fact that the estimate is relative to the team average for players from a certain country playing at a certain position, which are features that are correlated with the effect of valuation.

However, as Column 2 and 3 shows player valuations play the greatest role in the selection process (minutes shared on the pitch). When we have passes instead of pass per minute (Column 3), the same nationality coefficient is 6.12%, whereas the player value coefficient is 64.78%. This suggests a much smaller relative effect.

6.5. Robustness to Poisson

Table 7 compares the results from the fixed effects Poisson model with those from a model featuring log pass count as the dependent variable and estimated by OLS with fixed effects. The estimated same nationality coefficients are quite similar, also when excluding minutes shared on the pitch.

6.6. "Deep" collaboration

If same nationality is helpful for collaboration in general, one would expect it to be even more so when collaboration is more complex. In our setup, this implies that

Table 7: Comparing Poisson and OLS results

	pass_count (1) Poisson	ln_pass_per_min (2) OLS	pass_count (3) Poisson	Total pass count (ln) (4) OLS
Shared nationality (0/1)	0.0288*** (0.0067)	0.0202*** (0.0058)	0.0370*** (0.0076)	0.0372*** (0.0073)
Average length of passes (ln)	-1.159*** (0.0138)	-0.6950*** (0.0091)	-1.261*** (0.0146)	-0.7891*** (0.0112)
Observations	335,610	335,610	335,610	335,610
Pseudo R ²	0.80188	0.31618	0.74159	0.39079
wh_player_id1-time fixed effects	✓	✓	✓	✓
wh_player_id2-time fixed effects	✓	✓	✓	✓
teamid-time fixed effects	✓	✓	✓	✓

Poisson regression model (C1, C3) and OLS model (C2,C4). Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period. Match pass count is captured via team *time period fixed effects.

same nationality should matter more for more complex passing patterns. To investigate whether this is indeed the case, we re-estimate the baseline model focusing on complex passing patterns. To compute complex passing sequence, we first make a change to the dependent variable. Instead of adding up passes, we identify pass sequences and look at the number of passes in a sequence. A N-long pass sequence is simply a series of N consecutive passes between two players A and B. The simplest is a single pass: AB or BA.

We define complex pass sequences as a pass sequence that includes at least two passes (ABA, BAB, ABAB, BABA, ABABA, etc).¹⁸ On average, player pairs carry out 31.73 passes per season-half. A vast majority, 87%, are single passes, but 13% are complex pass sequences. These complex pass sequences have 3.54 passes on average. 90% of player pairs have made at least one complex pass in our sample.¹⁹

Before we estimate a model with complex pass sequences, we re-estimate our baseline model with using the count of pass sequences rather than passes, which allows for a neater comparison. The count of passes and pass sequences are 99% correlated. Then,

¹⁸We may identify pass sequences is possible because our raw event data has timestamps that allow us to capture them.

¹⁹A large share of players who spend significant time on the pitch but do not produce complex passes are goalkeepers. In modern football they do pass quite a lot, but rarely take part in a sequence of passes.

Table 8: Pass sequences and complicated pass sequences

Dep. var.	Count of pass sequences			
	all_count (1)	complex_count (2)	all_count (3)	complex_count (4)
Shared nationality (0/1)	0.0205*** (0.0058)	0.0450*** (0.0097)	0.0249*** (0.0062)	0.0531*** (0.0103)
Average length of passes (ln)	-0.8927*** (0.0116)	-2.079*** (0.0177)	-1.008*** (0.0129)	-2.412*** (0.0197)
Observations	335,554	334,811	335,610	323,735
Pseudo R ²	0.75035	0.56048	0.79217	0.62074
teamid-time fixed effects	✓	✓	✓	✓
p1_position-league_time fixed effects	✓	✓		
p2_position-league_time fixed effects	✓	✓		
p1_citizenship-league_time fixed effects	✓	✓		
p2_citizenship-league_time fixed effects	✓	✓		
wh_player_id1-time fixed effects			✓	✓
wh_player_id2-time fixed effects			✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. Models in Column 1 and 2 include total passes and player valuations. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same value is binary, 1 if similar valued players. Same nationality is equality of either nationality. Player values are start of the time period. Sequence count is the number of pass sequences, complex seq count is the number of at least 2 pass-long sequences. Match pass count is captured via team *time period fixed effects.

we count the complex sequences, i.e. those with at least two consecutive passes, and use the count of complex passes as dependent variable.

Table 8 presents the results. When we only count the number of passes in complex sequences, we find a stronger effect: the homophily premium is twice the size for complex passes in both models. This suggests that homophily is especially important for deeper collaboration tasks.

[MORE]

6.7. Nationality and language

Our data allowed us to define mother tongues of players, allowing for players to have a background with multiple languages.

Multiple citizenships also create a challenge in determining common language. To deal with it, we rely on Melitz and Toubal (2014) to ascertain whether or not two countries share one or more common official and widely spoken languages. For example,

the official and widely spoken languages in Morocco are Arabic and French. Accordingly, Morocco and Egypt share Arabic, Morocco and France share French, Egypt and France do not share any language.

We then assume that a player speaks the official and widely spoken languages of his country of citizenship at the beginning of his career. Hence, a player starting his career in Morocco shares a common language with players starting their careers in Egypt or France, whereas a player starting his career in Egypt (France) shares a common language with players starting their careers in Morocco but not with those starting their careers in France (Egypt).

Comparing nationality and language for the largest groups of players by nationality, Table 1 reveals an incomplete overlap that arises from former colonial links, or participation in the same legal nationality entity or empire: Spanish is spoken in Argentina and Uruguay, Portuguese is spoken in Brazil, French is spoken in Senegal and Cote d’Ivoire. Another reason for incomplete overlap is that some small countries share a language with neighbors: Croatia with Serbia; Switzerland with France, Germany and Italy, Belgium with France and the Netherlands. This suggests that speaking the same language is not only about verbal interaction, but also about common cultural traits.

To summarize, comparing nationality and language reveals an incomplete overlap between them allowing us to disentangle their effects. Based on the discussion here, we created three groups of player-pairs.

1. No same nationality and no same language (Example: Argentina and Brazil)
2. No same nationality but same language as mother tongue (example: Argentina and Spain)
3. Same nationality (example: two Argentina players)

In the *player – pairseason – half* dataset, 38.3% of observations are of the same nationality, 11.4% do not share a same nationality but has the same language as a mother tongue, while 50% do have the same mother tongue.

We introduce two binary variables, one for the case when two players have the same nationality, and one for the case when they share a language but not nationality. Accordingly, the shared nationality coefficient is not expected to be necessarily the same as in the previous tables as the control group now excludes players of different nationality who speak the same language.

Results reported in 9 show that shared language has about half the impact of shared nationality. This suggests that shared nationality has indeed multiple components, with language being only one of them.

7. Conclusions

We have investigated how homophily related to nationality and language affect collaboration in multinational teams. In doing so, we have exploited a newly assembled

Table 9: Dissecting same nationality: language

	pass_count (1)	minutes_shared (2)	pass_count (3)
Shared nationality (0/1)	0.0326*** (0.0071)	0.0107*** (0.0031)	0.0418*** (0.0080)
Shared only language (0/1)	0.0167* (0.0087)	0.0072** (0.0035)	0.0208** (0.0098)
Average length of passes (ln)	-1.159*** (0.0138)		-1.261*** (0.0146)
Observations	335,610	335,610	335,610
Pseudo R ²	0.80189	0.88506	0.74160
wh_player_id1-time fixed effects	✓	✓	✓
wh_player_id2-time fixed effects	✓	✓	✓
teamid-time fixed effects	✓	✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality, same language and different nationality: the base category is different language (and nationality). Match pass count is captured via team *time period fixed effects.

exhaustive dataset recording all passes by professional European football players in all teams competing in the top five men leagues over eight sporting seasons, together with full information on players’ and teams’ characteristics.

We have measured collaboration as the average number of passes per minute between a pair of players in a half-season.

To separately identify excess collaboration within nationality or language due to preferences (‘choice homophily’) from collaboration due to opportunities (‘induced homophily’), we have used a discrete choice model of players’ passing behavior as a baseline. Estimating the model with a rich set of player and play characteristics as well as player fixed effects, we have found strong evidence of homophily. Conditioning on players’ and teams’ characteristics, player pairs of same nationality exhibit an average number of passes per minute that is 2.8 percent higher than player pairs of different nationality. Same nationality is about as likely to lead to more passes as doubling the player pair’s valuation, which is a consensus measure of players’ skills. Shared language has about half the impact of same nationality. Pairs of same nationality are also more likely to engage in deeper collaboration, disproportionately participating in more complex pass sequences.

These findings show that homophily based on nationality and language is pervasive even in teams of very high skill individuals with clear common objectives and aligned incentives, and involved in interactive tasks that are well-defined and not particularly language-intensive.

Appendix

A. Football rules

A.1. Key football rules

This subsection describes the key rules in football (soccer)

Each game has a starting XI - 11 players who start. There are up to 3 substitutions per team/game. (Typically in the last third of the game.). There will be 11 players on the pitch unless some player gets a red card and is sent out (permanently) - this rarely happens.

There is freedom in composition, but mostly: 1 goalkeeper, 3-5 defenders (left, center or full- and right-backs), 0-3 forwards, and midfielders. In our data we have very specific positions such as left-back.

Football teams have squads of about 25-30 players. Spain: 22-28, England: 25-33. Churning is large: 20-40% of team changes season to season.

Transfers happen twice a year in Europe. Summer transfer window is the main opportunity to get new players, or sell existing one. It is between 1 July to 1 September. This is the main window (90% of deals in season). The winter window is shorter, between 1 Jan to 1 Feb. Much smaller. (Typically 1-2 players per team)- Includes loan deals.

Note that there are games during the window, bit of complication = measurement error

[MORE]

A.2. Nationality rules in leagues

In some leagues, there is no limit on use of players of any nationality on the pitch, while in others non-EU, especially South American players face some restrictions. Also, some leagues have rules regarding the squad - must have home grown (academy) players - this has very little effect on starting XI. For our five leagues, we have two types of regulations.

Spain, France, Italy do have restrictions on foreign players. Foreign is defined as non-EU. In Spain it is max 3, in France it is max 4 and in Italy, it is max 2 non-EU. For these countries, the non-EU definition varies marginally but include players from 70 countries “Cotonou” agreement + countries offered the same by home country²⁰. In Spain, South Americans get citizenship after 5 years, 2 if they can show Spanish ancestry. In Italy, ancestry also allows a fast track to citizenship, which has helped many people from Argentina and Uruguay.

²⁰<https://www.footballmanagerblog.org/2018/04/football-manager-squad-registration-rules.html> and on the list of countries, see https://en.wikipedia.org/wiki/Cotonou_Agreement

Non-EU restrictions bite for some African/Asian players but mostly South Americans. The result is that in Spain, France and Italy, this regulation will imply that two Brazilians or a Uruguay and Argentina players are less likely to play together than two Europeans.

England and Germany do not have non-EU player restriction. But both have preference for home grown / academy product players, especially Germany. In England, visa restrictions are managed in a way that gives a preference to players who play or have the potential to play for their national team.

[MORE] bigskip

B. Additional information on data and cleaning

B.1. Players

If more than one citizenship, typically the first one is the most important one, matching playing for a national team. We kept nationalities as defined by the FIFA, ie a nationality is what has a national team. Hence, we have English and Welsh players not Team GB or Team UK. Our results are robust to having a single UK team. Similarly Feroer Island and some other geographies are also treated as separate nation. For country of birth, we sometimes made edits to match current list of countries. Most importantly, for multi-ethnic countries dissolved since (such as Yugoslavia and Soviet Union) born players were given their current nationality if it was part of federation, if not we imputed the largest country (like Russia).

National team may include U21 and U19 as well. In a very few cases, players switch allegiances, but that means they played no more than 3 games for the team they left. We only included the last one if more than one.

There are several small measurement error issues, all are negligible in impact. First, nationalities may be handed out mid-career. These are typically based on heritage (Italian speaker Argentinians getting Italian citizenship), and it should be a formality rather than a culture change. Second, regarding player values, small measurement errors may come from some low key players having a single year estimate - we replicated it for both half-seasons. Junior team member players promoted to senior team mid-season and would thus have no player value, and we imputed a minimum value here.

B.2. Matching players from two sources and entity resolution

Football players came from two different sources: from the pass event data ('events') and from the player information ('passport') we developed an entity coreference algorithm to match players – find out which records in the two different datasets describe the same player - are coreferent²¹.

²¹This section is based on the algorithm developed for this project by Endre Borza, see <https://github.com/endreborza/encoref>

A baseline solution would be fuzzy matching: use simply the names of the players, with any additional information available, like date of birth, height, nationality and match them based on similarity.

There are several complications for a standard fuzzy matching algorithm for this purpose. First, even for ten thousand players in our sample, it takes a lot of computing power to calculate all possible similarities and find the best ones. Second, simply matching the players by themselves is not precise enough, mostly due to the noise in the data: player features are also not precise and unique²². Third, data quality problems mean that some players might have two or more different records in one dataset. Fourth, algorithmic checkups are really important, as re-examining and correcting the possible matches for over ten thousand players is simply not feasible.

Our improved solution relied on introducing motifs. The core concept of the improved solution is not to match simply players, but match motifs in a network of players, matches, seasons and teams. This way, already discovered coreferences can be utilized to narrow the search space, and noise in the data can be mitigated by relying on more than one similarity to establish a coreference²³.

The algorithmic matching is not perfect - as players may use different names, especially South American players, and accents may be incorrectly used as well. When the matching score was low, we checked players by hand - about 1% of total names.

B.3. Detailed cleaning steps and decisions

We dropped own passes. We dropped observations (N=4000), when a player had a single partner.

[MORE]

C. Additional Tables and Results

C.1. Additional descriptive tables

[TO ADD]

C.2. Results: robustness tests

This the same table as 5 and 6 , but with the *all_pass* variables.

²²This problem can be demonstrated with an example of teams: in one dataset the names of two clubs are *Athletic* and *Atletico Madrid*, while in the other *Athletic Bilbao* and simply *Atletico*. So the solution must be open to the possibility, that the two entities, *Athletic* and *Atletico* even though are very similar in name.

²³See more on the algorithm at https://github.com/sscu-budapest/football-data-project/blob/main/reports/coreference_description.md

Table .10: Looking at the role of minutes (A)

	pass_count (1)	minutes_shared (2)	pass_count (3)
Shared nationality (0/1)	0.0240*** (0.0062)	0.0384*** (0.0069)	0.0243*** (0.0060)
Pass pair player valuation (ln)	-0.0279*** (0.0043)	0.6165*** (0.0080)	0.0017 (0.0037)
P1 total passes	0.1036*** (0.0038)		1.007*** (0.0023)
P2 total passes	0.0696*** (0.0026)		1.001*** (0.0028)
Average length of passes (ln)	-1.033*** (0.0125)		-1.076*** (0.0126)
Observations	335,554	335,554	335,554
Pseudo R ²	0.75632	0.29533	0.71477
teamid-time fixed effects	✓	✓	✓
p1_position-league_time fixed effects	✓	✓	✓
p2_position-league_time fixed effects	✓	✓	✓
p1_citizenship-league_time fixed effects	✓	✓	✓
p2_citizenship-league_time fixed effects	✓	✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period. Match pass count is captured via team *time period fixed effects.

Table .11: Looking at the role of minutes - (B)

	pass_count (1)	minutes_shared (2)	pass_count (3)
Shared nationality (0/1)	0.0288*** (0.0067)	0.0089*** (0.0030)	0.0370*** (0.0076)
Average length of passes (ln)	-1.159*** (0.0138)		-1.261*** (0.0146)
Observations	335,610	335,610	335,610
Pseudo R ²	0.80188	0.88506	0.74159
wh_player_id1-time fixed effects	✓	✓	✓
wh_player_id2-time fixed effects	✓	✓	✓
teamid-time fixed effects	✓	✓	✓

Poisson regression model. Standard errors, clustered at player 1 level, are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. German, French, English, Italian, Spanish top soccer leagues, 8 seasons: 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January. Same nationality is equality of either nationality. Player values are start of the time period. Match pass count is captured via team *time period fixed effects.

[MORE]

D. Team level results

This appendix shows the correlation between passing intensity and team performance. The data has passes aggregated to the level of teams and time periods (season-halves). Thus, one observation is a team*period (i.e. Inter Milan in the first half of the 2015-16 season). We have N=1568 observations (16 time periods, i.e. 8 seasons and 2 season-halves per season; 4x20 + 1x18 teams). Team performance is measured as the average number of points won in the time period. Teams get 0 for a loss, 1 for a draw, 3 for a win.

We look at how team performance measured by average points is correlated with \ln_pass defined as $\log(\text{average pass count per game})$. We estimate a simple model of correlations:

$$\text{Average_points} = \beta * \ln_pass + FEs \quad (.1)$$

First, we only include league dummies, and show a cross-section correlation for a single half-season (2015-16, H1). Then, we estimate a panel fixed effects model adding league-season-half and team fixed effects. Column 1 and 2 has points per game, Column

Table .12: Team level performance and passes

	points_per_game (1)	points_per_game (2)	ln_points_per_game (3)
Total pass count (ln)	1.142*** (0.2039)	0.2471** (0.1168)	0.2095*** (0.0800)
Standard-Errors	HC Robust	cluster: teamid	cluster: teamid
Observations	98	1,568	1,568
Pseudo R ²	0.32060	0.72565	0.95605
leagueseason_half fixed effects	✓	✓	✓
team fixed effects		✓	✓

OLS regression models. Standard errors are in parentheses. *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$. Team-period level data. German, French, English, Italian, Spanish top soccer leagues. Column 1: First half of 2015/16 season. Columns 2 and 3: 8 seasons 2011-2019. Time period is defined as a half-season, 16-20 games before and after 1 January.

3 has $\log(\text{points per game})$ as dependent variable for easier interpretation. Table .12 presents the results.

The first column reports the cross-section OLS results showing a very strong cross-sectional correlation between points per game and pass frequency. In the panel fixed effect models of column 2 and 3, we see a smaller but economically significant relationship.

We find evidence that, when teams pass more, they also tend to win more. Conditioning on league-season-half specifics, in periods when teams pass 10% more than the average pass frequency, they tend to win a 2.5% of a point (or 2.1%) more points than average. This difference is equivalent to one position difference in a typical league's standings.

References

- AlShebli, Bedoor K., T. R. and Woon, W. L. (2018), ‘The preeminence of ethnic diversity in scientific collaboration’, *Nature communications* **9**, 1–10.
- Berge, L. (2018), Efficient estimation of maximum likelihood models with multiple fixed-effects: the r package fenmlm, Working Paper 13.
- Bertrand, M. and Duflo, E. (2017), Field experiments on discriminationa, in A. V. Banerjee and E. Duflo, eds, ‘Handbook of Field Experiments’, Vol. 1 of *Handbook of Economic Field Experiments*, North-Holland, chapter 10, pp. 309–393.
- Buchholz, M. (2021), ‘Immigrant diversity, integration and worker productivity: uncovering the mechanisms behind ‘diversity spillover’ effects’, *Journal of Economic Geography* **21**(2), 261–285.
- Calder-Wang, S., Gompers, P. A. and Huang, K. (2021), Diversity and performance in entrepreneurial teams, Working Paper 28684, National Bureau of Economic Research.
- Currarini, S., Jackson, M. O. and Pin, P. (2009), ‘An economic model of friendship: Homophily, minorities, and segregation’, *Econometrica* **77**(4), 1003–1045.
- Currarini, S., Jackson, M. O. and Pin, P. (2010), ‘Identifying the roles of race-based choice and chance in high school friendship network formation’, *PNAS* **107**(11), 4857–4861.
- Earley, C. P. and Mosakowski, E. (2000), ‘Creating hybrid team cultures: An empirical test of transnational team functioning’, *Academy of Management Journal* **43**(1), 26–49.
- Ertug, G., Brennecke, J., Kovacs, B. and Zou, T. (2021), ‘What does homophily do? a review of the consequences of homophily’, *Academy of Management Annals* **0**(ja), null.
- Fally, T. (2015), ‘Structural gravity and fixed effects’, *Journal of International Economics* **97**(1), 76–85.
- Freeman, R. B. and Huang, W. (2015), ‘Collaborating with People Like Me: Ethnic Coauthorship within the United States’, *Journal of Labor Economics* **33**(S1), 289–318.
- Head, K. and Mayer, T. (2014), Gravity equations: Workhorse, toolkit, and cookbook, in G. Gopinath, E. Helpman and K. Rogoff, eds, ‘Handbook of international economics’, Elsevier, chapter 3, pp. 131–195.

- Hinz, J., Stammann, A. and Wanner, J. (2021), State Dependence and Unobserved Heterogeneity in the Extensive Margin of Trade, CEPA Discussion Papers 36, Center for Economic Policy Analysis.
URL: <https://ideas.repec.org/p/pot/cepadp/36.html>
- Hjort, J. (2014), ‘Ethnic divisions and production in firms’, *The Quarterly Journal of Economics* **129**(4), 1899–1946.
- Ingersoll, K., Malesky, E. J. and Saiegh, S. M. (2017), ‘Heterogeneity and team performance: Evaluating the effect of cultural diversity in the world’s top soccer league’, *Journal of Sports Analytics* **3**(2), 67–92.
- Jackson, S. E., Joshi, A. and Erhardt, N. L. (2003), ‘Recent research on team and organizational diversity: SWOT analysis and implications’, *Journal of Management* **29**(6), 801–830.
- Joshi, A., Labianca, G. and Caligiuri, P. M. (2002), ‘Getting along long distance: understanding conflict in a multinational team through network analysis’, *Journal of World Business* **37**(4), 277–284.
- Kahane, L., Longley, N. and Simmons, R. (2013), ‘The Effects of Coworker Heterogeneity on Firm-Level Output: Assessing the Impacts of Cultural and Language Diversity in the National Hockey League’, *The Review of Economics and Statistics* **95**(1), 302–314.
- Lang, K. (1986), ‘A language theory of discrimination’, *Quarterly Journal of Economics* **101**(2), 363–382.
- Laurentsyeva, N. (2019), From friends to foes: National identity and collaboration in diverse teams, Working Paper 226.
- Lawrence, B. S. and Shah, N. P. (2020), ‘Homophily: Measures and meaning’, *Academy of Management Annals* **14**(2), 513–597.
- Lazear, E. (1999a), ‘Language and culture’, *Journal of Political Economy* **107**(6), S95–S126.
- Lazear, E. P. (1999b), ‘Globalisation and the market for team-mates’, *The Economic Journal* **109**(454), 15–40.
- McPherson, J. M. and Smith-Lovin, L. (1987), ‘Homophily in voluntary organizations: Status distance and the composition of Face-to-Face groups’, *American Sociological Review* **52**(3), 370–379.
- McPherson, M., Smith-Lovin, L. and Cook, J. M. (2001), ‘Birds of a feather: Homophily in social networks’, *Annual Review Sociology* **27**(1), 415–444.

- Melitz, J. and Toubal, F. (2014), ‘Native language, spoken language, translation and trade’, *Journal of International Economics* **93**(2), 351–363.
- Neeley, T. (2015), ‘Global teams that work’, *Harvard Business Review* .
- Nüesch, S. and Haas, H. (2013), ‘Are multinational teams more successful?’, *International Journal of Human Resource Management* **23**(15), 3105–3115.
- Ottaviano, G. I. and Peri, G. (2005), ‘Cities and cultures’, *Journal of Urban Economics* **58**(2), 304–337.
- Ottaviano, G. I. and Peri, G. (2006), ‘The economic value of cultural diversity: evidence from US cities’, *Journal of Economic Geography* **6**(1), 9–44.
- Santos-Silva, J. and Tenreyro, S. (2021), The log of gravity at 15, Discussion Paper 1, School of Economics, University of Surrey.
- Terenzini, P. T., Cabrera, A. F., Colbeck, C. L., Bjorklund, S. A. and Parente, J. M. (2001), ‘Racial and ethnic diversity in the classroom’, *Journal Higher Education* **72**(5), 509–531.
- Tovar, J. (2020), ‘Performance, Diversity And National Identity Evidence From Association Football’, *Economic Inquiry* **58**(2), 897–916.
- Weidner, M. and Zylkin, T. (2021), ‘Bias and consistency in three-way gravity models’, *Journal of International Economics* **132**, 103513.