# Lost in Translation: The Language Barrier to Trade

Kiran Manthry*
*University of Cambridge*

This dissertation decomposes bilateral trade flows of goods and services into the effects of physical and linguistic distance. I employ a standard OLS estimation alongside a Poisson Pseudo Maximum Likelihood (PPML) Estimation to deal with specification issues and arrive at the novel result that while the linguistic effect on services is far stronger than that on goods, physical distance hinders trade in goods more than in services. Data is collected from a variety of sources including the OECD BATIS and BTDxE series.

## I. INTRODUCTION

In the majority of gravity models, a dummy variable for common official language represents shared language. Whilst this simplifies the model, it does a disservice to the diversity and diffusion of language across the globe. For example, Kenya and India share a common official language in English, however only 8.5% of Kenyans speak English with 18% of Indians speaking English (Eberhard, Simons and Fennig (2020)). The existence of the common official language is a relic of colonial influence in this case. To find the true effect of language on trade, we must disentangle language from colonisation and common culture.

The language effect is defined in this paper as the reduction in trade resulting from a smaller pool of speakers who can converse with each other. Examining the extreme cases clarifies the effect that language could possibly have on trade. Suppose two countries have no common language and no speakers can converse. It is simple to see that trade can only occur with great difficulty. Meanwhile, if a language barrier does not exist, trade can occur with little hinderance.

Defining a language, a question of linguistic flavour, has its own complications which I will not delve into. I will use the language definition provided by Ethnologue (Eberhard, Simons and Fennig (2020)).

As we move to a more digital economy, in the wake of the Covid-19 pandemic, more people are questioning whether physical distance affects their ability to work. If employees can work without great constraint from physical distance, trade in services could follow the same trend. Especially when one considers services, devoid of transport costs (save for IT infrastructure facilitating virtual trade), distance could capture omitted variables of cultural and historical similarity arising from physical proximity.

Some have claimed the 'Death of Distance' in gravity models (J. F. Brun et al. (2005)) due to the rise of globalisation and plummeting transport costs. However several papers (Leamer and Levinsohn (1995),J. F. Brun et al. (2005)) find consistent, increasing bilateral elasticities to distance.

This relation between trade and language could be very relevant to the UK in a post-Brexit world. Currently, only 32% of 15-30-year-olds can write in a second language whilst the EU average is 89% (Bowler (2020)). It has been labelled a 'crisis' by some and could hold the country back when forging new trading relationships. The UK cannot rely on the English lingua franca as other languages may rise to prominence. With Brexit making trade with 215 million English speakers in the EU (Eberhard, Simons and Fennig (2020)) much more difficult, the UK may need to improve language instruction if it wants to maintain its large surplus in services.

Linguistic diversity can also have large impacts on trade. The Chinese Communist Party's attempts to homogenise Mandarin Chinese and phase out minority languages (Craig and Prakash (2015)) across the mainland has obvious social and political effects. The economic effects, however, may take a longer time to be realised. By eliminating frictions arising from dialectal variation, inter-province trade should, in theory, increase. If this is the case, then it could be a model that countries with large linguistic diversity want to follow. However, regional groups often want to preserve their language. In India, attempts to promote Hindi as the national language have been strongly opposed by Southern states, especially Tamil Nadu (India Today Web Desk (2019)). Policies like these are not always feasible.

This paper contributes to the literature by examining the language effect with the new BATIS (Balanced Trade in Services) and BDTIxE (Balanced Trade in Goods by Industry and End-use) datasets. This allows me to analyse a wide range of countries. I improve the measure of common shared language and arrive at the novel result of shared language promoting services trade more than goods trade. The distinction between goods and services sheds some insight into how the language barrier manifests itself as a trade friction.

Section II examines literature by first giving an overview of the empirical developments around trade models. Then the choices of variables are critiqued with special consideration given to how linguistic distance/proximity is represented. Lastly, theoretical considerations are discussed.

In Section III I set out the mathematical framework for this paper. This will follow the (Anderson and Van Wincoop (2003)), AvW hereafter, gravity model deriva-

---

* kiran.gopinathan2000@gmail.com

tion.

Section IV explains why a log-log OLS estimation is inconsistent in the presence of heteroskedasticity. Other estimation options are explored, but the PPML model is chosen.

I describe the data in Section V. This includes discussion of the source and its limitations.

Results and analysis are provided in Section VI along with comparison between my results and the literature's. I find that linguistic distance significantly hinders trade in services whilst having little effect on goods trade. Physical distance is much more significant for goods than it is for services. Section VII evaluates how robust the results are by checking for heteroskedasticity. The Ramsey RESET is used to test for omitted variables. Finally, The Park test tests the specification of the PPML estimator.

In Section VIII I evaluate the results, conclude with policy recommendations and suggest areas for future research.

## II. LITERATURE REVIEW

### A. Empirical developments for trade models

The gravity model evolved from initial discussions by (Isard (1954), Tinbergen (1962), Pöyhönen (1963)). It augmented Newton's gravity equation to represent trade in a simple log-log transformation, gaining popularity through an empirical rather than theoretical justification. This remained the case until (Anderson (1979)) provided theoretical justification using a Common Elasticity of Substitution (CES) model.

The gravity model suffered from self-selection bias as many trade values are 0 and were therefore left out of the logarithm transformed equation (Westerlund and Wilhelmsson (2011)). (Helpman, Melitz and Rubinstein (2008)) introduce firm heterogeneity and deal with null values in one fell swoop by introducing a Heckman-style correction. However, (Santos Silva and Tenreyro (2006)), SST hereafter, demonstrate that a log-log form is inconsistent regardless of sample selection issues as it violates the conditional mean independence assumption in the presence of heteroskedasticity. They utilise PPML estimation to obtain consistent estimates of the gravity equation. This has now become the workhorse estimation method for bilateral gravity models.

Recent studies have tried to incorporate a network approach to account for second order dependencies (Baskaran et al. (2011), Ward, Ahlquist and Rozenas (2013)). This accounts for the opportunity cost in trade, also known as the border puzzle (AvW). Suppose you are examining trading patterns for Canadian provinces. The border between the US and Canada encourages trade between Canadian states more than it does with US states as US states have more domestic trading options. This is particularly pertinent when considering implications for policy. If we take a small country that has the option of teaching its children English or Swahili (and thereby joining their respective lingua francas), joining the English lingua franca, in theory, may increase overall trade, but joining the Swahili lingua franca will increase trade more with individual partners due to there being fewer Swahili-speaking nations.

### B. Variable selection

(Helliwell (1999)) addressed multilateral trade resistance using a remoteness variable which takes a GDP weighted sum of the distances of country i's trading partners.

Domestic linguistic diversity contributes to multilateral resistance. If home linguistic diversity is high, the probability of direct communication with someone domestically may be similar to the probability of direct communication abroad, thus reducing the perceived trade barrier. (Melitz (2008)) measures the effect of domestic linguistic diversity, using Greenberg's Diversity index (Greenberg (1956)). It calculates the probability two people in a country have the same mother tongue. This captures ethnic diversity more than linguistic diversity. It is also less of a problem in countries with strong language provision. For example, Switzerland would have a high rating on the Greenberg index (0.699) and a very high common language prevalence[1] (0.775). Whereas, Djibouti has a lower diversity score (0.568) but also a much lower common language prevalence (0.549).

The dependent variable is usually some flow measure of trade. This choice is non-trivial. (Helliwell (1999)) uses OECD trade in merchandise goods and achieves a relatively small language effect (common language increasing trade by 76%) given the lack of control variables. One would expect a larger language barrier on services in comparison to goods. Linguistically intensive services would expect an even larger language barrier. (Hutchinson (2005)) considers consumer and producer goods finding similar results for both.

Distance is also a contentious point in the discussion of gravity models. Many models use distance between capital cities (Hutchinson (2005)). Whilst this can be useful, capitals are sometimes sparsely populated (in relation to other cities) and away from the majority of the population (take Washington DC for example). A population weighted distance between the most populous cities (Mayer (2006)) would be a more accurate representation of physical distance.

### C. Linguistic distance

Linguistic distance can be represented in many ways. Its representation must be considered carefully as to ex-

---

[1] Defined in Section V

amine the pure language effect. A dummy variable representing common official language a la (Helliwell (1999)) omits key information. Moreover a dummy variable supposes that the marginal cost of translation is 0 (Melitz (2008)).

(Chiswick and Miller (2005)), CM hereafter, monitor how long native speakers of languages take to learn other languages. This could represent the incentive to learn specific languages (particularly pertinent in countries suffering from so-called 'Brain Drain'). However, it neglects the large proportion of multilinguals. CM also focus on English as the locus, while this study aims to be inclusive of more lingua francas. (Gould (1994)) notes that immigrants bring with them knowledge of their home markets and language. The dissipation of this knowledge is dependent upon the immigrant's access to fluency, therefore tempering the migration effect. (Hutchinson (2002)) uses this metric to analyse the interaction between bilateral migration and language learning with language competence augmenting human capital. This measure only looks at 36 countries and does not account for the prevalence of other languages in the country, only taking one of the official languages, thereby including historical similarities as well as the pure language effect. For example, South Korea and Japan have the same linguistic distance score despite 10% of South Koreans speaking English compared with less than 1% of Japanese residents (Eberhard, Simons and Fennig (2020)). The CM index may be a root cause for some of the language prevalences, but if data on language prevalence is available, it makes little sense to use data on language learning.

Some have opted for a more linguistic approach (Irvine (2006)). This looks at the pure language similarity in terms of the language morphosyntax. It suffers from the same issues as the CM index in that it is highly related to historic similarity. (Lohmann (2011)) using Irvine's index finds a relatively small elasticity of 0.68. (Melitz and Toubal (2014)) use a similar statistic based on the similarity between 200 important words. This suffers from the same issue of measuring historic similarity rather than acting as a current mechanism for trade frictions.

(Wagner, Head and Ries (2002), Hutchinson (2005), Melitz (2008)) uses the probability that an exporter and importer would speak the same language. This involves multiplying the respective language prevalences in the two countries. (Melitz (2008)) compares this (called Direct Communication in his work) and a binary indicator for whether a country shares a language spoken by at least 20% of the population (Open Circuit Communication). Melitz finds that Direct Communications (elasticity of influence 0.18-0.32) have a stronger Influence than Open Circuit Communication.

This measure is, unfortunately, prone to double counting, especially in areas of heavy linguistic diversity. To account for this, (Melitz and Toubal (2014)) use the CSL (Common Shared Language) variable which adjusts for this. There is little discrepancy when the similarity is made up from just one language but it penalises dyads with similarity made up from a variety of languages - which are thereby more prone to double counting.

### D. Theoretical considerations

The obvious theoretical explanation for language hindering trade is that you are more likely to work with a firm that speaks your language-literally and figuratively. However, this becomes less clear when translation firms offer these services. There are also very low marginal costs to translation (once a language has been acquired). In theory, one employee being able to speak the language should allow a firm to trade across the language barrier.

(Babcock and Du-babcock (1996)) qualitatively look at the trading behaviour of Taiwanese expatriates. All business-based communication was spilt into 3 zones varying by the language competency of both parties. They find that the language of communication varies depending by Zone and topic of conversation. A constant factor was the improvement in interpersonal relations when executives made the effort to learn the language. This suggests that an increase in potential channels of communication allows for better business relations and more trade at the micro level.

(Melitz and Toubal (2014)) try to disentangle language from trust arising from shared ethnicity. Language also represents colonial ties, in turn representing shared legal origins and values (Platteau (1994)). To account for this, they use a dummy variable for shared legal origin and stock of immigrants.

(Liwiński (2019)) looks at the wage premium gained from being multilingual, focusing on the Polish labour market. They find a significant heterogenous return to language ability. Spanish, French and Italian yielded premiums of 32%, 22% and 15% respectively while English and German yielded 11% and 12%. The large supply of English and German speakers may be a contributory factor. This premium on multilingual wages could manifest itself as a transaction cost, thus lowering trade. Furthermore, a larger percentage of speakers who can converse would reduce this premium thus increasing trade.

## III. THEORETICAL SPECIFICATION

This paper follows the AvW specification of the gravity model. Goods are differentiated by origin country and supply is fixed. It starts by assuming homothetic CES utility preferences.

$$\max u = \left\{ \sum_i \beta_i^{\frac{1-\sigma}{\sigma}} c_{ij}^{\left( \frac{\sigma-1}{\sigma} \right)} \right\}^{\frac{\sigma}{\sigma-1}} \tag{1}$$

Where $\beta_i$ represents a taste parameter across goods and $c_{ij}$ is consumption of goods in region $j$ by country i. The constant elasticity of substitution is represented by

$\sigma$. This utility is maximised subject to the budget constraint:

$$\sum_i p_i t_{ij} c_{ij} = y_j \tag{2}$$

Where $p_i$ is the exporter supply price and $t_{ij}$ is the trade cost factor. We assume that the trade costs are borne by the exporter in an iceberg form [2].

This results in a demand equation of the form:

$$x_{ij} = \left(\frac{\beta_i p_j t_{ij}}{P_j}\right)^{1-\sigma} y_j \tag{3}$$

$P_j$ is the consumer price index in country $j$.

$$P_j = \left[\sum_i (\beta_i p_i t_{ij})^{1-\sigma}\right]^{\frac{1}{1-\sigma}} \tag{4}$$

It can also be seen as multilateral trade resistance as it is dependent on all trade resistances. $x_{ij}$ represents nominal demand for country i goods in country $j$.

Market clearing implies

$$y_i = \sum_j x_{ij} = \left(\frac{\beta_i p_j t_{ij}}{P_j}\right)^{1-\sigma} y_j \tag{5}$$

By defining world income (and share of it) as

$$\theta_j = \frac{y_j}{\sum_j y_j} = \frac{y_j}{y_w} \tag{6}$$

and assuming symmetric trade frictions, $t_{ij} = t_{ji}$, it can be shown that

$$x_{ij} = \frac{y_i y_j}{y_w} \left(\frac{t_{ij}}{P_i P_j}\right)^{1-\sigma} \tag{7}$$

This is the basic gravity equation.[3]

There are multiple equilibria subject to different ratios of multilateral trade resistance. Since $t_{ij}$ and $P_j$ are not separately identified, a normalisation must take place.

The key result of this equilibrium is that trade barriers have heterogenous effects on trade depending on the size of the countries in question. An increase in bilateral trade barriers reduces size-adjusted trade between large countries more than it does between small countries.

Hypothesis 1: Linguistic distance hinders trade in services more than trade in goods

Hypothesis 2: Physical distance hinders trade in goods more than trade in services

---

[2] Iceberg costs assume that a constant fraction of the export 'melts away'. Suppose we send $x$ goods from country $i$ to country $j$, a fraction equal to $(t_{ij} - 1)/t_{ij}$ is lost in the process of transportation. This can be taken as a kind of ad valorem tax.

[3] See AvW for full derivation

## IV.   METHODOLOGY

If the gravity equation held with the consistency of a physical law, the log transformation would be of little consequence. Economic data, however, does not have this blessing and must account for disturbances. By taking the logarithmic form, the disturbance term is altered, and the conditional mean independence condition may no longer hold.

Suppose we take the stochastic gravity model,

$$x_{ij} = \frac{y_i y_j}{y_w} \left(\frac{t_{ij}}{P_i P_j}\right)^{1-\sigma} \eta_{ij} \tag{8}$$

Trade costs can then be decomposed into distance $d_{ij}$ and other trade costs $b_{ij}$

$$t_{ij} = d_{ij}^\rho b_{ij} \tag{9}$$

Taking logarithms then yields the log-log gravity equation with two-way fixed effects.

$$\begin{aligned}
\ln x_{ij} = & -\ln y_w + \ln y_i + \ln y_j + (1-\sigma)\rho \ln d_{ij} \\
& + (1-\sigma)\ln b_{ij} - (1-\sigma)\ln P_i \\
& - (1-\sigma)\ln P_j + \ln(\eta_{ij})
\end{aligned} \tag{10}$$

$$\begin{aligned}
\ln x_{ij} = & \alpha_0 + \alpha_i + \alpha_j + \ln y_i + \ln y_j + \beta_1 \ln d_{ij} \\
& + \beta_1 \ln b_{ij} + \ln(\eta_{ij})
\end{aligned} \tag{11}$$

The $\ln(\eta_{ij})$ term must be independent of all the regressors. However, the expected value for a logarithmic function is also dependent on higher order moments. In the presence of heteroskedasticity, the variance of $\eta_{ij}$ depends on $y_i$ and therefore $E[\ln(\eta_{ij})]$ also does. This violation of the Gauss-Markov assumption results in biased estimates of the coefficients. SST showcase this bias in Monte Carlo simulations. Since most trade data is heteroskedastic, another method of estimation must be used.

Another problem associated with the log-log OLS specification is the presence of null values. Since the logged value does not exist, these data points often drop out of the sample. This creates significant self-selection bias (Westerlund and Wilhelmsson (2011)).

If OLS cannot be used, we must turn to other methods.

The simplest method adds a small constant to all values of trade (Wang and Winters (1991), Baldwin et al. (2008)). Although this does solve the problem of null values, the coefficient estimates become biased (Gómez-Herrera (2013)).

Non-linear Least Squares (NLS) would be consistent, but it is far too inefficient for this data. The NLS coefficient is characterised by in this case:

$$\widehat{\beta} = \arg\min_b \sum_{1=1}^n \left[y_i - e^{x_i b}\right] \tag{12}$$

Yielding the first order condition:

$$\sum_{i=1}^{n} \left[ y_i - e^{x_i \widehat{\beta}} \right] e^{x_i \widehat{\beta}} x_i = 0 \qquad (13)$$

Due to the nature of the exponential function, undue weight is given to large trade flows and higher variance. Weighted NLS could be used if more about the distribution of errors was known, but the process of performing a 2-stage WLS can be cumbersome and inefficient with many dependent variables (SST). A Tobit regression has been used by (Anderson and Marcouiller (2002), Baldwin et al. (2008)). This, however, lacks theoretical rationale and requires the same variables to predict the probability of censure and the trade flow itself (Gómez-Herrera (2013)).

In a similar vein to Tobit, (Helpman, Melitz and Rubinstein (2008)) create a 2-step Heckman style correction model with the exclusion restriction based on firm heterogeneity. This provides a convincing story behind null values and an intuitive means of estimating them.

(Egger and Lassmann (2015)) use a discontinuous spatial regression to analyse trade along the language borders of the Swiss cantons. This lends itself to transaction level data and is therefore not suitable to modelling bilateral trade flows between countries. Global linguistic borders also lack the definition of Swiss linguistic borders.

Gamma and Poisson Pseudo Maximum Likelihood have similar properties and have gained recent popularity. The gamma PML assumes that $V\left[y_i \mid x\right] \propto E\left[y_i \mid x\right]^2$. This assumption may give more weight to the higher variance observations (thus more prone to measurement error) (Manning and Mullahy (2001)).

If, instead, we assume that $V\left[y_i \mid x\right] \propto E\left[y_i \mid x\right]$, equal weighting is given to each observation. Assuming this and correct specification of the conditional mean $E\left[y_i \mid x\right] = e^{x_i \beta}$ takes us to the current workhorse gravity estimation method-the PPML method which this study opts for. Whilst Poisson models are usually reserved for count data, infrequent and integer in nature, (Gourieroux, Monfort and Trognon (1984)) show that Poisson ML estimators do not have to be restricted to count data. The suitability of the PPML estimator can be tested through a variant of the Park test. This is provided in the robustness checks.

The PPML two-way fixed effects model (with origin and destination fixed effects) is prone to the incidental parameters problem (Lancaster (2002)). As Weidner and Zylkin show, when the 3-way PPML fixed effects model is specified (with origin destination and dyad fixed effects), it is not asymptotically unbiased (Weidner and Zylkin (2021)). Furthermore, the linguistic data and bilateral distance are unchanging over time and will therefore be wiped out by dyadic fixed effects. The country-specific fixed-effects would also wipe out variables of interest such as linguistic diversity. For this reason, fixed effects models are not considered.

TABLE I. construction of linguistic distance

| Language Prevalence (% of speakers) | | | |
| --- | --- | --- | --- |
| | Country_i | Country_j | Joint Prevalence |
| Language | Belgium | United Kingdom | |
| German | 0.22 | 0.06 | 0.0132 |
| English | 0.38 | 0.91 | 0.3458 |
| Spanish | 0.06 | 0.06 | 0.0036 |
| French | 0.82 | 0.19 | 0.1558 |

## V. DATA DESCRIPTION AND CONTROL VARIABLES

Imports and exports for services come from the latest edition of the WTO-OECD BATIS dataset (Liberatore and Steen (2021)). This is a dataset containing directed bilateral trade flows for 202 economies (both OECD and non-OECD nations) between 2005 and 2019. All values are represented in millions of US dollars. It aims to create a full matrix in reported values by estimating missing values. Missing values were calculated using a gravity framework with PPML estimation. 56% of the values are from reported flows and 44% are estimated. 21% are gravity estimates. We could work solely with reported values, to ensure that the regression does not become a kind of 2SLS. However, this presents selection bias problems if something links the countries with unreported values. Therefore, we will use all values including the estimated ones on the assumption that the final balanced values are accurate.

Data for manufactured goods comes from the 4[th] edition of the BTDIxE dataset (OECD (2021)). The time frame varies by reporter country; however, most countries have figures from 2005 onwards. Like the BATIS dataset, it also breaks down the trade flows by end use.

There are certain asymmetries in both datasets (exports from Country A to B not matching imports from B to A). These usually occur as a result of reimports/reexports. This means that both figures for imports and exports can be examined separately.

Common Spoken Language is constructed using the 23[rd] editions of Ethnologue's Language-in-Country (LIC) dataset (Eberhard, Simons and Fennig (2020)). Compiling the work of countless ethnologists, this dataset displays the number of speakers of a given language in a given country there are 11,373 records on 7464 distinct languages. For the purposes of this dissertation, separate dialects of macrolanguages were considered as part of the same language. Melitz and Toubal construct the variable using the following formula:

$$\text{CSL} = \max(\alpha) + \{\alpha - \max(\alpha)\}\{1 - \max(\alpha)\} \qquad (14)$$

Where alpha is the product of language prevalences in a given country. We can construct an example with real data:

In this example, alpha is the sum of these joint prevalences (0.518) whilst max(alpha) is 0.3458 from the contribution of English. CSL would then be 0.458.

This method allows values above 1 due to double counting. The adjustment tries to correct for this. The adjustment penalises alphas made from lots of sources, thus reducing double counting. However, it only takes the maximum, so it does not differentiate between a dyad sharing 2 languages and one sharing 5 languages.

A variation on this measure could be the root of the sum of the squares of prevalences.

$$\text{CSL}x_{ij} = \sqrt{\sum_n \left(p_{in}p_{jn}\right)^2} \tag{15}$$

Where $n$ denotes a shared language in the dyad $ij$.

This would favour dyads sharing just one language whilst penalising those sharing more. This is provided and denoted as CSLx_ij. In this example, it would equal 0.379 To do better than this, individual data on multi-lingualism would be needed.

Prevalences of less than 1% were dropped from the dataset for ease of computation.

Linguistic Diversity (or lack of) is represented by OwnCSL_i and OwnCSL_j. It is the common shared language but computed for $i = j$. This would just be the root of the sum of squares of domestic language prevalence. This measure is preferred to the Greenberg Index due its focus on shared language as opposed to native language which carries ethnocultural ties. Own_CSL is therefore negatively correlated with Greenberg's definition of diversity.

Migration_ij comes from the world bank migration database (World Bank Group 2011). This provides decennial bilateral migration flows from 1960 to 2000. These values were summated over each dyad to provide a stock of migrants that have immigrated after 1960 . This may be prone to counting migrants who have since moved on from the country in question. Heterogenous migration patterns are also not accounted for. The longer migrants stay in their home country, the more they integrate and the less they identify with their home country. (Hutchinson (2005)) finds that a stock of immigrants has a positive effect, but length of stay has a negative effect. Migrants will also have stayed in their previous country for varying amounts of time, thus forming different links with this country. This could be addressed in future studies when more accurate migration data becomes available.

Common religion is formed in a similar vein. It takes the sum of the products of religious prevalences between two countries for Catholics, Protestants and Muslims. Due to the nature of religion, the prospect of double counting is unlikely. Therefore, there is no need for a downward adjustment. This measure is taken directly from (Disdier and Mayer (2007)). These religions this make up 55% of the world's population and most of the state religions (Central Intelligence Agency (2021)). Despite this, a better measure would take different denominations into account. Homogenising Islam without delineating the Shia and Sunni sects is likely to provide a biased result. Hinduism and Buddhism also play a large role in many South East Asian countries, so its omission is noteworthy.

Common coloniser also comes from this dataset. It is set to 1 if two countries shared the same coloniser after 1945. Whilst this does omit many notable countries, the Unites States for example, it is broadly representative and avoids the problem of multiple colonisers and piecewise colonisation.

ColoniserColony_ij is a dummy for whether the two countries were ever in a coloniser-colony relationship. Both these variables (common coloniser and Coloniser-Colony) have similar limitations to Migration_ij. Colonial relationships are heterogenous and involve varying lengths of stay and influence in the colony. Whilst length of stay is a large determinant of the institutional similarity imposed by imperialism, the nature of the independence may explain more of the trade flows. For example, the Algerian War for independence may have harmed Algeria's trading relationship with France more than Australia's gradual decoupling from the British Empire.

Therefore, conflict needs to be accounted for and this is done through data from the Correlates of War Dyadic Inter-State War Dataset (Maoz et al (2019)). The variable YearsAtWar_ij represents how many years two countries have been at war between 1816 and 2010. A war is defined by both countries in the dyad either deploying at least 1000 troops or suffering 100 battle related deaths. In the case of multi-state conflicts, each state involved will have a recorded conflict with the states opposed to it in the conflict. This variable also suffers from large heterogeneity. Some wars are bloodier than others and will thereby sever trade links to varying degrees. Countries also perform different roles in the war (aggressor or retaliator). The issue of changing territories and regimes complicates the variable further. To maintain simplicity, this paper just takes conflicts involving current states. The conflicts involving states prior to this are likely to be, on average, much further in the past with less bearing on current politics. Civil wars are also unaccounted for in this dataset. This means that domestic fractions and ethnic alliances are omitted as a source of multilateral resistance.

Common legal origins can build trust and help enforce contracts (Platteau (1994)). This variable comes from (La Porta, Lopez-de-silanes and Shleifer (2008)). it is a dummy variable with 1 denoting a common legal origin (French, British, German and Socialist).

Remoteness, as mentioned before, is included in an attempt to measure bilateral trade frictions in relation to multilateral frictions. Helliwell's remoteness index is slightly modified to give less weight to small countries with lower GDPs. Helliwell's measure is calculated as so:

$$\text{REM}_i = \sum_j \frac{D\text{ist}_{ij}}{Y_j} \tag{16}$$

As $Y_i$ approaches 0, $\text{REM}_i$ explodes. Very distant countries also have a large effect. (Head and Mayer (2014))

TABLE II. Summary of RTA variable

| Type of RTA | Freq. | Percent | Dummy variable name |
|---|---|---|---|
| CU | 5,502 | 2.11 | Rtad1 |
| CU & EIA | 11,867 | 4.55 | Rtad2 |
| EIA | 223 | 0.09 | Rtad3 |
| FTA | 20,061 | 7.69 | Rtad4 |
| FTA & EIA | 26,747 | 10.25 | Rtad5 |
| PSA | 193,322 | 74.07 | Rtad6 |
| PSA & EIA | 3,140 | 1.20 | Rtad7 |
| No Data | 124 | 0.05 | Omitted |
| Total | 260,986 | 100.00 | |

suggest a more suitable measure:

$$\text{REM}_i = \left( \sum_j \frac{Y_j}{Dist_{ij}} \right)^{-1} \qquad (17)$$

This measure is not adversely affected by distant countries or small countries. It does, nonetheless, have limitations in its crudeness. The physical presence of the countries does not fully represent their trading relationship. It also says little about the network in which countries trade. More work needs to be done to account for multilateral trade resistance, but that is beyond the scope of this paper.

GDP is measured in current thousands of US dollars and is obtained by the World Development Indicators arm dataset produced by the World Bank (World Bank, 2021). This is time variant and covers 217 economies.

Traditionally, distance between capital cities has been used as in (Hutchinson (2005)). A slightly better measure would take weighted distance between most populous cities. Capital cities are often away from the epicentre of economic activity. The measure provided weights the distance by share of total population (Mayer (2006))

$$DistCES_{ij} = \left( \sum_{k \in i} \frac{pop_k}{pop_i} \sum_{l \in j} \frac{pop_l}{pop_j} d_{kl}^\theta \right) \qquad (18)$$

In this dataset, the $\theta$ parameter is set to $-1$.

Contiguity data comes from CEPIIs geodist dataset (Conte, Cotterlaz and Mayer (2020)) and is set to 1 if countries are contiguous.

Regional trade agreements can also take many forms with varying effects on trade. This dataset distinguishes between 8 types:

With CU being "Customs Union", EIA "Economic Integration Agreement" and FTA "Free Trade Agreement" (WTO (2021)). This provides a reasonable coverage of the different types of agreement; however, it does not control for how long the agreement has been in place. Newer trade deals have had less time to lay the groundwork for trade.

## VI. RESULTS

The headline result is that linguistic proximity is highly significant for services whilst it is barely significant for manufactured goods. Distance is still significant for services; however, the coefficients are significantly lower than the corresponding coefficients for manufactured goods. This suggests that the effects of distance have been conflated with linguistic distance. Contiguity promotes bilateral trade of goods whilst having no significant effect for services, thus further supporting our hypothesis.

The coefficients in the PPML model represent the log difference in the expected count. Since we are not working with count data, this is merely the semi-elasticity. For example, in the services export PPML, the coefficient of 1.540 represents the fact that an increase of 0.01 in CSLx_ij (in levels) increases bilateral trade by 1.56%. For the logged values, it represents an elasticity. The PPML coefficients are thereby comparable to OLS coefficients as they represent the same partial effects.

The partial effect of common language would vary depending on language and dyad prevalence. To put this in real terms, if the French speaking population of the United Kingdom increased from 19% to 20%, exports from the United Kingdom to Belgium would increase by 0.6%[4] whilst imports would increase by 0.5% on average.

The coefficients in this study generally seem larger than those in the literature. (Helliwell (1999)), when exploring merchandise trade, finds a semi-elasticity of 0.565 for their common language dummy with specific language elasticities ranging from 0.842 to (-0.110) for English and French, respectively. The coefficient here is larger with a smaller standard error.

(Melitz (2008)) finds coefficients between 0.83 and 1 for Direct Communication (constructed as mentioned before) after controlling for country-specific fixed effects. While (Melitz and Toubal (2014)) find an elasticity of 0.775.

Curiously, common religion has a negative coefficient. This could be due to the prevalence of religious wars. This measure of religious proximity is also rather crude with little differentiation of sects and lack of coverage. This self-selection of Abrahamic religions may bias the coefficient downwards. The trade between India and Sri Lanka, for instance, may be positively affected by a share of Hinduism, but it would not be treated as such in this metric. This goes against much of the literature (Linders et al. (2005)).

Linguistic diversity is only significant when considering service imports. This is at odds with (Melitz (2008))'s

---

[4] French speakers make up 82% of Belgium's population, so CSLx would increase from 0.379 to 0.383. This increase of 0.004 would increase the exponent by 0.004*1.540 thus increasing exports by. Thus, resulting in a 0.6% increase

TABLE III. PPML results[a]

| Variables | (1) Service Exports | (2) Manufacturing Trade Exports | (3) Service Imports | (4) Manufacturing Trade Imports |
|---|---|---|---|---|
| y2015 | 0.0740*** | -0.0263** | 0.0671*** | -0.0131 |
| | (0.00641) | (0.0109) | (0.00771) | (0.0131) |
| y2016 | 0.104*** | -0.0512*** | 0.0778*** | -0.0436*** |
| | (0.00803) | (0.0128) | (0.0125) | (0.0149) |
| y2017 | 0.0663*** | -0.0751*** | 0.0488*** | -0.0799*** |
| | (0.00927) | (0.0128) | (0.0142) | (0.0143) |
| y2018 | 0.0598*** | -0.0862*** | 0.0310** | -0.0908*** |
| | (0.00949) | (0.0141) | (0.0137) | (0.0155) |
| y2019 | 0.0909*** | -0.119*** | 0.0559** | -0.130*** |
| | (0.0109) | (0.0164) | (0.0235) | (0.0196) |
| OwnCSL_i | -0.00570 | -0.141 | -0.492*** | -0.151 |
| | (0.138) | (0.147) | (0.141) | (0.117) |
| OwnCSL_j | -0.0720 | -0.128 | 0.463*** | -0.323** |
| | (0.127) | (0.109) | (0.143) | (0.148) |
| Log(remoteness_i) | -0.0741 | 0.118** | -0.0836 | 0.205*** |
| | (0.0572) | (0.0514) | (0.0578) | (0.0604) |
| Log(remoteness_j) | -0.0451 | 0.200*** | -0.0631 | 0.165*** |
| | (0.0567) | (0.0534) | (0.0619) | (0.0568) |
| CSLx_ij | 1.540*** | -0.0310 | 1.333*** | 0.0426 |
| | (0.116) | (0.114) | (0.165) | (0.118) |
| lmMigration_ij | 0.0367*** | 0.0502*** | 0.0533*** | -0.0242** |
| | (0.00835) | (0.0113) | (0.0104) | (0.00979) |
| CommonReligion_ij | -0.385*** | -0.451*** | -0.349*** | -0.355*** |
| | (0.101) | (0.0949) | (0.111) | (0.0926) |
| CommonColoniser_ij | 0.578*** | 0.331** | 0.621*** | 0.470*** |
| | (0.138) | (0.158) | (0.162) | (0.149) |
| CommonLegalOrigins_ij | 0.177*** | 0.0928* | 0.159** | 0.108** |
| | (0.0527) | (0.0516) | (0.0623) | (0.0539) |
| ColoniserColony_ij | 0.157* | -0.132 | 0.0321 | -0.0317 |
| | (0.0904) | (0.111) | (0.0960) | (0.0990) |
| Log(GDP_i) | 0.745*** | 0.799*** | 0.721*** | 0.869*** |
| | (0.0170) | (0.0182) | (0.0173) | (0.0256) |
| Log(GDP_j) | 0.736*** | 0.759*** | 0.728*** | 0.876*** |
| | (0.0164) | (0.0184) | (0.0201) | (0.0265) |
| Log(DistCES_ij) | -0.496*** | -0.755*** | -0.459*** | -0.745*** |
| | (0.0389) | (0.0453) | (0.0460) | (0.0469) |
| YearsAtWar_ij | -0.00913*** | -0.00354 | -0.00542 | -0.00383 |
| | (0.00342) | (0.00385) | (0.00346) | (0.00380) |
| Contig_ij | -0.132 | 0.400*** | -0.147 | 0.351*** |
| | (0.0909) | (0.0724) | (0.0970) | (0.0810) |
| rtad1 | 0.595 | 0.596* | 0.122 | 1.070*** |
| | (0.428) | (0.321) | (0.343) | (0.395) |
| rtad2 | 0.948** | 0.895*** | 0.779** | 1.282*** |
| | (0.422) | (0.325) | (0.340) | (0.395) |
| rtad3 | 0.272 | -0.310 | 0.474 | 0.0190 |
| | (0.452) | (0.396) | (0.362) | (0.446) |
| rtad4 | 0.724* | 0.427 | 0.476 | 0.878** |
| | (0.425) | (0.326) | (0.337) | (0.396) |
| rtad5 | 0.775* | 0.952*** | 0.581* | 1.234*** |
| | (0.426) | (0.312) | (0.334) | (0.377) |
| rtad6 | 0.666 | 0.640** | 0.524 | 0.836** |
| | (0.417) | (0.311) | (0.320) | (0.373) |
| rtad7 | -0.179 | 0.365 | -0.505 | 0.690* |
| | (0.447) | (0.337) | (0.348) | (0.395) |
| Constant | -23.06*** | -5.604*** | -23.21*** | -7.997*** |
| | (1.610) | (1.680) | (1.624) | (1.741) |
| Observations | 110,014 | 110,014 | 124,117 | 124,117 |
| R-squared | 0.748 | 0.810 | 0.651 | 0.777 |

[a] *Robust standard errors in parentheses*
  *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

TABLE IV. OLS results[a]

| Variables | (1) Logged Service Exports | (2) Logged Manufacturing Trade Exports | (3) Logged Service Imports | (4) Logged Manufacturing Trade Imports |
|---|---|---|---|---|
| y2015 | 0.0569*** | 0.0362** | 0.0524*** | 0.106*** |
| | (0.00686) | (0.0147) | (0.00639) | (0.0140) |
| y2016 | 0.101*** | 0.0295* | 0.0746*** | 0.128*** |
| | (0.00726) | (0.0161) | (0.00716) | (0.0152) |
| y2017 | -0.0314*** | -0.155*** | -0.0819*** | -0.135*** |
| | (0.00789) | (0.0163) | (0.00817) | (0.0158) |
| y2018 | -0.0165** | -0.247*** | -0.0790*** | -0.237*** |
| | (0.00814) | (0.0172) | (0.00819) | (0.0168) |
| y2019 | -0.0221** | -0.140*** | -0.130*** | -0.208*** |
| | (0.00929) | (0.0183) | (0.00992) | (0.0184) |
| OwnCSL_i | 0.0554 | 0.309*** | -0.497*** | -0.618*** |
| | (0.0403) | (0.0640) | (0.0386) | (0.0570) |
| OwnCSL_j | -0.143*** | -0.296*** | 0.271*** | 0.554*** |
| | (0.0408) | (0.0557) | (0.0396) | (0.0581) |
| Log(remoteness_i) | -0.203*** | -0.0569* | -0.111*** | -0.109*** |
| | (0.0255) | (0.0306) | (0.0256) | (0.0362) |
| Log(remoteness_j) | -0.285*** | 0.126*** | -0.574*** | 0.0974*** |
| | (0.0247) | (0.0402) | (0.0225) | (0.0317) |
| CSLx_ij | 1.910*** | 0.877*** | 1.596*** | 0.653*** |
| | (0.0593) | (0.0779) | (0.0584) | (0.0802) |
| lmMigration_ij | 0.0384*** | 0.112*** | 0.0539*** | 0.115*** |
| | (0.00387) | (0.00554) | (0.00378) | (0.00573) |
| CommonReligion_ij | -0.468*** | -0.268*** | -0.393*** | -0.286*** |
| | (0.0442) | (0.0614) | (0.0418) | (0.0608) |
| CommonColoniser_ij | 0.506*** | 0.838*** | 0.646*** | 0.887*** |
| | (0.0413) | (0.0607) | (0.0373) | (0.0577) |
| CommonLegalOrigins_ij | 0.0239 | 0.0591* | -0.0608*** | -0.00738 |
| | (0.0214) | (0.0318) | (0.0212) | (0.0319) |
| ColoniserColony_ij | 0.702*** | 0.513*** | 0.706*** | 0.564*** |
| | (0.0705) | (0.107) | (0.0701) | (0.102) |
| Log(GDP_i) | 0.786*** | 1.266*** | 0.790*** | 1.005*** |
| | (0.00592) | (0.00839) | (0.00538) | (0.00806) |
| Log(GDP_j) | 0.729*** | 0.860*** | 0.726*** | 1.240*** |
| | (0.00591) | (0.00855) | (0.00550) | (0.00832) |
| Log(DistCES_ij) | -0.435*** | -1.089*** | -0.439*** | -0.922*** |
| | (0.0173) | (0.0250) | (0.0167) | (0.0250) |
| YearsAtWar_ij | 0.0161*** | 0.00744 | 0.0163*** | -0.00603 |
| | (0.00412) | (0.00818) | (0.00404) | (0.00882) |
| Contig_ij | 0.132** | 0.473*** | 0.0754 | 0.336*** |
| | (0.0653) | (0.0985) | (0.0670) | (0.107) |
| rtad1 | 0.422 | 1.325*** | 0.910** | 1.806*** |
| | (0.257) | (0.457) | (0.397) | (0.436) |
| rtad2 | 1.480*** | 1.756*** | 1.596*** | 2.445*** |
| | (0.252) | (0.448) | (0.393) | (0.425) |
| rtad3 | 0.649** | 0.856* | 1.314*** | 1.606*** |
| | (0.279) | (0.505) | (0.425) | (0.463) |
| rtad4 | 0.693*** | 1.012** | 1.046*** | 1.806*** |
| | (0.251) | (0.447) | (0.392) | (0.422) |
| rtad5 | 0.943*** | 1.223*** | 1.088*** | 1.841*** |
| | (0.250) | (0.445) | (0.391) | (0.421) |
| rtad6 | 0.490** | 0.602 | 0.777** | 0.984** |
| | (0.249) | (0.444) | (0.390) | (0.420) |
| rtad7 | -0.0288 | 1.260*** | 0.486 | 1.773*** |
| | (0.274) | (0.459) | (0.402) | (0.433) |
| Constant | -32.05*** | -21.05*** | -36.07*** | -26.53*** |
| | (0.787) | (1.189) | (0.823) | (1.162) |
| Observations | 104,809 | 110,011 | 118,286 | 124,117 |
| R-squared | 0.719 | 0.657 | 0.736 | 0.686 |

[a] *Robust standard errors in parentheses*
*** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

result of large pro-trade effects for home linguistic diversity. It does nonetheless have the correct signs (common shared language at home makes you less likely to import due to the existence of the language barrier abroad).

## VII.   ROBUSTNESS CHECKS

Efficiency should not be a problem given the size of the dataset. The most important condition to be met is asymptotic efficiency.

The data may not be suitable for PPML estimation. If the data is not distributed approximately following the Poisson distribution, then the coefficients will be asymptotically biased. We can test how well suited the data is between different maximum likelihood estimators using the Park test.

The relationship between the conditional variance and the conditional mean can be expressed in this way:

$$V\left[y_i \mid x\right] = \lambda_0 E\left[y_i \mid x\right]^{\lambda_1} \qquad (19)$$

$$\left(y_i - \hat{y}_i\right)^2 = \lambda_0 \left(\hat{y}_i\right)^{\lambda_1} + \xi_i \qquad (20)$$

It then follows that (after a Taylor expansion around $\lambda_1 = 1$ )

$$\left(y_i - \hat{y}_i\right)^2 = \lambda_0 \hat{y}_i + \lambda_0 \left(\lambda_1 - 1\right) \ln\left(y_i\right) \hat{y}_i \qquad (21)$$

For the PPML estimator to be adequate, $\lambda_1 = 1$. Therefore, the equation above can be estimated by PPML with the significance of the coefficient $\lambda_0\left(\lambda_1 - 1\right)$ representing a lack of fit for the Poisson distribution.

Table V shows that they all pass the test at the 5% level with only 1 failing at the 10% level. Therefore, the null hypothesis that $\lambda_1 = 1$ is not rejected at the 5% level, thus deeming PPML a suitable estimator.

To further prove the necessity of the PPML estimator, all of the OLS regressions fail the Breusch-Pagan test for heteroskedasticity at the 0.001% level. As mentioned previously, this creates biased estimates due to the logarithmic transformation.

To test for omitted variables, the Ramsey RESET test is used.

TABLE VI. Ramsey RESET $p$-values

| Measure of trade flow | PPML | OLS |
|---|---|---|
| Service exports | 0.0156 | 0.000 |
| Goods exports | 0.0248 | 0.000 |
| Service imports | 0.0232 | 0.001 |
| Goods imports | 0.0091 | 0.000 |

Given the large range of covariates and data explored, a confident value for the RESET test is hard to come by. Passing at the 1% level suggests a lack of omitted variables. The failings of the OLS test may be due to the specification problem covered earlier This suggests that the PPML model is, broadly, correctly specified.

The errors in the model could be temporally correlated. To correct for this, time fixed effects have been included.

To deal with the dyadic correlation, clustered standard errors are used.

## VIII.   CONCLUSION

The main results agree with conventional wisdom and the stated hypotheses to varying degrees. For services, linguistic distance is very important, but physical distance still relevant while linguistic distance has no impact on goods, with physical distance becoming more significant.

This result could nonetheless be contested. To fully answer the counterfactual, multilateral resistance needs to be addressed properly as it is not adequately covered by GDP weighted remoteness. Network economics is still relatively underutilised in the modelling of bilateral trade flows (Baskaran et al. (2011)). An empirical, dynamic, network-based approach could better capture the heterogeneity of multilateral resistance. Dynamic panel methods could be utilised to address serial correlation. Without this, it is difficult to make strong predictions and policy recommendations for the long run.

Sticking with the UK and Belgium example, the model predicts a 0.6% increase in exports and a 0.5% increase in imports. However, this does not account for the production function or make any kind of welfare statement.

The language barrier could also be seen as a barrier to development. Given that developed economies tend to specialise in services over goods, to make this transition, less developed countries will have to forge strong linguistic links through language education or migration. Linguistic diversity having little effect suggests that domestic linguistic unification policies may be slightly misguided on the economic front. These policies tend to alienate those who feel strongly about their language without providing much economic benefit in terms of trade. It could, nonetheless, improve social harmony, thus supplementing social capital, however more research would need to be done to explore that channel of influence.

The exact mechanism by which linguistic distance prohibits trade is still yet to be uncovered. Future studies differentiating between linguistically intensive and non-linguistically intensive exports could help address this. It could also be enhanced by data on language proficiencies. Certain services may require mastery of language whereas others may only require conversational language.

This paper, whilst emphasising the regional disparities in terms of language, has treated each individual language homogenously. (Helliwell (1999)) finds significant heterogeneity between languages. Moreover, acquisition of sibling languages is likely to be much easier than languages from a different family.

Language acquisition is seen as an exogenous endowment to countries in most models. Endogenizing the

TABLE V. OLS results[a]

| Variables | (1) Service Exports | (2) Manufacturing Trade Exports | (3) Service Imports | (4) Manufacturing Trade Imports |
|---|---|---|---|---|
| $\hat{y}_i$ | 1.892*** | 1.810*** | 1.623*** | 1.664*** |
| | (0) | (0) | (0) | (0) |
| $\ln(y_i)\hat{y}_i$ | 0.0109 | 0.0104 | 0.0383* | 0.0182 |
| | (0.231) | (0.238) | (0.0639) | (0.165) |
| Observations | 104,809 | 110,011 | 118,286 | 124,117 |
| R-squared | 0.580 | 0.615 | 0.314 | 0.528 |

[a] *Robust p-value in parentheses*
  *** $p < 0.01$, ** $p < 0.05$, * $p < 0.1$

learning of languages may provide more insight to governments seeking to implement foreign language policy.

The impact of script is yet to be explored. Ataturk's abandonment of the Turkish Ottoman script in favour of Latin script, may have been motivated by other social and political factors (Çolak (2004)), however, its impact on trade is yet to be explored. Erdogan's attempts to reverse this will therefore risk undoing these changes. Kazakhstan face a similar trade-off in their switch from the Cyrillic to a modified Latin script. Whilst being part of the former Soviet Union leaves them closer to Russia in an institutional sense, the language is closer to Turkish, making the phonemes more difficult to represent in a Cyrillic script (Kimanova (2011)).

ANDERSON, J. E. (1979): "A theoretical foundation for the gravity equation", *American Economic Review*, 69(1), 106-116.

ANDERSON, J. E. AND VAN WINCOOP, E. (2003): "Gravity with gravitas: A solution to the border puzzle", *American Economic Review*, 93(1), 170-192.

ANDERSON, J. AND MARCOUILLER, D. (2002): "Insecurity And The Pattern Of Trade: An Empirical Investigation", *The Review of Economics and Statistics*, 84(2), 342-352.

BABCOCK, R. D. AND DU-BABCOCK, B. (1996): "Communication Zones in International Business Communication", 372-412.

BALDWIN, R. ET AL. (2008): Study on the Impact of the Euro on Trade and Foreign Direct Investment, Directorate General Economic and Financial Affairs (DG ECFIN), European Commission.

BASKARAN, T. ET AL. (2011): "The Heckscher-Ohlin Model and the Network Structure of International Trade", *International Review of Economics and Finance*, 20, 135-145.

BOWLER, M. (2020): "A Languages Crisis?", Higher Education Policy Institute.

BRUN, J. F. ET AL. (2005): "Has distance died? Evidence from a panel gravity model", *World Bank Economic Review*, 19(1), 99-120.

C, P. ET AL. (2011): "Where on Earth is everybody? The evolution of Global Migration 1960-2000", *World Bank Economic Review*, 25(1), 12-56.

CENTRAL INTELLIGENCE AGENCY (2021): The World Factbook 2021. Available at: https://www.cia.gov/the-world-factbook/.

CHISWICK, B. R. AND MILLER, P. W. (2005): "Linguistic distance: A quantitative measure of the distance between English and other languages", *Journal of Multilingual and Multicultural Development*, 25(1), 26(1), 1-11.

ÇOLAK, Y. (2004): "Language policy and official ideology in early Republican Turkey", *Middle Eastern Studies*, 40(6), 67-91.

CONTE, M., COTTERLAZ, P. AND MAYER, T. (2020): "The CEPII Gravity Database".

CRAIG, J. M. AND PRAKASH, D. (2015): "Language Education in Xinjiang: A paradox for cultural homogenization and strengthened Ethnic identity".

DISDIER, A. C. AND MAYER, T. (2007): "Je t'aime, moi non plus: Bilateral opinions and international trade", *European Journal of Political Economy*, 23(4), 1140-1159.

EBERHARD, D. M., SIMONS, G. F. AND FENNIG, C. D. (2020): "Ethnologue Global Dataset Twenty-third edition data", 1-23.

EGGER, P. H. AND LASSMANN, A. (2015): "The Causal Impact of Common Native Language on International Trade: Evidence from a Spatial Regression Discontinuity Design", *Economic Journal*, 125(584), 699-745.

GÓMEZ-HERRERA, E. (2013): "Comparing alternative methods to estimate gravity models of bilateral trade", *Empirical Economics*, 44(3), 1087-1111.

GOULD, D. M. (1994): "Immigrant Links to the Home Country: Empirical Implications for U . S . Bilateral Trade Flows", *The Review of Economics and Statistics*, 76(2), 302-316.

GOURIEROUX, A. C., MONFORT, A. AND TROGNON, A. (1984): "Pseudo Maximum Likelihood Methods: Applications to Poisson Models Published by: The Econometric Society".

GREENBERG, J. H. (1956): "The Measurement of Linguistic Diversity", *Linguistic Society of America*, 32(1), 109-115.

HEAD, K. AND MAYER, T. (2014): "Chapter 3 - Gravity Equations: Workhorse,Toolkit, and Cookbook", in Gopinath, G., Helpman, E., and Rogoff, K. B. T.-H. of I. E. (eds) Handbook of International Economics. *Elsevier*, 131-195.

HELLIWELL, J. F. (1999): "Language and Trade", in New Canadian Perspectives: Exploring the Economics of Language.

HELPMAN, E., MELITZ, M. AND RUBINSTEIN, Y. (2008): "Es-

timating Trade Flows: Trading Partners and Trading Volumes*", *The Quarterly Journal of Economics*, 123(2), 441-487.

HUTCHINSON, W. K. (2002): "Does ease of communication increase trade? Commonality of language and bilateral trade", *Scottish Journal of Political Economy*, 49(5), 544-556.

HUTCHINSON, W. K. (2005): "'Linguistic Distance' as a Determinant of Bilateral Trade", *Southern Economic Journal*, 72(1), p. 1.

INDIA TODAY WEB DESK (2019): No intention of imposing Hindi, says Centre after Tamil Nadu leaders warn of protests, India Today. Available at: https://www.indiatoday.in/india/story/no-intention-of-imposing-hindi-new-education-policy-centre-tamil-nadu-protests-1540527-2019-06-02.

IRVINE, A. (2006): "Measuring Linguistic Similarity within Language Groups: Revisiting Tradutional Typology in the Language Information Age."

ISARD, W. (1954): "Location Theory and Trade Theory: Short-Run Analysis", *The Quarterly Journal of Economics*, 68(2), 305-320.

KIMANOVA, L. (2011): "Analysis of Arguments in the Public Debate on the Alphabet Change in bilingual Kazakhstan", *Gaziantep University Journal of Social Sciences*, 10(3), 1021-1035.

LANCASTER, T. (2002): "Orthogonal Parameters and Panel Data", 647-666.

LEAMER, E. E. AND LEVINSOHN, J. B. T.-H. OF I. E. (1995): "Chapter 26 International trade theory: The evidence", 1339-1394.

LIBERATORE, A. AND STEEN, W. (2021): "The OECD-WTO Balanced Trade in Services Database (Bpm6/EBOPS 2010)", *(January)*, 1-33. Available at: https://www.oecd.org/sdd/its/balanced-trade-in-services.htm.

LINDERS, G. M. ET AL. (2005): "Cultural and Institutional Determinants of Bilateral Trade Flows", SSRN Electronic Journal.

LIWIŃSKI, J. (2019): "The wage premium from foreign language skills", *Empirica*, 46(4), 691-711.

LOHMANN, J. (2011): "Do language barriers affect trade?", *Economics Letters*, 110(2), 159-162.

MANNING, W. G. AND MULLAHY, J. (2001): "Estimating log models: to transform or not to transform?", *Journal of Health Economics*, 20(4), 461-494.

MAYER, T. (2006): "Notes on CEPII's distances measures", *(May)*, 20(4), 1-5.

MELITZ, J. (2008): "Language and foreign trade", *European Economic Review*, 52(4), 667-699.

MELITZ, J. AND TOUBAL, F. (2014): "Native language, spoken language, translation and trade", *Journal of International Economics*, 93(2), 351-363.

OECD (2021): "BTDIxE Bilateral Trade in Goods by Industry and End-use, ISIC Rev.4", Available at: https://stats.oecd.org/Index.aspx?DataSetCode=BTDIXE_I4.

PLATTEAU, J. (1994): "Behind the market stage where real societies exist - part I: The role of public and private order institutions", *The Journal of Development Studies*, 30(3), 533-577.

LA PORTA, R., LOPEZ-DE-SILANES, F. AND SHLEIFER, A. (2008): "The Economic Consequences of Legal Origins", 285-332.

PÖYHÖNEN, P. (1963): "A Tentative Model for the Volume of Trade between Countries", *Weltwirtschaftliches Archiv*, 90, 93-100.

SANTOS SILVA, J. M. C. AND TENREYRO, S. (2006): "The log of gravity", *Review of Economics and Statistics*, 88(4), 641-658.

TINBERGEN, J. (1962): "The log of gravityShaping the world economy; suggestions for an international economic policy".

WAGNER, D., HEAD, K. AND RIES, J. (2002): "Immigration and the Trade of Provinces", 49(5), 507-525.

WANG, Z. K. AND WINTERS, L. (1991): The Trading Potential of Eastern Europe. C.E.P.R. Discussion Papers.

WARD, M. D., AHLQUIST, J. S. AND ROZENAS, A. (2013): "Gravity's Rainbow: A dynamic latent space model for the world trade network", *Network Science*, 1(1), 95-118.

WEIDNER, M. AND ZYLKIN, T. (2021): "Bias and Consistency in Three-way Gravity Models", Papers 1909.01327, arxiv.org, revised March 2021.

WESTERLUND, J. AND WILHELMSSON, F. (2011): "Estimating the gravity model without gravity using panel data", *Applied Economics*, 43(6), 641-649.

WTO (2021): Regional Trade Agreements Database, http://rtais.wto.org/UI/About.aspx

ZEEV MAOZ, PAUL L. JOHNSON, JASPER KAPLAN, FIONA OGUNKOYA, AND AARON SHREVE (2019): "The Dyadic Militarized Interstate Disputes (MIDs) Dataset Version 3.0: Logic, Characteristics, and Comparisons to Alternative Datasets", Journal of Conflict Resolution (forthcoming)