

# Where to live? English proficiency and residential location of UK migrants.\*

Yu Aoki,<sup>a†</sup> Lualhati Santiago<sup>b</sup>

<sup>a</sup>*Department of Economics, University of Aberdeen, Dunbar Street, Aberdeen, AB24 3QY, United Kingdom*

<sup>b</sup>*Division of Analysis, Microdata and Engagement, Directorate of Macroeconomic Statistics and Analysis, Office for National Statistics, 2 Marsham Street, London SW1P 4DF, United Kingdom*

## ABSTRACT

This article analyses the causal effect of English proficiency on residential outcomes of migrants using a natural experiment. Based on the phenomenon that young children learn a new language more easily than older children, we construct an instrument for English proficiency exploiting age at arrival in the United Kingdom for childhood migrants. Using a unique dataset, we construct various measures of the extent of residential segregation aimed at capturing different dimensions of enclaves, and find a negative impact of better English skills on residency in a language enclave, but a *positive* impact on ethnic-enclave residency. We also find strong evidence of an impact of poorer English proficiency on living in a neighbourhood of lower quality. *Keywords:* Language skills, residential segregation, enclave, neighbourhood quality. *JEL codes:* J15, R23, Z13.

## 1. Introduction

Globally, 272 million international migrants were recorded in 2019, where Oceania recorded the highest proportion of immigrant population, 21 per cent, followed by North America, 16 per cent, and Europe, 11 per cent (United Nations, 2019). We live in an increasingly diverse society and the social integration of migrants is becoming an important policy objective in many developed countries. Although it is widely believed that proficiency in the language spoken in the host country is an important factor for promoting integration, there is limited

---

\*Acknowledgments: We gratefully acknowledge the permission of the Office for National Statistics to use the Longitudinal Study, as is the help provided by staff of the Centre for Longitudinal Study Information & User Support (CeLSIUS). CeLSIUS is supported by the ESRC Census of Population Programme (Award Ref: ES/V003488/1). Financial support from the Scottish Institute for Research in Economics and Carnegie Trust for the Universities of Scotland is gratefully acknowledged. We also thank the participants of the conferences of the EALE/SOLE online, the ESPE Barcelona online, the Scottish Economic Society online, and of the seminars at the University of Warwick, University of Curtin Australia, and Institute for Labour Law and Industrial Relations in the European Union. The authors alone are responsible for the interpretation of the data.

This work contains statistical data from ONS which is Crown Copyright. The use of the ONS statistical data in this work does not imply the endorsement of the ONS in relation to the interpretation or analysis of the statistical data. This work uses research datasets which may not exactly reproduce National Statistics aggregates.

<sup>†</sup>Corresponding author.

*E-mail addresses:* y.aoki@abdn.ac.uk (Aoki), Lualhati.Santiago@ons.gov.uk (Santiago)

knowledge on its causal impact on residential location outcomes in the host country. We add to this knowledge by investigating the causal effects of language skills on a variety of residential outcomes that measure the extent of residential segregation and the quality of the neighbourhood migrants reside in.

Our paper makes mainly three contributions to the literature on immigrant outcomes in a host country. First, we construct measures of the extent of residential segregation of migrants aimed at capturing the concept of enclave along four dimensions: Main language spoken by residents (language enclave), ethnicity (ethnic enclave), country of birth (country-of-birth enclave), and world region of birth (region-of-birth enclave). We analyse which of these four dimensions of enclave is affected by the language skills of migrants, when they are making a residential location decision. Distinguishing different types of enclave is important, because language skills can have heterogeneous impact on different types of enclave. For example, migrants who are fluent in English may live outside of a language enclave, if the reason for living in an enclave is simply for linguistic convenience. However, migrants proficient in English may decide to live in an ethnic enclave if they value other aspects of living in an enclave, such as offering employment opportunities, cultural amenities or protection from possible discrimination they might face outside of the enclave. Which aspect of enclave migrants value, and thus how English proficiency affects the residency in different types of enclave is an empirical question.

Second, we study the quality of the neighbourhood migrants live in, where the quality is measured at a small geographical area of an average of 1,500 individuals. We can conduct this analysis by linking a unique dataset from the Office for National Statistics (ONS) Longitudinal Study, which contains individual-level data from the England and Wales 2011 Census, to the indices measuring neighbourhood quality in England. The various measures of neighbourhood quality we exploit capture different quality dimensions (i.e., the extent of income, employment, and health deprivation of residents), allowing us to analyse residential environments in which migrants with different levels of English proficiency live. It is important to analyse this, since lower language proficiency might have amplifying negative effects, if it not only directly affects migrants' social and labour market outcomes, as shown in the existing literature to be summarised in Section 2, but also indirectly affects their outcomes through neighbourhood effects in their residential areas. We are not aware of any other studies that have provided arguably causal evidence on the impact of language proficiency on the quality of the neighbourhood in which migrants reside.

Third, we also make a methodological contribution in identifying the causal effect of language on residential outcomes. A major challenge to identify the causal effect is the endogeneity of language skills. First, there may be reverse causality. For example, poorer English skills may lead an individual to live in an area with a higher concentration of individuals speaking

their native language, while at the same time living close to individuals who speak their native language may make it more difficult to improve their English skills. Second, there may be unobserved heterogeneity across individuals that is correlated with both English skills and residential outcomes (e.g., ability). Third, the self-reported measure of English skills used in our analysis may contain measurement error. To address these possible endogeneity concerns, the seminal papers Bleakley and Chin (2004; 2010), use an instrumental variable (IV) strategy where age at arrival in the host country is exploited to construct an instrument for English language skills. The idea of using age at arrival is based on the “critical period hypothesis of language acquisition” by Lenneberg (1967), suggesting that individuals exposed to a new language within the critical period of language acquisition (i.e., childhood) can learn it more easily than those exposed to it outside of this critical period. This hypothesis implies that migrants who arrived in the host country when they were young would on average have better English skills than migrants who arrived at an older age.

However, age at arrival is unlikely to be a valid instrument for English skills on its own because it may influence an immigrant’s residential outcomes through other channels than language acquisition; for example, through cultural assimilation. To overcome this problem, Bleakley and Chin (2004; 2010) interact age at arrival with a dummy variable for being born in a non-Anglophone country, based on the following idea. After arriving in the host country, migrants born in non-Anglophone countries would be exposed to a new language in addition to the new country environment, while those born in Anglophone countries would be exposed to the same new country environment but not to a new language. This implies that, conditional on individual characteristics, any difference observed in the outcomes of early- and late-arrivers born in Anglophone countries reflects age-at-arrival effects only, while this same difference observed in migrants born in non-Anglophone countries reflects those same age-at-arrival effects and an additional effect, the language effect. Thus, the difference in an outcome of early- and late-arrivers for those born in non-Anglophone countries in excess of the equivalent difference for those born in Anglophone countries can be attributed to the effect of language.

One might be concerned, however, that Anglophone and non-Anglophone migrants face different age-at-arrival effects, and thus Anglophone migrants may not be able to be used to partial out the age-at-arrival effects. This is an issue that a large number of papers using this type of identification strategy face (e.g., Bleakley and Chin, 2004, 2010; Akbulut-Yuksel et al., 2011; Miranda and Zhu, 2013; Guven and Islam, 2015; Yao and van Ours, 2015). To mitigate this type of concerns, instead of using the interaction of age at arrival with coming from a non-Anglophone country as an instrument for English skills, we use this interaction term as a control variable in our empirical model. We then instead construct an instrument exploiting

the variation in linguistic distance from English,<sup>1</sup> within non-Anglophone countries. Precisely, among non-Anglophone countries, there is variation in how close their native languages are to English. For example, Dutch is linguistically more similar to English than Vietnamese. We use this heterogeneity in similarity of migrants' native languages to English, and construct our instrument by interacting age at arrival with linguistic distance between the origin-country language and English. We are not aware of any other studies that exploit the interaction of age at arrival with coming from a non-Anglophone country as a control variable, rather than as an instrumental variable. We also account for year-of-arrival fixed effects and parental education, which are possibly important confounding factors that have not been accounted for in the causal literature on immigrant's residential outcomes.

Our IV estimates indicate that language skills have a heterogeneous impact on residency in different types of enclave: Poorer English skills significantly lead migrants to live in areas with a higher concentration of individuals who speak their native language (i.e., language enclave). In contrast, we find positive effects of better English proficiency on residency in an ethnic enclave and a region-of-birth enclave, suggesting that *better* English skills lead migrants to live in those types of enclave. Our results highlight the importance of distinguishing different dimensions of enclave and suggest different mechanisms of the impact of language skills at play. Turning to the quality of the neighbourhood where migrants live, we find strong evidence that poorer English skills lead migrants to live in a neighbourhood of lower quality. Our supplementary analysis finds that better educational attainment as a result of better English skills may be a key channel through which language affects neighbourhood quality outcomes, but a different mechanism is likely at play for enclave residency.

The rest of the paper is structured as follows: Section 2 highlights the differences of our work from the existing literature, after reviewing the literature on residential segregation and on neighbourhood quality. Section 3 presents our identification strategy, while Section 4 describes datasets and sample specifications, and explains how we construct our main variables, such as various measures of residential segregation. Section 5 presents our empirical results, and Section 6 conducts a series of robustness checks aimed at addressing various concerns that can threaten the validity of our identification strategy. Finally, Section 7 discusses policy implications and conclusions.

---

<sup>1</sup>Clarke and Ispording (2017) is the first to exploit this variation in their study of language skills and immigrants' health.

## 2. Literature review

Starting from studies on residential segregation, its relation with language proficiency has been extensively studied by researchers in economics and other disciplines. In a seminal paper, Lazear (1999) proposes a model of cultural and language assimilation of migrants that inversely relates an immigrant's language proficiency to the proportion of local population who speak their same native language. This model predicts that an immigrant residing in an area with a large proportion of people who speak their native language has less incentive to learn a new language. On the other hand, the model of spatial assimilation developed by Massey (1985) suggests that ethnic enclaves are a natural first stage for migrants when entering a country, but they leave the enclaves once they have integrated to the new country's culture.

Empirically, a large number of papers have investigated the correlation between host-country language proficiency and ethnic residential segregation (e.g., Logan et al. 2002; Dustmann and Fabbri 2003; Bauer et al. 2005; Iceland and Scopilliti 2008; Beckhusen et al. 2013). Broadly, they find that having lower English language skills is positively correlated with ethnic-enclave residence. For example, Dustmann and Fabbri (2003), in their analysis of the determinants of language skills, find strong negative correlations between ethnic minority concentrations and English language skills of ethnic minority migrants in the UK. Although it is not numerous, there is also some research on residence in language enclaves (e.g., Chiswick and Miller, 1995, 2005). For example, Chiswick and Miller (2005) study the relation between living in a language enclave and English proficiency of migrants in the United States (US), and find that English proficiency is negatively associated with a higher extent of minority language concentration. A limitation of these studies is that it is not clear which direction causality runs: Namely, whether poor language skills cause migrants to live in enclaves, or whether they have poor language skills because they live in an enclave. In fact, there are studies investigating the opposite relation, i.e., the effect of enclave residence on language skills (e.g., Cutler et al., 2008; Danzer and Yaman, 2016; Laliberte, 2019), indicating the relevance of reverse causality in this context. Bleakley and Chin (2010) is the first paper to address this potential endogeneity issue using an IV for English proficiency, which is an interaction between age at arrival in the US and an indicator for being born in a non-Anglophone country. They find weak evidence of the effects of English proficiency on ethnic and country-of-birth enclave residence, unlike previous studies that found strong correlations between host-country language proficiency and enclave residence.

We add to the existing literature by analysing the causal effects of language skills on four different measures of segregation based on language, ethnicity, country of birth, and world region of birth in a comparable manner. Distinguishing different measures of segregation is

important as language skills could have heterogeneous effects on different types of enclave. For example, if the reason for living in an enclave is simply for linguistic convenience, migrants who are fluent in English would have no incentive to live in a language enclave. However, migrants fluent in English may choose to live in an ethnic enclave if they value other aspects of living in an enclave, such as offering employment networks, cultural identity and protection from possible discrimination. Which aspect of living in an enclave is important for migrants, and consequently, how English proficiency affects residence in different types of enclave is an empirical question. To the best of our knowledge, we are the first to study the causal impact on residency in language and region-of-birth enclaves, and more generally residency in four different types of enclave using the same dataset in a comparable manner.

Another contribution to highlight is that we use administrative boundaries as the geographical unit when constructing our measures of enclave, unlike Bleakley and Chin (2010) who base their analysis on public-use microdata areas (PUMAs), which are census-created geographies that contain no fewer than 100,000 individuals. PUMAs and counties coincide only around five per cent of the times. This is an important distinction when analysing the impact of language skills, because for language skills to affect people's residential locations, we require an assumption that people can interact and meet other people from their same language, ethnicity, or country or world region of origin. In this regard, the use of administrative boundaries is arguably more suitable than census-created boundaries, as transport communications are likely to exist within an administrative area, allowing individuals to interact more easily. For example, Bailey et al. (2020) find evidence that travel time has more of an importance than distance in formulating and maintaining social links across individuals, and suggest public transit infrastructure as an important determinant of social connectedness.

In contrast to the relation between language skills and residential segregation which has been intensively studied, we are not aware of any research that has analysed the relation between language proficiency and the quality of the neighbourhood migrants live in. Existing research related to this topic is the studies on neighbourhood quality which concerns its impact on social and labour market outcomes (e.g., Katz et al., 2001; Edin et al., 2003; Oreopoulos, 2003; Weinberg et al., 2004; Kling et al., 2007; Bertrand et al., 2000; Sanbonmatsu et al., 2011; Ludwig et al., 2012; Damm, 2014; Weinhardt, 2014). For example, Damm (2014) studies the effect of neighbourhood quality on labour market outcomes of refugees in Denmark, using a natural experiment which quasi-randomly assigned refugees to municipalities. Her results indicate no impact of living in a low-quality neighbourhood on labour market outcomes. Interestingly, however, the skill level of non-Western immigrant neighbours and the employment rate of co-national neighbours are found to have a significant impact on the labour market outcomes of refugees. Based on these results, Damm (2014) suggests that residence-based networks on

job information are ethnically stratified.

In our study, instead of looking at the consequences of living in a neighbourhood of a certain quality, we analyse whether language proficiency of migrants leads them to live in a neighbourhood of a certain quality. By doing so, our paper bridges two strands of literature, one on language proficiency of migrants and the other on neighbourhood quality. It is important to analyse residential environments of migrants with different English skills, as lower language proficiency can have amplifying negative effects if it not only directly affects immigrants' social and labour market outcomes,<sup>2</sup> but also indirectly affects their outcomes via neighbourhood effects in their residential environment. For example, using US census tract data, Topa (2001) finds positive spatial correlations of unemployment across neighbouring tracts, and suggests that there are local spillover effects of unemployment due to (presence or absence of) a local exchange of job information. This implies that, if lower language skills lead migrants to live in a high-unemployment neighbourhood, the employment opportunities for these migrants could be reduced due to neighbourhood effects, in addition to the direct effect of poor language proficiency itself on unemployment found in the existing literature (e.g., Gonzalez, 2005; Clausen et al., 2009).

### 3. Identification strategy

We estimate the causal effect of English language proficiency on residential outcomes of childhood migrants in England and Wales by regressing these outcomes on a measure of English proficiency, controlling for various individual characteristics and fixed effects. We begin by presenting the empirical framework used in Bleakley and Chin (2010). Subsequently, we will introduce our empirical specification and highlight the differences from the specification in this seminar paper:

$$outcome_{ica} = \alpha_0 + \alpha_1 prof_{ica} + X'_{ica} \xi + \gamma_c + \delta_a + u_{ica} \quad (1)$$

where  $outcome_{ica}$  represents the outcome of individual  $i$  born in country  $c$  who arrived in the UK at age  $a$ , and  $prof_{ica}$  is a measure of English language skills. The individual characteristics,

---

<sup>2</sup>There are a large number of studies on the relation between language proficiency and social and labour market outcomes. For example, there is a strand of literature on the relation with earnings (e.g., Chiswick 1998; Dustmann 1994; Chiswick and Miller 1995; Dustmann and van Soest 2001, 2002; Shields and Price 2002; Dustmann and Fabbri 2003; Bleakley and Chin 2004; Di Paolo and Raymond 2012; Miranda and Zhu 2013; Budria and Swedberg 2015), with employment (e.g., Miller and Neo 1997; Gonzalez 2005; Clausen et al. 2009; Yao and van Ours 2015), with education (e.g., Glick and White 2003; Bleakley and Chin 2010; Akbulut-Yuksel et al. 2011; Aoki and Santiago 2018), and with health (e.g., Miranda et al. 2011; Bauer et al. 2012; Kimbro et al. 2012; Lee et al. 2013; Guven and Islam 2015; Clarke and Ispording 2017; Aoki and Santiago 2018), among others.

$X_{ica}$ , and the parameter  $\xi$  are  $K \times 1$  vectors, where  $K$  is the number of variables capturing individual characteristics such as age and sex.  $\gamma_c$  and  $\delta_a$  are country-of-birth and age-at-arrival fixed effects, respectively, and  $u_{ica}$  is the disturbance term.

The main coefficient of interest is  $\alpha_1$ , which measures the effect of English skills on the residential outcomes of migrants. An econometric challenge to estimate equation (1) is the endogeneity of English skills. First, the residential locations of migrants may affect their English skills (reverse causality). Second, unobserved individual characteristics, such as ability, may be correlated with both English skills and our outcome variables. Third, our self-reported measure of language proficiency may contain measurement error. Thus, using OLS to estimate  $\alpha_1$  is unlikely to produce a causal estimate of the effect of English proficiency. To identify the causal effect, equation (1) can be estimated using the IV estimator, which requires an IV that gives exogenous variation in English skills. To construct an IV for language skills, we exploit age at arrival in the host country, following Bleakley and Chin (2004; 2010). The idea of using age at arrival is based on the “critical period of language acquisition” hypothesis (Lenneberg, 1967), which states that an individual exposed to a new language during the critical period of language acquisition (childhood) can learn the language relatively easily.<sup>3</sup> This hypothesis implies that, among migrants from a non-Anglophone country, those who arrive in the UK at a young age can learn English relatively easily, while those who arrive at an older age find it harder to learn English.

For a variable to serve as an IV for English skills, the following assumptions are required: (i) it does not appear in equation (1) and (ii) it is uncorrelated with any other determinants of the residential outcomes of migrants apart from proficiency in English. However, age at arrival per se is unlikely to satisfy these assumptions because it affects the extent of assimilation apart from language acquisition. For example, age at arrival may affect immigrant residential outcomes through knowledge about living conditions in different neighbourhoods in the host country. To overcome this problem, Bleakley and Chin (2004; 2010) use as an instrument, the interaction of age at arrival with an indicator variable for being born in a non-Anglophone country:

$$\phi_{ica} \equiv \max(0, a_i - \text{cutoff}) \times I(c \text{ is non-Anglophone}) \quad (2)$$

where  $a_i$  is age at arrival for individual  $i$ ;  $\text{cutoff}$  is the value of cut-off age;<sup>4</sup> the function  $\max(0, a_i - \text{cutoff})$  corresponds to the additional years after cut-off age for those who arrived in the host country after cut-off age, and zero otherwise; and  $I(c \text{ is non-Anglophone})$  is an in-

---

<sup>3</sup>Lenneberg (1967) observes that, until early teens, individuals have an innate flexibility for the organisation of brain functions necessary for the acquisition of a language. If basic language skills have not been acquired by puberty, they tend to remain deficient for the rest of their life because the ability to adjust to physiological demands for verbal acquisition declines sharply after puberty due to physiological changes in the brain.

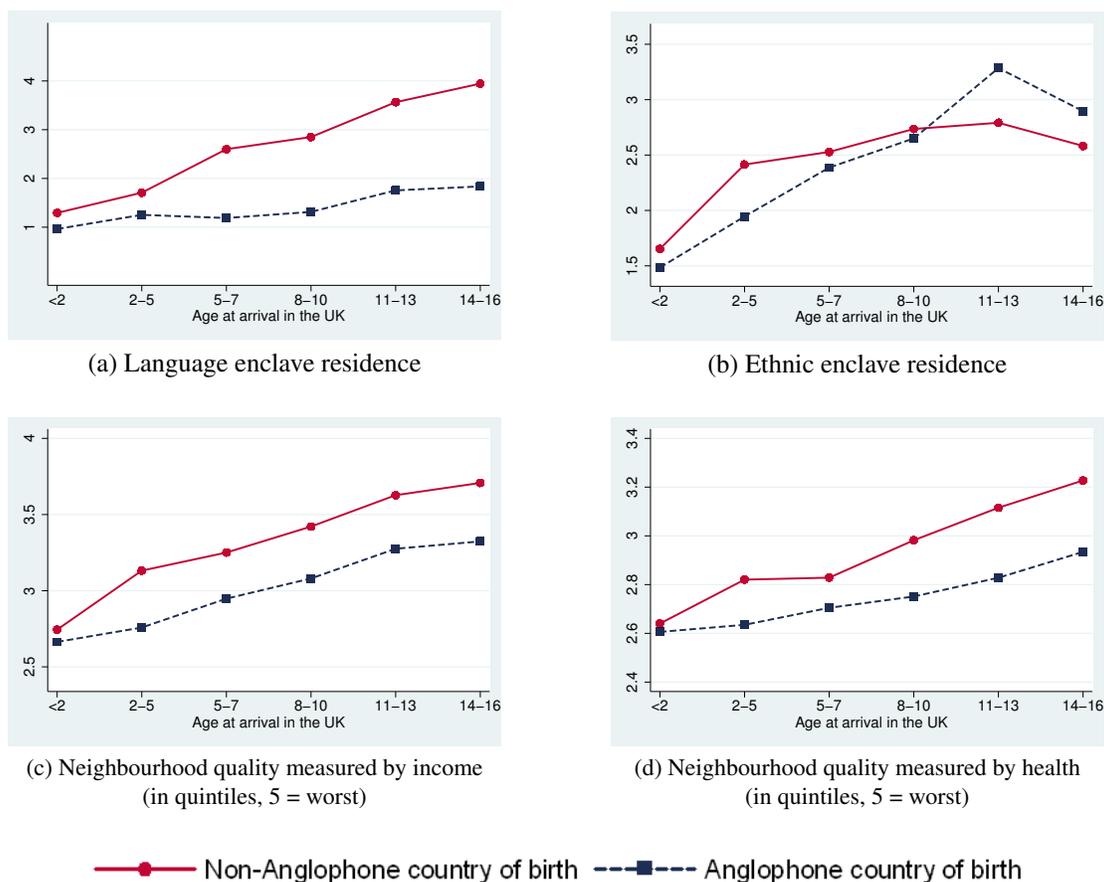
<sup>4</sup>Bleakley and Chin (2004) use the cut-off age of 11, and Bleakley and Chin (2010) use the age of 9.

indicator variable that takes the value of one if the individual is from a non-Anglophone country, and zero otherwise. The idea behind equation (2) is as follows: on arrival in the host country, all migrants are exposed to a new environment, but only those born in non-Anglophone countries encounter a new language. Thus, conditional on individual characteristics, differences in outcomes of early- and late-arrivers from Anglophone countries would only reflect age-at-arrival effects, while differences in outcomes of migrants from non-Anglophone countries would reflect those same age-at-arrival effects and an additional effect, the language effect. Therefore, a difference in the outcomes between early- and late-arrivers born in non-Anglophone countries in excess of the corresponding difference for migrants born in Anglophone countries can arguably be attributed to the effect of language.

Figure 1 plots mean immigrant residential outcomes by age at arrival, where these outcomes are extent of language residential segregation (panel (a)) and of ethnic residential segregation (panel (b)), and quality of the neighbourhood migrants live in, measured by income of residents (panel (c)) and health of residents (panel (d)).<sup>5</sup> The solid and dashed lines correspond to migrants from non-Anglophone and Anglophone countries, respectively. Focusing on the first two panels (a) and (b), early arrivers follow similar trends, but for late arrivers, the two series diverge. Later arrivers from non-Anglophone countries tend to cluster in areas with a higher concentration of residents who speak their native language (panel (a)), but not in areas with a higher concentration of residents with the same ethnicity (panel (b)), relatively speaking to Anglophone migrants. The last two panels (c) and (d) indicate that later arrivers from non-Anglophone countries tend to live in the neighbourhoods, where residents are relatively more deprived. An interesting observation from Figure 1 is that migrants from Anglophone countries also exhibit age-at-arrival effects. This observation implies that, apart from the effect of language, age at arrival is likely to have direct effects on migrants' residential outcomes, confirming that age at arrival per se is not a valid instrument and it is important to control for age-at-arrival fixed effects in the empirical specification.

---

<sup>5</sup>As there are numerous outcome variables, we do not report graphs for every outcome to save space. Instead, we report the relation between age at arrival and each immigrant outcome (i.e., reduced-form estimates) in Table 2.

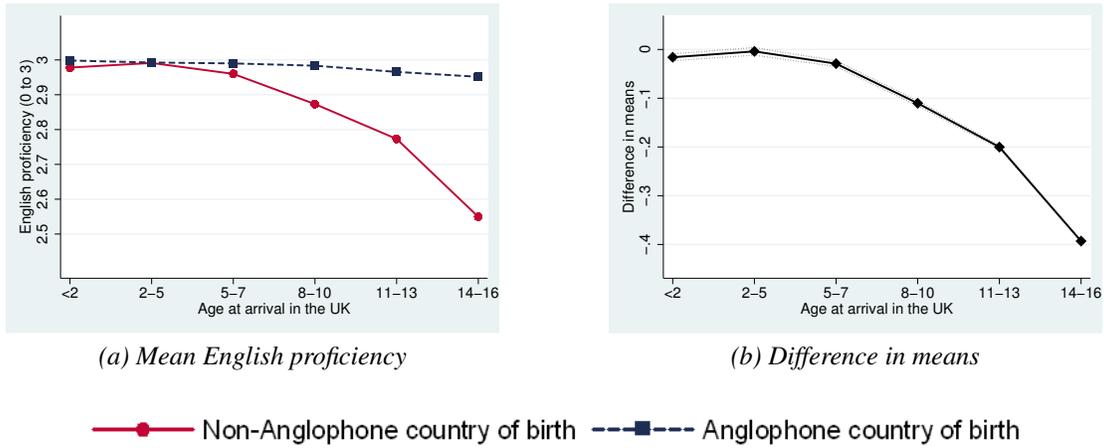


**Fig. 1.** Immigrant outcomes by age at arrival.

Notes: Immigrant outcomes are plotted by age at arrival where the outcomes are extent of residential segregation in terms of the main language spoken by residents (panel (a)) and of ethnicity of residents (panel (b)); and quality of the neighbourhood where migrants live measured by income (panel (c)) and health (panel (d)) of residents. Each outcome is regression adjusted for age and sex. The sample corresponds to childhood migrants aged 20 or over at the time of the 2011 Census.

Source: ONS Longitudinal Study.

Figure 2 shows the relation between English language proficiency and age at arrival of migrants who arrived in the UK during their childhood. Figure 2 shows that migrants born in Anglophone countries score between 2.9 and 3 in the ordinal measure of English proficiency, where 3 corresponds to “speaks very well”, and are generally proficient in English irrespective of their age at arrival. In contrast, for those born in non-Anglophone countries, the later they arrived, the poorer their English is on average, which is consistent with the critical period hypothesis. Panel (b) indicates that the two series start diverging at around ages five to seven, and the difference becomes more salient from ages eight to 10.



**Fig. 2.** Age at arrival and English proficiency.

Notes: Panel (a) plots the average ordinal measure of English proficiency, where 3, 2, 1, and 0 correspond to speak English "very well", "well", "not well", and "not at all", respectively. Panel (b) plots the difference in the two sets of means reported in panel (a). The sample corresponds to childhood migrants aged 20 or over at the time of the 2011 Census.

Source: ONS Longitudinal Study.

Based on the observed pattern in Figure 2, we could have used the Bleakley and Chin's instrument,  $\phi_{ica}$ , in equation (2). However, there are potentially important issues with using  $\phi_{ica}$ . The first issue concerns cohort effects: it is possible that cohort quality have changed over time and in a different manner between migrants from non-Anglophone and Anglophone countries. If this change coincided with the critical age cut-off,  $\phi_{ica}$  will capture not only the effect of age at arrival on English proficiency, but also of differential changes in cohort quality between the two sets of migrants. As equation (1) controls for observable individual characteristics, any observable differences have been accounted for. However, if there are any unobservable cohort differences, they could potentially bias resulting IV estimates. This is a potential issue that a large number of papers using this type of identification strategy face (e.g., Akbulut-Yuksel et al., 2011; Miranda and Zhu, 2013; Guven and Islam, 2015; Yao and van Ours, 2015).

In an attempt to mitigate this type of concerns, instead of using  $\phi_{ica}$  as an instrument for English skills, we use it as a control variable in our empirical model. Controlling for  $\phi_{ica}$  partially addresses this concern, by absorbing systematic differences in cohort outcomes between the two sets of migrants which coincided with the critical age cut-off, if any. We then instead construct our instrument, exploiting the variation in linguistic distance from English within non-Anglophone countries, following Clarke and Ispording (2017). More specifically, among non-Anglophone migrants, there is heterogeneity in how similar their native languages are to English, and this may affect their capacity to become proficient in English. For example, an immigrant with a language that is more similar to English (e.g., Dutch) will find it easier to

learn English than an immigrant with a native language that is very different to English (e.g., Vietnamese). Exploiting this heterogeneity, we construct our instrument,  $\theta_{ica}$ , by interacting age at arrival with linguistic distance,  $ldist_c$ , between English and the native language in an origin country,  $c$ :

$$\theta_{ica} \equiv \max(0, a_i - 8) \times ldist_c \quad (3)$$

This specification is different from  $\phi_{ica}$  in equation (2),<sup>6</sup> which does not allow for the age-at-arrival effects to differ by linguistic distance. Not surprisingly, the cost we must pay for using  $\phi_{ica}$  as a control and  $\theta_{ica}$  as an instrument, is that part of the variation in our instrument,  $\theta_{ica}$ , will be absorbed by  $\phi_{ica}$ , since  $\theta_{ica}$  and  $\phi_{ica}$  are highly correlated. We will see in Section 5, however, that there is still enough variation left in our data to allow us to estimate our empirical model. We are not aware of any other studies of the causal effects of language which uses  $\phi_{ica}$  as a control variable rather than an instrument.

Next, we implement a further change to the specification used in Bleakley and Chin (2010). Namely, we add to our model, year-of-arrival fixed effects that vary by individual. Recent literature has documented substantial heterogeneity in year-of-arrival effects among migrants, with more recent waves of migrants doing worse. Borjas (2015), for example, provides suggestive evidence to this fact, arguing that more recent immigrant cohorts come from a different part of the skill distribution. This is an important concern if the two sets of migrants arrived in the UK in systematically different years, as it can lead for the two sets of migrants to achieve different outcomes in the UK. Another important implication here is that, by accounting for year-of-arrival fixed effects, we are implicitly accounting for length of stay in the UK, since we can obtain length of stay by subtracting the year of arrival from the year at which data are obtained, 2011.<sup>7</sup> This is important since English skills at a specific point in time of individuals, who arrived in the UK at the same age but in different years, might differ because of the time they had to learn the language or the conditions they faced throughout their stay.

All in all, our main equation of interest is specified as follows:

$$outcome_{ica} = \beta_0 + \beta_1 prof_{ica} + \beta_2 \phi_{ica} + X'_{ica} \zeta + \iota_i + \chi_c + \alpha_a + \varepsilon_{ica} \quad (4)$$

where  $outcome_{ica}$  represents the residential outcome of individual  $i$  born in country  $c$  who arrived in the UK at age  $a$ ;  $prof_{ica}$  is a measure of English language skills; and  $\phi_{ica}$  is as defined in

<sup>6</sup>Age eight is chosen as the cut-off value because Figure 2 indicates that the two series start diverging at ages five to seven, and the gap becomes more salient at ages eight to 10 in our dataset.

<sup>7</sup>Length of stay can also be obtained by subtracting age at arrival from age at the time of the 2011 census, both of which are included in our model.

equation (2).  $\beta$ s are parameters to be estimated, and in particular  $\beta_1$  is our parameter of interest. The individual characteristics  $X_{ica}$  including age fixed effects and sex, and the parameter  $\zeta$  are  $K \times 1$  vectors, where  $K$  is the number of variables capturing individual characteristics.  $\iota_i$ ,  $\chi_c$ , and  $\alpha_a$  are year-of-arrival, country-of-birth, and age-at-arrival fixed effects, respectively, and  $\varepsilon_{ica}$  is the disturbance term. Our first-stage equation, which relates English proficiency and age at arrival, using our instrument  $\theta_{ica}$  in equation (3), can be specified as follows:

$$prof_{ica} = \beta_{f0} + \beta_{f1}\theta_{ica} + \beta_{f2}\phi_{ica} + X'_{ica}\zeta_f + \iota_{fi} + \chi_{fc} + \alpha_{fa} + \varepsilon_{fica} \quad (5)$$

where the additional letter  $f$  in subscripts refer to 'first stage'.

## 4. Data and sample

### 4.1. Data

To analyse the impact of English language skills on immigrant outcomes, we use an individual-level dataset from the ONS Longitudinal Study of England and Wales, which contains linked census and life events data for a one-per-cent sample of the population of England and Wales. Unless indicated otherwise, our individual characteristics are collected from the 2011 Census sample within the ONS Longitudinal Study, including information on English proficiency, which is a self-reported ordinal variable that takes values 3, 2, 1, and 0, corresponding to speaks English “very well”, “well”, “not well”, or “not at all”, respectively. From this variable, we derive our measure of English skills, which is the indicator variable taking the value of one if the respondent speaks English “very well”, and zero otherwise. We also extract our measure of parental education from the ONS Longitudinal Study by tracking the individuals in our dataset through all censuses contained in the Longitudinal Study. Once we have identified their parents, we assign them to the individuals in our sample.

To create our instrument for English proficiency, we exploit two census variables, country of birth and age at arrival in the UK,<sup>8</sup> and a measure of linguistic distance between English and origin-country language. We measure linguistic distance using a variation of the Levenshtein distance computed by Ispording and Otten (2014). Following a procedure to evaluate phonetic similarity between different languages, developed by the Max Planck Institute for Evolutionary Anthropology, they compute the extent of similarity between languages in percentages. The

---

<sup>8</sup>Age at arrival in the UK is derived from the date that a person last arrived to live in the UK and their age. Short visits away from the UK are not counted in determining the date that a person last arrived. The age of arrival is only applicable to usual residents who were not born in the UK and does not include usual residents born in the UK who have emigrated and since returned.

measure of linguistic distance they construct is a standardised and continuous measure of the distance between languages based on phonetic similarity, where a higher number indicates a greater linguistic distance. Despite its purely descriptive nature that does not require any prior knowledge on language relations, this measure of linguistic distance is highly correlated with other linguistic distance measures such as those developed by linguists based on language families. We assign linguistic distance based on the official language in the country of birth of migrants. In the case of migrants born in a multilingual country, we assign the predominant native language of the country. For migrants born in a country where English is an official language and the predominant language spoken, we assign linguistic distance of zero. The following sub-sections detail the construction of our outcome variables.

### ***Residential segregation***

After linking the ONS Longitudinal Study to the local-authority level data from the 2011 Census obtained from ONS Nomis,<sup>9</sup> we construct the measure of residential segregation using an index of relative clustering following Borjas (2000), defined as:

$$Relative\ Clustering\ Index_{ij} = \frac{N_{ij}/N_j}{N_i/N} \quad (6)$$

where  $i = 1, \dots, I$  represents the languages and  $j = 1, \dots, J$  represents the local authorities.  $N_{ij}$  is the total number of persons reporting language  $i$  as their main language and living in local authority  $j$ ,  $N_j$  is the total number of persons living in local authority  $j$ ,  $N_i$  is the total number of persons reporting language  $i$  as their main language in England and Wales, and  $N$  is the total population in England and Wales. This relative clustering index is based on the 'exposure index', corresponding to the numerator of equation (6), which gives the fraction of people in a local authority reporting a particular language as their main language. Although the exposure index is widely used in the literature that studies immigrant enclaves (e.g., Borjas 2000; Edin et al. 2003; Bauer et al. 2005), a problem with this index is that it can underweight the available contacts for small ethnic groups. The relative clustering index is a better measure (Bertrand et al., 2000), since it deflates the exposure index by the proportion of people reporting a particular language  $i$  in the whole of England and Wales (i.e., the denominator of equation (6)). The relative clustering index in equation (6) captures the share of individuals reporting the same native language in the local authority where an immigrant lives in. It takes value one if the proportion of people speaking language  $i$  living in local authority  $j$  is the same as the proportion of people speaking that language in England and Wales. If the relative clustering index is

<sup>9</sup>The 2011 Census data for local authorities can be downloaded from ONS Nomis: <https://www.nomisweb.co.uk/>.

greater than one, then the group of individuals speaking that language is overrepresented in that particular local authority, whereas if the index is smaller than one, the group is underrepresented in that particular local authority.

In addition to measuring immigrant segregation based on their main language, we measure it based on their ethnicity, country of birth, and world region of birth. Each of these measures captures residential segregation along different dimensions: An ethnic group includes anyone who reports having a particular ethnic group, irrespective of whether they were born in the UK, whereas a country-of-birth group only includes individuals born in a particular country. A world-region-of-birth group includes not only those born in the same country, but also those born in the same world region of birth (e.g., North Africa, South America, etc.), since migrants may congregate not necessarily only with compatriots but also with individuals from surrounding countries of their country of origin. These different measures of segregation allow us to investigate whether and how much English language skills affect these different dimensions of immigrant residential segregation.

The geographical unit we use for our analysis is the local authority district, which is an administrative division in the UK. There were 348 local authority districts in England and Wales at the 2011 Census, which sum to at least 2,000 population with an average size of approximately 160,000 individuals. Using this geographical unit has some advantages. First, it is large enough: This is important because an individual does not necessarily interact with his immediate neighbours, but may have different networks of people (e.g., family, friends and colleagues) with whom they can interact frequently provided they have easy access to them, which happens if they live within a reasonable distance. In addition, choosing small areas could create measurement error problems in the case of immigrant groups with few observations. The second advantage of using local authorities is that they are not too large, as is the case with regions, which are too large to allow us to make the assumption that individuals could interact and meet other individuals from their same language, country or ethnic group. The third advantage is that local authority districts are administrative divisions. This is very important as it ensures that transport communications are likely to exist and be easily accessible. This latter motive makes an administrative division better than a census-created division for the purpose of capturing possible interactions with other group members. For example, Bailey et al. (2020) find evidence that travel time has more of an importance than distance in formulating and maintaining social links across individuals, and suggest public transit infrastructure as an important determinant of social connectedness. In this respect, we provide an alternative approach to Bleakley and Chin (2010), who also analyse the impact of English skills on residential segregation, but use census-created geographies that contain no fewer than 100,000 individuals called PUMAs, which do not coincide with administrative geographic boundaries. Using administrative boundaries is

arguably better as it makes it more likely that both workplace and residential interactions are taken into account, and both types of interactions can affect the decisions of migrants about where to live.

### *Neighbourhood quality*

We measure neighbourhood quality using data from the English Indices of Deprivation 2015, which are published by Ministry of Housing, Communities and Local Government (2015). These indices measure relative neighbourhood quality at a small-area level, called the ONS Lower-layer Super Output Areas (LSOAs). LSOAs are small areas designed to be of similar population size with a minimum of 1,000 individuals and a maximum of 3,000 (between 400 and 1,200 households), which have an average of approximately 1,500 residents or 650 households. We have matched our individuals to these indices corresponding to the area in which they were living at the time of the 2011 Census.

Three domains of the English Indices of Deprivation are exploited: Income deprivation, employment deprivation, and health deprivation. Income deprivation measures the proportion of population experiencing low income in the neighbourhood, while employment deprivation measures the proportion of working-age population who are involuntarily excluded from the labour market. For each of these domains, multiple indicators are used to measure the extent of deprivation,<sup>10</sup> are constructed as non-overlapping counts, and are summed together to pin down the total at-risk population for the domain.<sup>11</sup> This total at-risk population is then used to calculate the proportion of population experiencing that form of deprivation. Health deprivation is intended to measure age and sex specific premature mortality and the population experiencing the impairment of quality of life due to poor physical or mental health. Unlike income and employment deprivation, a single ordinal measure of deprivation, i.e., the proportion of at-risk population, cannot be calculated. Thus, factor analysis is used to generate appropriate weights for combining the standardised indicators into a single ordinal score of health deprivation.<sup>12</sup>

---

<sup>10</sup>The indicators used to measure the extent of income deprivation are adults and children in income support families; in income-based job-seeker's allowance families; in income-based employment and support allowance families; in pension credit families; and in working tax credit or child tax credit families not already counted (i.e., those who are not in receipt of the first four support/allowances); and asylum seekers in England in receipt of subsistence support and/or accommodation support. For employment deprivation, the following indicators are used: women aged 18 to 59 and men aged 18 to 64 who claim job-seeker's allowance, employment and support allowance, incapacity benefit, severe disablement allowance, and carer's allowance.

<sup>11</sup>For detailed definitions of indicators used to construct each measure of deprivation as well as further technical details of calculations, refer to the English Indices of Deprivation 2015 Technical Report (Ministry of Housing, Communities and Local Government, 2015).

<sup>12</sup>For health domain, indicators used to measure deprivation are an age and sex standardised measure of premature death; standardised morbidity-disability ratio; and standardised rate of emergency admission to hospital; and a composite measure on mood and anxiety disorders based on the rate of adult suffering from mood and anxiety disorders, and on the data on hospital episodes, suicide mortality and health benefits.

For each of the deprivation domains, quintiles are calculated, ranking the 32,844 LSOAs in England from least deprived to most deprived and dividing them into five equal groups. We create one variable for each domain, and each of these variables takes values 1 to 5, where 1 corresponds to the least deprived area and 5 corresponds to the most deprived area. Note that, although we would have liked to use a finer level of ranking such as deciles, quintiles are the smallest ranking we could get access to when dealing with data at the LSOA-level due to data confidentiality.

For the analysis of neighbourhood quality, we only use a sample of individuals who were living in England at the time of the 2011 Census. This is because, although there are the Welsh Indices of Deprivation, these indices measure relative levels of neighbourhood quality within Wales, and thus the Indices of England and that of Wales are not directly comparable.

## **4.2. Sample**

### ***Age restriction***

Our sample consists of individuals in the ONS Longitudinal Study dataset who were present in the 2011 Census, aged 20 or older at the time of the 2011 Census, and are childhood migrants, defined as individuals born outside of the UK who moved into the UK at age 16 or earlier. We impose this age-at-arrival restriction and assume that these childhood migrants did not make a migration decision on their own, but moved into the country following their parents or guardians.

### ***Country classification***

To implement our identification strategy, we include two types of migrants in our sample: (i) individuals born in a non-Anglophone country where English is not an official language (treatment group) and (ii) individuals born in an Anglophone country (control group). We classify a country as Anglophone if English is an official language and the predominant language spoken in the country.<sup>13</sup> We exclude from our sample individuals born in countries where English is an official language but not the predominant language spoken because it is not clear to what extent they were exposed to English prior to their arrival in the UK. This rule drops migrants from countries such as India and Pakistan who account for significant proportions of UK migrants.

The Appendix and Table 1 present a list of countries of birth for the migrants in our sample and summary statistics, respectively, for Anglophone and non-Anglophone countries by age-at-arrival group. Table 1 presents summary statistics where mean values and standard deviations are reported in parentheses. Panel A presents individual characteristics. A key observation is

---

<sup>13</sup>The World Almanac and Book of Facts 2011 is used to classify countries.

**Table 1** Immigrant characteristics and residential outcomes.

	Born in non-Anglophone country			Born in Anglophone country		
	Arrived aged 0 - 8	Arrived aged 9 - 16	Total	Arrived aged 0 - 8	Arrived aged 9 - 16	Total
<i>A. Individual characteristics</i>						
English proficiency, ordinal measure	2.971 (0.203)	2.673 (0.614)	2.805 (0.500)	2.993 (0.099)	2.963 (0.210)	2.979 (0.162)
English proficiency, = 1 if speaks very well	0.977 (0.150)	0.744 (0.437)	0.847 (0.360)	0.995 (0.073)	0.967 (0.178)	0.982 (0.134)
Linguistic distance	0.931 (0.104)	0.962 (0.076)	0.948 (0.091)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
Age	38.634 (16.025)	33.907 (15.675)	35.990 (16.002)	44.320 (14.029)	47.050 (16.699)	45.606 (15.404)
Female	0.511 (0.500)	0.509 (0.500)	0.510 (0.500)	0.513 (0.500)	0.540 (0.498)	0.526 (0.499)
White	0.625 (0.484)	0.428 (0.495)	0.515 (0.500)	0.691 (0.462)	0.378 (0.485)	0.544 (0.498)
Black	0.093 (0.291)	0.200 (0.400)	0.153 (0.360)	0.140 (0.347)	0.346 (0.476)	0.237 (0.425)
Asian	0.115 (0.319)	0.209 (0.406)	0.167 (0.373)	0.133 (0.339)	0.233 (0.423)	0.180 (0.384)
Other single race	0.118 (0.322)	0.117 (0.321)	0.117 (0.322)	0.008 (0.087)	0.009 (0.094)	0.008 (0.090)
Multiracial	0.045 (0.208)	0.039 (0.194)	0.042 (0.201)	0.026 (0.160)	0.032 (0.176)	0.029 (0.168)
<i>B. Enclave residency</i>						
Language enclave	2.038 (4.471)	3.781 (6.482)	3.013 (5.749)	1.269 (3.668)	1.810 (3.925)	1.524 (3.800)
Country-of-birth enclave	4.238 (6.765)	5.421 (6.897)	4.900 (6.863)	2.290 (3.421)	3.866 (5.202)	3.033 (4.422)
World-region-of-birth enclave	2.374 (2.616)	2.898 (2.728)	2.667 (2.692)	1.833 (2.124)	2.735 (2.569)	2.258 (2.387)
Ethnic enclave	2.290 (2.434)	2.712 (2.267)	2.525 (2.351)	1.988 (2.211)	3.023 (2.904)	2.476 (2.613)
<i>C. Neighbourhood quality index (in quintiles, 5 = worst)</i>						
Income deprivation index	3.083 (1.454)	3.658 (1.336)	3.404 (1.418)	2.820 (1.384)	3.274 (1.386)	3.034 (1.403)
Employment deprivation index	2.901 (1.402)	3.369 (1.377)	3.162 (1.407)	2.707 (1.372)	3.069 (1.394)	2.878 (1.394)
Health deprivation index	2.779 (1.364)	3.173 (1.367)	2.999 (1.380)	2.670 (1.377)	2.867 (1.369)	2.763 (1.377)

Notes: The sample consists of individuals in the ONS Longitudinal Study dataset aged 20 or over who lived in England and Wales at the time of the 2011 Census, and were born outside the UK who arrived in the UK at age 16 or before. The number of observations varies by panel and column: Panels A and B have 2,005; 2,545; 4,550; 3,158; 2,815 and 5,973 observations in the first to sixth columns, respectively, except for ethnic enclave (1,999; 2,526; 4,525; 3,148; 2,809 and 5,957) due to 40 missing values of ethnicity. Panel B has 1,972; 2,494; 4,466; 3,087; 2,768 and 5,855 observations. Source: Authors' calculations based on the ONS Longitudinal Study.

that, for Anglophone migrants, the proportions of individuals who speak English “very well” for early- and late-arrivers are high (97 to 100 per cent) and similar as one would expect. In contrast, for migrants born in non-Anglophone countries, late-arrivers show a lower share of people who speak English “very well” (74 per cent) than early-arrivers (98 per cent). This latter group has a proficiency level similar to migrants born in Anglophone countries. Linguistic distance (from English) is zero for Anglophone countries by construction, and it takes positive values for non-Anglophone countries. Turning to residential outcomes, late-arrivers born in non-Anglophone countries live in the areas with higher concentrations of people who speak their same native languages and from the same countries of birth (panel B), and in the neighbourhoods of lower quality measured in terms of income, employment and health of residents (panel C).

## 5. Results

We begin by estimating equation (4) using the OLS estimator.<sup>14</sup> Column (1) of Table 2 reports the OLS estimates of the effect of English proficiency on the residential outcomes of childhood migrants in England and Wales, after controlling for individual characteristics, the interaction of age at arrival with the dummy for coming from a non-Anglophone country as defined in equation (2), and year-of-arrival, country-of-birth and age-at-arrival fixed effects. The results for enclave residence in panel A, and for neighbourhood quality in panel B all indicate that poorer English skills are significantly associated with living in an enclave, irrespective of the types of enclave, and in a neighbourhood of lower quality measured in terms of income, employment and health of residents.

A problem with the OLS estimates of the effects of English proficiency is that they are biased if English proficiency is endogenous. To address this potential endogeneity issue, equation (4) is estimated using the IV estimator, where we use as an instrument for English skills, the *interaction* of the excess age at arrival from age eight with linguistic distance between English and origin-country language (see equation (3)). The first-stage estimates indicate that, for migrants born in a linguistically distant country, each year past age eight at arrival significantly decreases their likelihood of speaking English “very well” by about 0.148 on average (column

---

<sup>14</sup>Our measure of English language skills is an indicator variable for speaking English “very well” as described in Section 4. We also try using an indicator for speaking English “very well” or “well”. When the indicator is constructed to take the value of one if one speaks English “very well” or “well”, the variation of this variable is not as large since the majority of individuals in our sample reported to speak English either “very well” or “well”. First-stage and IV estimates become smaller and larger, respectively, in absolute terms, although we do not observe any qualitative changes in the results. In addition, we try using the original Census ordinal measure of English skills ranging from zero to three, where three corresponds to speaks English “very well”. The results using this alternative measure of English language skills, presented in the online Appendix, are again qualitatively similar to our main results.

**Table 2**  
OLS, IV, reduced-form, and first-stage estimates.

Dependent variable:	Enclave, neighbourhood quality			Dummy for English ability
	OLS	IV	Reduced-form	First-stage
	(1)	(2)	(3)	(4)
<i>A. Enclave residency</i>				
Language enclave	-4.420*** (0.873)	-3.574* (2.016)	0.528* (0.290)	-0.148*** (0.015)
Country-of-birth enclave	-1.330*** (0.437)	0.807 (1.908)	-0.119 (0.290)	-0.148*** (0.015)
World-region-of-birth enclave	-0.278** (0.128)	2.762*** (0.943)	-0.408*** (0.129)	-0.148*** (0.015)
Ethnic enclave	-0.335* (0.173)	2.987*** (1.010)	-0.444*** (0.087)	-0.148*** (0.015)
<i>B. Neighbourhood quality index (in quintiles, 5 = worst)</i>				
Income deprivation	-0.416*** (0.075)	-0.914* (0.538)	0.134* (0.075)	-0.147*** (0.015)
Employment deprivation	-0.372*** (0.072)	-1.376* (0.717)	0.202** (0.093)	-0.147*** (0.015)
Health deprivation	-0.263*** (0.069)	-1.502** (0.596)	0.221*** (0.076)	-0.147*** (0.015)

Notes: Standard errors are clustered by country of birth. OLS and IV are the estimates of  $\beta_1$  in equation (4). Reduced-form and first-stage are the estimates of the coefficient on the interaction of linguistic distance between origin-country language and English, with age at arrival. Rows in each panel correspond to regressions for different outcomes. Every regression controls for dummies for age, gender and race, and year-of-arrival, country-of-birth, and age-at-arrival fixed effects. Sample sizes for panels A and B are 10,522 and 10,320, respectively, except for ethnic enclave (10,482 observations due to 40 missing values). The F-statistics on the excluded instrument in column 4 range from 97 to 103. Source: Authors' calculations based on the ONS Longitudinal Study.

(4), Table 2). The magnitude of this estimate implies that a person's likelihood of speaking English very well would be lowered by approximately one if the person arrived from a non-Anglophone country at age 15 instead of at age eight. It is important for the identification that our instrument is not weak as a weak instrument is known to bias the IV estimator toward the probability limit of the corresponding OLS estimator. Stock et al. (2002) compute the critical value for the test of weak instruments based on the first-stage F-statistic, and suggest that an F-statistic above roughly 10 makes IV inferences reliable. According to their test for weak instruments, our instrument is not weak as the first-stage F-statistics on the excluded instrument range between 97 and 103.

Column (3) of Table 2 presents the reduced-form estimates of the effects of the instrument on residential outcomes. The results for residential segregation, reported in panel A, indicate that, for those born in a linguistically distant country who arrived in the UK after age eight, each additional year that passes before they arrive in the UK is correlated with living in an area with a higher concentration of people who speak same native language. In line with this reduced-form estimate, the causal estimate in column (2) shows that a poorer English proficiency significantly leads migrants to live in an area with a higher concentration of people who speak their native language. In contrast, we find a positive and significant impact on residence in world-region-of-birth and ethnic enclaves.

To understand our findings, let us take Spanish as an example. Our findings suggest that a white Spanish person who speaks English well clusters with white people (their ethnicity) or people from Europe (their world region of birth), but not necessarily with other Spanish speakers. Our results are in contrast to the majority of correlation studies which find negative associations between poor language skills and ethnic concentration (citations), and to that of Bleakley and Chin (2010) who find an insignificant impact of English proficiency on ethnic enclave residency. Both our study and that of Bleakley and Chin (2010) find an insignificant impact on residency in a country-of-birth enclave. This could be due to the fact that, for the group of individuals whose number of compatriots is small, there may not be enough fellow citizens to cluster with (e.g., Koreans). In this regard, world-region-of birth enclaves may be a better measure to capture concentrations of individuals with similar cultural backgrounds in the current context.

Regarding the magnitude of the effect, to facilitate an interpretation of our estimate for the effect on language residential segregation  $-3.574$ , consider the following hypothetical situation: Suppose that there is an immigrant born in a Spanish-speaking country who does not speak English very well and lives in the local authority with the relative language clustering index of 3.6 (meaning that there are roughly 3.6 times as many Spanish-speaking migrants in the local authority as one would have expected if the Spanish-speaking population had distributed itself

randomly across England and Wales). If this immigrant had spoken English “very well”, they would have lived in a local authority where Spanish-speaking migrants are neither over- nor under-represented (i.e., local authority with the relative clustering index of one).

Panel B reports neighbourhood quality outcomes. IV estimates in column (2) indicate that poor English proficiency leads migrants to live in a neighbourhood of lower quality, where residents are more deprived in terms of income, employment and health. The magnitudes of the effects are non-negligible: a one-standard-deviation increase in English skills of late arrivers from a non-Anglophone country lowers the quality of neighbourhood they live in, measured by quintiles, by roughly half of a unit. There appears to be strong evidence of the effects of language skills on neighbourhood quality.

When comparing OLS and IV estimates, there is no clear pattern for the residential segregation outcomes,<sup>15</sup> whereas IV estimates are larger in absolute terms (more negative) for neighbourhood quality outcomes. For example, IV estimate is roughly four times larger for neighbourhood quality measured by employment deprivation. It is possible that an omitted variable, such as ability, biases the OLS estimator downward (making the negative effect more negative), but at the same time measurement error possibly correlated with our measure of language proficiency biases the OLS estimator upward (making the negative effect less negative). For example, migrants surrounded by native English speakers may report their proficiency being poor, while those surrounded by other non-Anglophone migrants may report their proficiency being fluent irrespective of their true English proficiency. In fact, self-reported categorical language measures are found to contain substantial measurement error (Dustmann and van Soest, 2001). Bearing in mind that the estimated effects have negative signs, if the upward bias caused by measurement error, known as attenuation bias, outweighs the downward bias caused by unobserved characteristics, IV estimates will be larger in absolute terms (more negative) than OLS estimates, which can help explain the relatively larger IV effects for neighbourhood quality outcomes.

### **5.1. Role of education**

This subsection investigates a possible channel through which English proficiency affects the residential outcomes of migrants: namely, education. Apart from the direct effects of English proficiency on immigrant outcomes by facilitating communication with native residents and colleagues, English proficiency may also have indirect effects by improving the educational attainments of migrants (Aoki and Santiago, 2018). To further investigate this channel, in ad-

---

<sup>15</sup>The previous study using the US census self-reported measure of English proficiency to analyse the impact of English proficiency on country-of-birth and ethnic segregation outcomes did not find any clear pattern in the relative sizes of IV and OLS estimates either (Bleakley and Chin, 2010).

dition to English skills, we control for the measures of education in equation (4). As education is likely to be endogenous in our empirical model, the estimates of the impact of English skills no longer have causal interpretations. Nevertheless, we present these results in Table 3 to provide suggestive evidence of the role that education plays in explaining the effects of English skills. The first column reports the coefficient estimate on English skills, while the second to fourth columns report estimates on the measures of education. Education is measured by a set of dummy variables that take the value of one if the person has a compulsory qualification, a post-compulsory qualification or an academic degree as the highest level of education qualification obtained, respectively, and zero otherwise. The dummy variable for no qualifications is omitted from equation (4).

Table 3 indicates that the role of education varies by outcome. Concerning residential segregation, panel A indicates that the point estimates of the impact of English proficiency do not change much after controlling for education. It is also noticeable that holding an academic degree has no significant impact on the residential segregation outcomes. In contrast, panel B, which reports the results for neighbourhood quality, indicates that the magnitudes of the effects of English skills have been greatly diminished (e.g., roughly between three quarter to half of the original magnitudes), while education, in particular, holding an academic degree, has a non-negligible impact on the neighbourhood quality outcomes. Taken together, our sub-analysis supports the possibility that a better educational attainment as a result of better English skills is a key channel through which better language skills lead migrants to live in a neighbourhood of higher quality, but it is unlikely to be a key channel for residential segregation outcomes.

## **6. Robustness checks**

Having found the significant impact of English proficiency on the residential outcomes of UK migrants, we address various concerns that can threaten the validity of our identification strategy. The first concern to be addressed is that our main results are driven by differences in the background characteristics associated with linguistic distance, which are different from language. Two strategies are employed to deal with this issue: first, we explicitly account for various factors potentially associated with linguistic distance, including cultural distance between the UK and origin country (interacted with age at arrival); and second, we drop a set of countries that have special ties to the UK, in an attempt to make our sample less heterogeneous than the original sample. The second concern to be addressed is that our main results are driven by differences in parental characteristics. Third, we explore the possibility that migrants from linguistically distant countries from the UK, in which English is commonly used, find it easy

**Table 3**

Effects of English proficiency and education qualifications on immigrant outcomes.

	English skills	Compulsory education	Post-compulsory education	Academic degree
	(1)	(2)	(3)	(4)
<i>A. Enclave residency</i>				
Language enclave	-3.372 (2.181)	-0.815 (1.078)	-1.155 (1.072)	0.078 (0.214)
Country-of-birth enclave	0.941 (2.094)	-1.817 (1.538)	-2.192 (1.583)	0.210 (0.210)
World-region-of-birth enclave	2.869*** (1.020)	-0.840* (0.474)	-1.055** (0.516)	0.081 (0.080)
Ethnic enclave	3.212*** (1.128)	0.000 (0.000)	-0.305*** (0.104)	0.004 (0.073)
<i>B. Neighbourhood quality index (in quintiles, 5 = worst)</i>				
Income deprivation	-0.500 (0.595)	-0.330* (0.196)	-0.481** (0.211)	-0.364*** (0.041)
Employment deprivation	-0.982 (0.766)	-0.121 (0.227)	-0.226 (0.246)	-0.389*** (0.045)
Health deprivation	-1.177* (0.607)	-0.245 (0.226)	-0.296 (0.249)	-0.352*** (0.048)

Notes: Standard errors are clustered by country of birth. The estimates shown in column 1 of each panel are the IV estimates of  $\beta_1$  in equation (4), including all controls described in Table 2, in addition to the following controls for education: dummy variables for having compulsory, post-compulsory, and degree qualifications, respectively, as the highest level of qualification obtained. Estimates for these education controls are shown in columns 2 to 4. Sample sizes for panels A and B are 10,522 and 10,320, respectively, except for ethnic enclave (10,482 observations due to 40 missing values).

Source: Authors' calculations based on the ONS Longitudinal Study.

to acquire English skills, irrespective of linguistic distances between their mother tongues and English.

For our IV strategy to identify the causal effects of language skills, we require an assumption that, *within* those born in non-Anglophone countries, linguistic distance impacts an English language acquisition process only, conditional on individual characteristics, and age-at-arrival, country-of-birth, and year-of arrival fixed effects. Note that our specification allows for the age-at-arrival effects to differ between Anglophone and non-Anglophone migrants, as we control for a dummy variable for coming from a non-Anglophone country, interacted with age at arrival, in our specification. In this regard, our identification assumption is weaker than the assumption used in Bleakley and Chin (2004; 2010) and other studies that adopt the same identification strategy, which requires age-at-arrival effects to be the same for Anglophone and non-Anglophone migrants. Nevertheless, one could still question the validity of the assumption required for our identification strategy: within non-Anglophone migrants, linguistic distance may not only capture a different English language acquisition process, but also other aspects of assimilation process that might differently affect residential outcomes. So far, in addition to controlling for individual characteristics including age, sex and race, we have controlled for three types of fixed effects: (i) country-of-birth fixed effects which account for the country-of-origin specific effects, which do not vary by age at arrival; (ii) year-of-arrival fixed effects which absorb the systematic differences in cohort characteristics of those who arrived in the UK in different years; and (iii) age-at-arrival fixed effects which absorb the corresponding differences of those who arrived in the UK at different ages. If migrants with different linguistic distances from English arrived in the UK at systematically different ages (e.g., linguistically distant migrants arriving at an early age), the last fixed effects will capture the effects on outcomes stemming from this difference.

To further address the concern of differences in the background characteristics associated with linguistic distance, we now start adding further potentially important controls. The first control to be added is the measure of cultural distance between the UK and origin country (interacted with age at arrival). It is possible that linguistic distance is associated with other aspects of cultural differences between the UK and origin country, and that migrants born in a country that is culturally more distant to the UK find it more difficult to adapt to the new UK environment, resulting in different residential outcomes to be achieved. Importantly, this adverse effect could become more severe as age at arrival gets older, which may create a direct link between our instrument and residential outcomes. To account for this possibility, we use, as a summary measure of cultural distance, genetic distance between the origin-country population and UK population, obtained from Spolaore and Wacziarg (2016). They argue that genetic distance, measuring the closeness of populations in terms of genes, reflects time since the pop-

**Table 4**

IV effects of English proficiency using additional controls and alternative samples.

	Base (1)	Control for cultural distance (2)	Control for religious distance (3)	Drop Europe (4)	Drop Common -wealth (5)	Control for parental education (6)	Control for proficiency index (7)
<i>A. Enclave residency</i>							
Language enclave	-3.574* (2.016)	-3.045 (1.938)	-3.649 (2.257)	-4.878* (2.601)	-4.398** (2.067)	-1.146 (3.352)	-1.041 (1.714)
Country-of-birth enclave	0.807 (1.908)	1.729 (2.011)	0.218 (2.465)	0.909 (2.715)	-1.048 (1.899)	0.159 (3.836)	3.758* (1.960)
Region-of-birth enclave	2.762*** (0.943)	2.577*** (0.926)	2.461** (1.001)	2.160 (1.537)	2.002** (0.791)	2.658** (1.270)	3.171*** (1.037)
Ethnic enclave	2.987*** (1.010)	2.415*** (0.848)	2.298*** (0.867)	1.196 (1.187)	2.454*** (0.818)	3.058** (1.346)	2.861*** (1.055)
<i>B. Neighbourhood quality index (in quintiles, 5 = worst)</i>							
Income deprivation	-0.914* (0.538)	-1.073* (0.575)	-1.253** (0.543)	-0.689 (0.786)	-0.912* (0.499)	0.501 (0.680)	-0.795 (0.510)
Employment deprivation	-1.376* (0.717)	-1.415* (0.729)	-1.621** (0.732)	-0.363 (0.675)	-1.343** (0.645)	-0.411 (0.987)	-1.247* (0.673)
Health deprivation	-1.502** (0.596)	-1.499** (0.610)	-1.977*** (0.645)	0.111 (0.790)	-1.430*** (0.530)	0.311 (0.788)	-1.553** (0.610)

Notes: Standard errors are clustered by country of birth. The estimates shown are the IV estimates of  $\beta_1$  in equation (4). Every column controls for the variables specified in Table 2. Additionally, columns 2, 3, 6 and 7 control for an additional variable each, which is an interaction of age at arrival with cultural distance from the UK (column 2), with religious distance from the UK (column 3), parental education (column 6), and English proficiency index (column 7), respectively. Columns 4 and 5 correspond to different sample specifications: sample excluding European migrants (column 4); and excluding Commonwealth migrants (column 5).

ulations shared the same ancestors. Over time, the ancestors transmit to their descendants not only their biological traits (i.e., genes) but also their cultural traits, such as habits and values, and this transmission occurs with variation. Populations that are genetically far from each other had more time to diverge in terms of cultural traits, and this divergence can in turn create barriers to human interactions. Spolaore and Wacziarg (2016) document that genetic distance of populations is significantly positively correlated with a wide array of measures of cultural differences. Column (2) of Table 4 indicates that, qualitatively, the results are not sensitive to the inclusion of this additional control. Note that column (1) reports the base results from Table 2 for comparison.

In a similar vein, we add a further control, which is a measure of religious distance between the UK and origin country (interacted with age at arrival). It is possible that linguistic distance is correlated with religious distance from the UK. The differences in norms associated with different religions, such as those related to educational attainment and marriage, could affect residential outcomes of migrants. Importantly, this effects of religious distance may vary by age at arrival, in which case our instrument will capture the compound effects of English skills and religious heterogeneity (or its associated heterogeneity in norms). Column (3) indicates that the results are robust to the inclusion of this additional variable.

We now take a different approach to address the concern that the main results are driven by different background characteristics associated with linguistic distance (that vary by age at arrival). Namely, we restrict our sample to a set of countries that may be less heterogeneous from each other than the original sample. To this end, we drop countries that have special ties with the UK. Unavoidably, we must pay the cost of losing observations, although the more we restrict our sample, the less heterogeneous the sample becomes. First, migrants from Europe might find it easier to adapt to the UK environment because European countries share commonality with the UK in culture and institutions, due to a long history of interactions across European countries. Similarly, migrants from Commonwealth countries might find it easier to adapt to the UK because of, for example, a similarity in their legal systems. The special ties these countries have with the UK might affect assimilation process, which could subsequently affect their residential outcomes. The results that exclude European and Commonwealth migrants are reported in columns (4) and (5) of Table 4, respectively. The results are broadly similar to our main results when Commonwealth countries are omitted. Interestingly, when migrants from European countries are omitted, the estimates of English proficiency become insignificant for ethnic-enclave residency and for the neighbourhood quality indexes. Considering the fact that non-European migrants in our dataset are generally from low-income countries, it might be that better English skills do not necessarily make them live in a neighbourhood of better quality, due to their better acceptance towards living in a relatively lower quality neighbourhood within the

UK, which could still be better conditions than those in their origin countries. As a result, these migrants might be more willing to live in a lower quality neighbourhood, and this tendency may magnify as age at arrival increases, because late arrivers are likely to be more affected by their origin-country standard of living. Non-European migrants might also enjoy relatively inexpensive rents in those neighbourhoods. Another possible explanation is that non-European migrants might be more financially constrained, such that whether they speak English very well or not has less of an impact on their residential outcomes.

We now consider another important factor, which can give an alternative explanation to our findings in the previous section, parental background. Precisely, parental characteristics of migrants from the two sets of countries might be different, and parents with different characteristics might have made different decisions regarding the timing of migration to the UK. For example, parents from linguistically distant countries might have recognised a possible barrier that their children would face if they migrate when their children are older, and may have chosen to migrate when their children were younger. At the same time, these parents might be different from the parents of migrants from linguistically close countries, in a way that can affect the future residential outcomes of childhood migrants. If this is the case, our IV estimates may reflect not only the effects of English skills but also the effects of different parental characteristics. To address this type of concerns, we control for parental education, measured by the dummy variable that takes the value of one if any of the two parents of the migrants has college education or above, and zero otherwise.<sup>16</sup> A limitation of this exercise is that, due to missing information on parental education, sample sizes decrease by roughly 40 per cent. Despite this limitation, we control for this possibly important confounding factor in column (6) of Table 4. We are not aware of any other studies on the causal effects on residential outcomes that explicitly account for parental education, which is a potentially very important confounding factor. The results are broadly similar to our main results for enclave residency in panel A, but for neighbourhood quality in panel B, the results become insignificant. To investigate whether this change in estimation results are driven by a change in sample sizes or the inclusion of parental education, we estimate the model with the smaller samples used in column (6) *without* controlling for parental education. Results (not reported) are very similar to those in column (6), implying that differences in the results are likely driven by a change in sample size.

Finally, it is possible that there is variation in the commonality of English within non-Anglophone countries due to, for example, the extent and quality of English language education. If that is the case, migrants from a linguistically distant country where English plays an important role may find it easy to acquire English skills, irrespective of linguistic distance

---

<sup>16</sup>Ideally, we would have liked to control for more detailed measures of parental education, but the indicator for college education or above is the only measure that we can construct from our dataset.

between their mother tongue and English, potentially biasing our IV estimates. To tackle this concern, we control for English Proficiency Index 2018 as the measure of average English skills in source countries.<sup>17</sup> A limitation of this exercise is that, because the English Proficiency Index is not available for all non-Anglophone countries in our sample, the sample sizes of non-Anglophone migrants are reduced by roughly 35 per cent. Nevertheless, we conduct this exercise using the best available data to explore this potential issue which can bias our estimates. The results reported in column (7) indicate that the results are broadly similar to the main results.

## 7. Conclusion

Inflows of migrants have increased in the countries in the Organisation for Economic Co-operation and Development (OECD) and European Union (EU) over the past two decades (OECD/EU, 2015), and the social integration of migrants is becoming an increasingly important policy objective. It is crucial to understand the key factors that influence migrant integration, and we focus on one of these possible factors: language skills. To date, the impact of language skills on earnings and employment of migrants has been intensively studied, but there is limited knowledge on the causal impact on residential outcomes of migrants in a host country. In this paper, we construct the measures of four different types of residential segregation — language, country-of-birth, world-region-of-birth and ethnic segregation — and analyse the impact of English skills on these measures of residential segregation. It is important to distinguish different types of enclave because English skills can have different effects on residency in different types of enclave. If migrants live in an enclave simply for linguistic convenience, those who are proficient in English would have no incentive to live in a language enclave. On the other hand, if they live in an enclave for other reasons than linguistic convenience, such as better job opportunities, availability of cultural amenities, and protection from possible discrimination, those fluent in English may still choose to live in an ethnic enclave. We also bridge the literature on language proficiency of migrants and on neighbourhood quality. We do so by analysing the effect of English skills on the quality of the neighbourhood migrants live in, where the neighbourhood quality is measured along three dimensions, i.e., the extent of income, employment and health deprivation of residents.

Our analysis is conducted using a unique dataset from the 2011 Census of England and Wales that we link to the measures of neighbourhood quality in England, allowing us to gain insight into the residential environments which migrants with different English skills live in.

---

<sup>17</sup>The English Proficiency Index 2018 is calculated by the international education company, EF Education First, based on the data on the scores of their English test collected from 1.3 million test takers across the world in 2017.

To overcome a possible endogeneity issue of English skills, we rely on an IV strategy, where linguistic distance between English and the origin-country language is used to construct an instrument. The use of linguistic distance is based on the idea that it is plausible that migrants whose mother tongue is linguistically closer to English (e.g., Dutch) find it easier to learn English, and vice versa. Moreover, the critical period hypothesis of language acquisition by Lenneberg (1967) documents that people exposed to a new language within the critical period of language acquisition (i.e., childhood) learn it more easily relative to those exposed outside of this critical period. This hypothesis implies that migrants born in a linguistically distant country who arrived in the UK at a younger age would on average have better English language skills than late arrivers. We incorporate these ideas into our analysis, and construct our instrument for English proficiency by interacting age at arrival with linguistic distance between English and the origin-country language. A methodological contribution we make here is that, to account for different assimilation process between Anglophone and non-Anglophone migrants, the interaction of age at arrival with coming from a non-Anglophone country, which is used as an instrument for English skills in Bleakley and Chin (2004; 2010) and other papers that follow their identification strategy (e.g., Akbulut-Yuksel et al., 2011; Miranda and Zhu, 2013; Guven and Islam, 2015; Yao and van Ours, 2015). We also account for year-of-arrival fixed effects which are important controls that have not been accounted for in Bleakley and Chin (2004) and other studies with the same identification strategy.

Our results suggest that better English skills lead migrants to live in an area with a lower concentration of individuals who speak their own native language. In contrast, we have found *positive* effects of English skills on the residency in region-of-birth and ethnic enclaves. This last effect is in contrast to the majority of the findings of previous correlation studies, showing the negative associations between better language skills and residence in ethnic enclaves. An exception is Bleakley and Chin (2010), who find, using US data, insignificant causal relations between English proficiency and residency in an ethnic enclave. Our findings imply that, for example, a Spanish-speaking white immigrant born in Spain who speaks English very well tend to live in an area with a low concentration of Spanish speakers, but in an area with a high concentration of individuals from Europe (their world region of birth) or of other white people (their ethnicity). Regarding a country-of-birth enclave, we find an insignificant effect, as is the case for Bleakley and Chin (2010). Our analysis has found heterogeneous effects of English skills on different types of enclave, suggesting that different mechanisms are at play depending on the types of enclave.

Turning to neighbourhood quality outcomes, we find that poorer English skills lead migrants to live in the neighbourhood of relatively poor quality, measured in terms of income, employment and health of residents. Our sub-analysis that explores the role of education suggests

that a lower educational attainment as a result of poor English skills is likely a key mechanism through which the poor language skills lead migrants to live in the neighbourhood of low quality. This is particularly the case for a lower neighbourhood quality captured by low income and employment deprivation of residents. In contrast, the sub-analysis on residency in different types of enclave indicates that education is not a key channel through which language skills impact residential segregation.

The results based on our IV strategy suggest that the impact of English skills on residential segregation varies depending on the type of enclave. Helping migrants improve their English skills, via for example providing English language courses, could be effective in reducing residential segregation by promoting migrants to live in less segregated areas with lower concentrations of people speaking their own native language. However, better language proficiency has been found to make migrants cluster in areas with higher concentrations of individuals from same world region of birth and ethnicity, suggesting that it is not only linguistic convenience, but also other aspects, such as the availability of cultural amenities or good job opportunities, are likely to play an important role in determining migrant residential location outcomes. To gain a further understanding as to why migrants, when they speak English well and presumably have a wide option of residential locations, cluster in ethnic enclaves and world-region-of-birth enclaves is an interesting and important future avenue for research.

Another important finding from our IV analysis is that poor English skills lead migrants to live in a lower quality neighbourhood in which residents are deprived of employment, among others. This could imply that poor language skills can have amplifying negative effects on migrant's labour market outcomes through neighbourhood effects (e.g., Topa, 2001), in addition to through its direct effects as found in a number of existing literature (e.g., Gonzalez, 2005; Clausen et al., 2009; Miranda and Zhu, 2013; Yao and van Ours, 2015). Our supplementary analysis finds that lower educational attainment as a result of poor English proficiency is likely the key channel through which language skills affect neighbourhood quality outcomes. In this regard, helping migrants improve their English skills could be an effective policy to improve their residential environments, among others, by improving their educational outcomes (Aoki and Santiago, 2018). When considering a design of language support, it would be beneficial to target those who arrived in the UK after age eight, because individuals who arrived in the UK before age eight appear to catch up with the level of proficiency of Anglophone migrants by the time they become adults anyways. It is also likely to be an efficient use of resources to target younger migrants, among those who arrived after age eight, because the earlier migrants are exposed to English, the easier it is for them to learn the language.

## Appendix

### Immigrants by country of birth.

<i>A. Anglophone countries</i>			<i>B. Non-Anglophone countries</i>		
<i>A1. Arrived aged 0 - 8</i>	<i>N</i>	<i>%</i>	<i>B1. Arrived aged 0 - 8</i>	<i>N</i>	<i>%</i>
Ireland	555	17.6	Cyprus	238	11.9
Kenya	346	11.0	Somalia	122	6.1
United States	239	7.6	Italy	121	6.0
South Africa	233	7.4	Turkey	110	5.5
Canada	223	7.1	France	79	3.9
Australia	204	6.5	Malaysia	77	3.8
Singapore	186	5.9	Germany	61	3.0
Jamaica	165	5.2	Iran	58	2.9
Malta	150	4.7	Egypt	55	2.7
Uganda	149	4.7	Iraq	52	2.6
Nigeria	117	3.7	Netherlands	51	2.5
Zambia	76	2.4	Vietnam	48	2.4
Zimbabwe	67	2.1	Spain	44	2.2
New Zealand	60	1.9	Portugal	43	2.1
Gibraltar	51	1.6	Belgium	41	2.0
Ghana	47	1.5	Yemen	40	2.0
Guyana	36	1.1	Saudi Arabia	39	1.9
Isle of Man	30	0.9	Malawi	33	1.6
Trinidad and Tobago	27	0.9	Aafghanistan	32	1.6
Mauritius	25	0.8	Libya	32	1.6
Total top 20	2,986	94.6	Total top 20	1,376	68.6
<i>A2. Arrived aged 9 - 16</i>	<i>N</i>	<i>%</i>	<i>B2. Arrived aged 9 - 16</i>	<i>N</i>	<i>%</i>
Ireland	536	19.0	Somalia	359	14.1
Kenya	496	17.6	Turkey	195	7.7
Jamaica	347	12.3	Cyprus	175	6.9
Uganda	217	7.7	Afghanistan	122	4.8
Nigeria	193	6.9	Poland	119	4.7
South Africa	156	5.5	Vietnam	90	3.5
Zimbabwe	119	4.2	China	89	3.5
Ghana	107	3.8	Portugal	77	3.0
United States	66	2.3	Iraq	71	2.8
Guyana	57	2.0	Yemen	67	2.6
Canada	44	1.6	Italy	66	2.6
Australia	41	1.5	Iran	61	2.4
Singapore	39	1.4	Kosovo	58	2.3
Sierra Leone	39	1.4	Germany	55	2.2
Zambia	31	1.1	France	47	1.8
Trinidad and Tobago	29	1.0	Malaysia	44	1.7
St Lucia	26	0.9	Malawi	39	1.5
Mauritius	24	0.9	Ethiopia	28	1.1
New Zealand	23	0.8	Rrussia	28	1.1
Barbados	22	0.8	Lithuania	28	1.1
Total top 20	2,612	92.8	Congo (Democratic Republic)	28	1.1
			Total top 20	1,846	72.5

Notes: Panels A and B present Anglophone and non-Anglophone countries, respectively. *N* refers to the number of individuals by country of birth for the top 20 countries present in our sample for those who arrived in the UK between age 0 and 8 (upper panels) and between 9 and 16 (lower panels).

Source: Authors' calculations based on the ONS Longitudinal Study.

## References

- Akbulut-Yuksel, M., Bleakley, H., Chin, A., 2011. The effects of english proficiency among childhood immigrants: Are hispanics different?, in: Leal, D.L., Trejo, S.J. (Eds.), *Latinos and the Economy*. Springer New York. Immigrants and Minorities, Politics and Policy, pp. 255–283.
- Aoki, Y., Santiago, L., 2018. Speak better, do better? Education and health of migrants in the UK. *Labour Economics* 52, 1 – 17.
- Bailey, M., Farrell, P., Kuchler, T., Stroebel, J., 2020. Social connectedness in urban areas. *Journal of Urban Economics* 118, 103264.
- Bauer, A.M., Chen, C.N., Alegria, M., 2012. Prevalence of physical symptoms and their association with raceethnicity and acculturation in the United States. *General Hospital Psychiatry* 34, 323 – 331.
- Bauer, T., Epstein, G., Gang, I., 2005. Enclaves, language, and the location choice of migrants. *Journal of Population Economics* 18, 649–662.
- Beckhusen, J., Florax, R., Graaff, T., Poot, J., Waldorf, B., 2013. Living and working in ethnic enclaves: Language proficiency of immigrants in U.S. metropolitan areas. *Papers in Regional Science* 92, 305–328.
- Bertrand, M., Luttmer, E.F.P., Mullainathan, S., 2000. Network effects and welfare cultures. *The Quarterly Journal of Economics* 115, 1019–1055.
- Bleakley, H., Chin, A., 2004. Language skills and earnings: Evidence from childhood immigrants. *The Review of Economics and Statistics* 86, 481–496.
- Bleakley, H., Chin, A., 2010. Age at arrival, English proficiency, and social assimilation among US immigrants. *American Economic Journal: Applied Economics* 2, 165–92.
- Borjas, G.J., 2000. Ethnic enclaves and assimilation. *Swedish Economic Policy Review* 7, 89–122.
- Borjas, G.J., 2015. The slowdown in the economic assimilation of immigrants: Aging and cohort effects revisited again. *Journal of Human Capital* 9, 483–517.

- Budria, S., Swedberg, P., 2015. The impact of language proficiency on immigrants' earnings in Spain. *Revista de Economía Aplicada* 23, 62–91.
- Chiswick, B.R., 1998. Hebrew language usage: Determinants and effects on earnings among immigrants in Israel. *Journal of Population Economics* 11, 253–271.
- Chiswick, B.R., Miller, P.W., 1995. The endogeneity between language and earnings: International analyses. *Journal of Labor Economics* 13, 246–288.
- Chiswick, B.R., Miller, P.W., 2005. Do enclaves matter in immigrant adjustment? *City & Community* 4, 5–35.
- Clarke, A., Isphording, I.E., 2017. Language barriers and immigrant health. *Health Economics* 26, 765–778.
- Clausen, J., Heinesen, E., Hummelgaard, H., Husted, L., Rosholm, M., 2009. The effect of integration policies on the time until regular employment of newly arrived immigrants: Evidence from Denmark. *Labour Economics* 16, 409–417.
- Cutler, D.M., Glaeser, E.L., Vigdor, J.L., 2008. When are ghettos bad? Lessons from immigrant segregation in the United States. *Journal of Urban Economics* 63, 759 – 774.
- Damm, A.P., 2014. Neighborhood quality and labor market outcomes: Evidence from quasi-random neighborhood assignment of immigrants. *Journal of Urban Economics* 79, 139 – 166. *Spatial Dimensions of Labor Markets*.
- Danzer, A.M., Yaman, F., 2016. Ethnic concentration and language fluency of immigrants: Evidence from the guest-worker placement in Germany. *Journal of Economic Behavior & Organization* 131, 151 – 165.
- Di Paolo, A., Raymond, J.L., 2012. Language knowledge and earnings in Catalonia. *Journal of Applied Economics* 15, 89 – 118.
- Dustmann, C., 1994. Speaking fluency, writing fluency and earnings of migrants. *Journal of Population Economics* 7, 133–156.
- Dustmann, C., Fabbri, F., 2003. Language proficiency and labour market performance of immigrants in the UK. *Economic Journal* 113, 695–717.
- Dustmann, C., van Soest, A., 2001. Language fluency and earnings: Estimation with misclassified language indicators. *The Review of Economics and Statistics* 83, 663–674.

- Dustmann, C., van Soest, A., 2002. Language and the earnings of immigrants. *Industrial and Labor Relations Review* 55, 473–492.
- Edin, P.A., Fredriksson, P., Aslund, O., 2003. Ethnic enclaves and the economic success of immigrants: Evidence from a natural experiment. *The Quarterly Journal of Economics* 118, 329–357.
- Glick, J., White, M., 2003. Academic trajectories of immigrant youths: Analysis within and across cohorts. *Demography* 40, 759–783.
- Gonzalez, L., 2005. Nonparametric bounds on the returns to language skills. *Journal of Applied Econometrics* 20, 771–795.
- Güven, C., Islam, A., 2015. Age at Migration, Language Proficiency, and Socioeconomic Outcomes: Evidence From Australia. *Demography* 52, 513–542.
- Iceland, J., Scopilliti, M., 2008. Immigrant residential segregation in U.S. metropolitan areas, 1990-2000. *Demography* 45, 79–94.
- Isphording, I.E., Otten, S., 2014. Linguistic barriers in the destination language acquisition of immigrants. *Journal of Economic Behavior & Organization* 105, 30 – 50.
- Janssen, S., 2010. *The world almanac and book of facts 2011*. Infobase Learning, New York.
- Katz, L.F., Kling, J.R., Liebman, J.B., 2001. Moving to opportunity in Boston: Early results of a randomized mobility experiment. *The Quarterly Journal of Economics* 116, 607–654.
- Kimbro, R.T., Gorman, B.K., Schachter, A., 2012. Acculturation and self-rated health among Latino and Asian immigrants to the United States. *Social Problems* 59, pp. 341–363.
- Kling, J.R., Liebman, J.B., Katz, L.F., 2007. Experimental analysis of neighborhood effects. *Econometrica* 75, 83–119.
- Laliberte, J.W., 2019. Language skill acquisition in immigrant social networks: Evidence from Australia. *Labour Economics* 57, 35 – 45.
- Lazear, E.P., 1999. Culture and language. *Journal of Political Economy* 107, S95–S126.
- Lee, S., O'Neill, A., Ihara, E., Chae, D., 2013. Change in self-reported health status among immigrants in the United States: Associations with measures of acculturation. *PLOS ONE* 8, 76494.
- Lenneberg, E.H., 1967. *Biological foundations of language*. Wiley, New York.

- Logan, J.R., Zhang, W., Alba, R.D., 2002. Immigrant enclaves and ethnic communities in New York and Los Angeles. *American Sociological Review* 67, 299–322.
- Ludwig, J., Duncan, G.J., Gennetian, L.A., Katz, L.F., Kessler, R.C., Kling, J.R., Sanbonmatsu, L., 2012. Neighborhood effects on the long-term well-being of low-income adults. *Science* 337, 1505–1510.
- Massey, D.S., 1985. Ethnic residential segregation: A theoretical synthesis and empirical review. *Sociology and Social Research* 69, 315–350.
- Miller, P.W., Neo, L.M., 1997. Immigrant unemployment: The Australian experience. *International Migration* 35, 155–185.
- Ministry of Housing, Communities and Local Government, 2015. English indices of deprivation. URL: <https://www.gov.uk/government/collections/english-indices-of-deprivation>. accessed on 4 July 2016.
- Miranda, A., Zhu, Y., 2013. English deficiency and the native-immigrant wage gap. *Economics Letters* 118, 38 – 41.
- Miranda, P.Y., Gonzalez, H.M., Tarraf, W., 2011. Pathways between acculturation and health: Does the measure matter? *Hispanic Journal of Behavioral Sciences* 33, 524–539.
- OECD\European Union, 2015. Indicators of immigrant integration 2015: Settling in. OECD Publishing, Paris.
- Oreopoulos, P., 2003. The long-run consequences of living in a poor neighborhood. *The Quarterly Journal of Economics* 118, 1533–1575.
- Sanbonmatsu, L., Katz, L., Ludwig, J., Gennetian, L., Duncan, G., Kessler, R., Adam, E., McDade, T., Lindau, S., 2011. Moving to opportunity for fair housing demonstration program: Final impacts evaluation. US Department of Housing and Urban Development.
- Shields, M.A., Price, S.W., 2002. The English language fluency and occupational success of ethnic minority immigrant men living in English metropolitan areas. *Journal of Population Economics* 15, 137–160.
- Spolaore, E., Wacziarg, R., 2016. Ancestry and development: New evidence. Discussion Papers Series, Department of Economics, Tufts University 0820. Department of Economics, Tufts University.

- Stock, J.H., Wright, J.H., Yogo, M., 2002. A survey of weak instruments and weak identification in generalized method of moments. *Journal of Business & Economic Statistics* 20, 518–29.
- Topa, G., 2001. Social interactions, local spillovers and unemployment. *The Review of Economic Studies* 68, 261–295.
- United Nations, 2019. International migrant stock 2019. URL: <https://www.un.org/en/development/desa/population/migration/data/estimates2/estimates19.asp>. Accessed on 29 February 2020.
- Weinberg, B.A., Reagan, P.B., Yankow, J.J., 2004. Do neighborhoods affect hours worked? Evidence from longitudinal data. *Journal of Labor Economics* 22, 891–924.
- Weinhardt, F., 2014. Social housing, neighborhood quality and student performance. *Journal of Urban Economics* 82, 12–31.
- Yao, Y., van Ours, J.C., 2015. Language skills and labor market performance of immigrants in the Netherlands. *Labour Economics* 34, 76 – 85.



## Online Appendix

Alternative measure of English skills.

Dependent variable:	Enclave, neighbourhood quality			English ability ordinal measure
	OLS	IV	Reduced-form	First-stage
	(1)	(2)	(3)	(4)
<i>A. Enclave residency</i>				
Language enclave	-2.953*** (0.682)	-2.519* (1.457)	0.528* (0.290)	-0.209*** (0.062)
Country-of-birth enclave	-0.944*** (0.333)	0.569 (1.349)	-0.119 (0.290)	-0.209*** (0.062)
World-region-of-birth enclave	-0.219* (0.111)	1.947*** (0.687)	-0.408*** (0.129)	-0.209*** (0.062)
Ethnic enclave	-0.235* (0.124)	2.112*** (0.736)	-0.444*** (0.087)	-0.210*** (0.061)
<i>B. Neighbourhood quality index (in quintiles, 5 = worst)</i>				
Income deprivation	-0.300*** (0.056)	-0.638 (0.387)	0.134* (0.075)	-0.211*** (0.063)
Employment deprivation	-0.273*** (0.056)	-0.960* (0.520)	0.202** (0.093)	-0.211*** (0.063)
Health deprivation	-0.180*** (0.051)	-1.048** (0.429)	0.221*** (0.076)	-0.211*** (0.063)

Notes: Standard errors are clustered by country of birth. OLS and IV are the estimates of  $\beta_1$  in equation (4), where the ordinal measure for English skills, ranging between 0 and 3, is used as the measure of English proficiency. First-stage and reduced-form are the estimates of the coefficients on the instrument, which is an interaction of linguistic distance between a home-country language and English, with age at arrival. Refer to Table 2 for controls included. Sample sizes for panels A and B are 10,522 and 10,320, respectively, except for ethnic enclave (10,482 observations due to 40 missing values).

Source: Authors' calculations based on the ONS Longitudinal Study.