

# Speaking English in a Globalizing World: Information Technology and Education in India

Gauri Kartini Shastry\*

Harvard University

PRELIMINARY AND INCOMPLETE

March 2007

## Abstract

I study how the impact of globalization on returns to education and school enrollment varies with the elasticity of the skilled labor supply. I exploit variation in the cost of learning English across districts in India, driven by linguistic diversity that made it necessary for individuals to learn additional languages. In India, the two common choices for a second language are English and Hindi, the native lingua franca. Individuals whose native language is linguistically further from Hindi have lower relative opportunity costs of learning English, mainly because they find Hindi harder to learn but also because they often suffer psychic costs when using Hindi, a language that many non-native Hindi speakers feel was imposed on them. I first show that linguistic distance from Hindi increases the probability of learning English, even in 1961. Using newly collected data on information technology (IT), I show that districts with lower costs of learning English experienced greater growth in IT after trade reforms in the early 1990s. In addition, these districts experienced greater growth in relative employment of educated workers but smaller growth in skilled wage premiums, due to the greater skilled labor supply elasticity. Finally, I show that these districts experienced greater increases in school enrollment.

---

\*Correspondence: shastry@fas.harvard.edu. I am grateful to David Cutler, Esther Duflo, Caroline Hoxby, Michael Kremer for their advice and support, David Clingingsmith for access to additional data and all participants of the Research in Labor Economics and Development Economics lunch workshops at Harvard for their comments. In addition, I thank Filipe Campante, Davin Chor, Quoc-Anh Do, Eyal Dvir and Michael Katz for useful conversations and assistance. Finally, I thank Daniel Tortorice for invaluable support and encouragement. All errors are mine.

# 1 Introduction

While most economists agree that free trade has significant benefits over autarky for all countries involved, i.e. that free trade increases the size of the pie, there is some debate over how the benefits of trade are distributed within a country. While several recent empirical studies have found that trade liberalization in Latin America has caused sizeable increases in inequality and skill wage premiums,<sup>1</sup> there is less evidence from Asian countries and it is more mixed.<sup>2</sup> In contrast, standard Heckscher-Ohlin trade theory unambiguously predicts that globalization should reduce inequality and skill premia. Under the simplest model with two goods, two countries and two factors (skilled and unskilled labor), the country abundant in unskilled labor (the poor country) should specialize in unskilled-labor-intensive industries after trade liberalization. This increases demand for unskilled labor and drives down the skilled wage premium. While there are numerous theoretical extensions to help reconcile the theory with these empirical findings,<sup>3</sup> an additional, under-emphasized, dimension is whether labor supply and education responds to the increased wage inequality. The fact that countries in Latin American countries experienced a much smaller increase in the supply of skilled workers relative to East-Asian economics may explain the mixed findings above.<sup>4</sup>

In this paper, I explore this dimension of how the effects of liberalization vary within a developing country and provide evidence on the labor supply response to globalization as well as how the effect on skill premiums varies with labor supply elasticity. I exploit variation in the labor supply elasticity of skilled workers due to historically-driven differences in policies regarding language of instruction. Many countries, particularly those with a colonial past, have struggled with the question of whether to encourage their people to retain diverse local

---

<sup>1</sup>See Goldberg and Pavcnik (2004) for a review of this literature. The literature includes, e.g., Hanson and Harrison (1999), Feenstra and Hanson (1997), Feliciano (1993) and Cragg and Epelbaum (1996) on Mexico, Robbins, Gonzales and Menendez (1995) on Argentina, Robbins (1995a) on Chile, Robbins (1996a) and Attanasio, Goldberg and Pavcnik (2004) on Colombia, Robbins and Gindling (1997) on Costa Rica and Robbins (1995b, 1996b) on Uruguay.

<sup>2</sup>See, e.g., Wood (1997) for a survey, Lindert and Williamson (2001) and Wei and Wu (2001).

<sup>3</sup>See, e.g. Feenstra and Hanson (1996, 1997), Kremer and Maskin (2006) and other extensions discussed in Goldberg and Pavcnik (2004).

<sup>4</sup>See Attanasio and Szekely (2000), Sanchez-Paramo and Schady (2003).

languages, choose a single native language or promote a global lingua franca, such as English. In particular, the choice of medium of instruction in public schools is one with far-reaching consequences.<sup>5</sup> On one hand, there are costs of promoting a non-native global language. Native language instruction can strengthen national identity, particularly important in young countries made up of numerous ethnic groups. In addition, instruction in a foreign language may impose costs on poor households if they find such education less accessible. On the other hand, there may be benefits of promoting a global lingua franca if coveted white-collar jobs in government or business use that language.<sup>6</sup> Instruction in a global language in public schools may increase economic opportunities for the poor.<sup>7</sup> Most importantly, promoting the learning of English may also allow more people to benefit from globalization and technological progress. The ability to integrate better with the world economy may bring more of the benefits of trade liberalization to places that promote English. Since much technological progress happens in English (for example, in information technology), the ability to speak English may facilitate the adoption of new technologies.

This paper examines the increased returns to an English education due to economic liberalization and technological progress in the 1990s and demonstrates the effect on educational attainment in India, by exploiting exogenous variation in the cost of learning English across Indian districts. Using most measures of variation in the cost of learning English, such as the number of individuals who learn English, would be highly problematic because the cost of learning English is endogenous. State or local governments that care more about the benefits of promoting a global lingua franca such as access to global opportunities may promote the teaching of English, but also pursue other policies that increase trade-related jobs. I exploit variation in costs of learning English that is driven by linguistic and historical forces that are exogenous to such outward-oriented or forward-looking policies and also to reverse causality. In fact, this variation in the cost of acquiring English caused people to learn English

---

<sup>5</sup>See Human Development Report (2004) and Angrist, Chin and Godoy (2006), e.g.

<sup>6</sup>See Lang and Siniver (2006).

<sup>7</sup>See Angrist and Lavy (1997) and Munshi and Rosenzweig (2006).

for non-trade related reasons even in 1961, long before trade liberalization could have been contemplated. The variation in costs is driven by historical linguistic diversity that made it necessary to learn a second language even to communicate with others in the same district. The common choices for a second language in India are English and Hindi, the native lingua franca. Individuals whose mother tongue is linguistically further from Hindi have a lower opportunity cost of learning English because they find Hindi more difficult to learn, but also because they are more likely to suffer psychic costs when speaking Hindi, a language that many non-native Hindi speakers feel was imposed on them as a national language. Over time, these historical tendencies led to the growth of institutions that promote the learning of English in districts where native languages are further from Hindi. I first show that this relationship holds; linguistic distance from Hindi increases the number of native speakers who learn English and predicts the percent of schools that teach English.

Next, I demonstrate how the impact of globalization during the 1990s has varied by pre-existing differences in linguistic distance from Hindi. I examine data that I gathered and coded on the information technology (IT) sector in India, an industry that grew primarily due to economic liberalization and technological progress in the 1990s and hires educated, English-speaking workers. Information technology includes both software firms and business-process outsourcing such as call centers and data entry firms. I show that IT firms were more likely to locate in districts with lower costs of learning English. I also find that these districts have greater IT employment. I then posit that in the new more open Indian economy, there is a greater payoff to being educated and English speaking. I provide evidence using micro-level data that, in the 1990s, districts with lower costs of learning English experienced a greater increase in employment for educated workers but a smaller increase in the average skilled wage premium. A simple theoretical model provides intuition for these results. Separating out workers by industry, I show that wage premiums in certain industries (financing, insurance, computer related activities, research and development, other business activities) rose faster in these districts as well. As suggestive evidence for the trade in services channel,

I do not find a corresponding rise in wage premium for education workers in other industries (other services, such as public service, education, health; manufacturing; agriculture; hotel and restaurants; wholesale and retail; transportation services; or communications, such as post, courier and telecommunications). However, I do find an increase in all wages in the transportation and communications industries lending more credibility to trade-related growth as an important mechanism behind these differential trends. I show how pre-existing differences in linguistic diversity explain changes in employment and wages between 1987 and 1999; thus, I estimate not just correlations but the differential impact of these lower costs of English during a period of liberalization.

Finally, these pre-existing differences in linguistic distance to Hindi also explain differential changes in school enrollment trends during the 1990s relative to pre-existing trends. I demonstrate that districts where the average person speaks a language that is linguistically further from Hindi experienced greater increases in urban school enrollment from 1993 to 2002 even relative to pre-existing trends in school enrollment.

Thus, the contributions of this paper are two-fold, corresponding to the two motivating themes described above. The most conservative interpretation centers on the benefits of language policies that advance the study of a global language in a developing country. In particular, districts that had more English speakers for these reasons found themselves in a better position to take advantage of the opportunities from trade after liberalization. In addition, the paper contributes to the literature on trade liberalization because it confirms a possible explanation for mixed evidence on the impact of globalization on wage inequality and finds evidence of longer term consequences to this rise in wage inequality. While trade liberalization may increase wage inequality in the short run, this effect could be dampened in the longer term as factor supply responds.

Besides the two strains of trade literature and the education literature described above, this paper is related to Munshi and Rosenzweig (2006). Using a household survey from a suburb of Mumbai, India, the authors show that increases in the returns to English dwarfed

increases in the returns to education and that enrollment rates in English-medium schools rose in the 1980s and 1990s. My paper differs in a number of ways. First, I show that educational attainment overall rises, not just in English instruction schools. Second, I use data from all over India exploiting exogenous variation in the cost of learning English. Lastly, I go one step further to explore the trade-related mechanism through which the returns to English have risen. Edmonds, Pavcnik and Topalova (2005) also study the relationship between economic liberalization and educational attainment in India and find an adverse impact on schooling. While the authors isolate the effect on school enrollment of reduced family income due to import competition, I focus on the impact of increases in job opportunities and returns to English education from exports and integration with world markets.

The paper is organized as follows. Section 2 provides background information on trade liberalization and information technology in India. Section 3 describes the linguistic diversity, costs of learning English, and medium of instruction. I show that linguistic distance from Hindi increases the tendency to learn English, but does not predict other economic measures prior to economic liberalization. Section 4 describes a simple theoretical model to provide intuition for the empirical findings. Section 5 discusses the empirical methodology, while section 6 describes the data. Section 7 examines IT firm location and employment decisions and provides evidence on employment of educated workers and returns to education. Finally, in section 8, I show that districts where native languages are farther from Hindi experienced greater increases in school enrollment. Section 9 concludes.

## **2 Background on trade liberalization and IT**

Throughout much of the post-colonial period, India heavily protected its economy. While some small steps towards integrating with world markets were taken in the late 1970s and 1980s, even as late as 1990, tariff and non-tariff barriers posed significant obstacles for trade. The average tariff was 79% and sixty-five percent of all imports were subject to non-tariff

barriers (Panagariya 2003). A balance-of-payments crisis due to extensive borrowing in 1991 resulted in a shift towards policies favoring a more open economy. Reforms ended most import licensing requirements for capital goods and reduced tariff rates substantially, although mostly for non-agricultural goods. Service sectors which had previously been heavily regulated by the government saw significant changes. The 1994 National Telecommunications Policy and 1999 New Telecom Policy opened cellular and other telephone services to both private and foreign investors. Foreign direct investment (FDI) in e-commerce was free of all restrictions and foreign equity in software and electronics was granted automatic approval, particularly for IT firms set up exclusively to export (Panagariya 2004). This service sector liberalization, along with technological progress, led to the remarkable growth in the outsourcing of services in the information technology sector, to India.

By 2004, India was the single largest destination for foreign companies to purchase IT services, contributing about two-thirds of global software outsourcing and half of business process outsourcing. In 2005, IT outsourcing accounted for 5% of India's GDP and was forecasted to contribute 17% to India's projected growth to 2010 (The Economist 2006). Employment growth has also been strong over the past decade; from 56,000 professionals in 1990-91, the sector employed 813,500 in 2003, implying an annual growth rate of more than twenty percent (NASSCOM 2004). In particular, the IT sector increased job opportunities for young, educated workers; the median age of IT professionals is 27.5 years and 81% of them have at least a bachelor's degree (NASSCOM 2004). An entry-level job in a call sector can earn on average Rs. 10,000 (\$230) considered very high for a first job (The Economist 2005). In addition, the excitement regarding the growth of the IT sector is palpable. IT firms advertise heavily in newspapers and on job search websites.

Their young age, export focus and reliance on foreign capital make IT firms relatively free to locate based on other inputs. One of the principle factors in the location decision of IT firms is manpower, i.e. the availability of an educated, English-speaking population. I show below that IT firms choose to locate in places with lower costs of learning English.

### 3 Linguistic distance from Hindi and identification

The 1961 Census of India documented 1652 mother tongues spoken in India from five distinct language families native to India. These language families are quite diverse; while linguists assert familial relationships between languages as far apart as English and Hindi (both are Indo-European), they are unwilling to connect many languages native to India such as Hindi and Kannada, the language spoken in Bangalore. Figures 1 and 2 present maps of India with the density of native speakers of 114 languages. This linguistic diversity has had implications for bilingualism (Clingsmith 2006); most individuals, especially urban, educated individuals, need to learn a second language to communicate at a local level. According to the 1991 census, 19.4% of Indians are multilingual. The two most common languages learned are Hindi, the native lingua franca, and English, due to the British colonial history. As of the 1991 census, sixty percent of all multilingual people not native in Hindi learned Hindi as a second language. For English, this fraction was only slightly smaller at 56%.<sup>8</sup> The next most popular second language, Kannada, was learned by only 6% of the multilingual population. In fact, 83% of all multilinguals speak either Hindi or English.

An individual with a more obscure mother tongue has to choose between Hindi and English. Mechanically, an individual whose mother tongue is linguistically close to Hindi will find it easier to learn Hindi relative to someone whose mother tongue is farther from Hindi. The history of language in India amplifies this tendency because of the controversial decision to make Hindi the national language. During the British occupation, English was established as the language of government, the medium of instruction and the language of the elite. After India became independent in 1947, a nationalist movement to make an indigenous Indian language the official language favored Hindi, since it was spoken by more people than any other native language. This movement was opposed by non-native Hindi speakers, but after much debate, Hindi was written into the constitution as the language

---

<sup>8</sup>Of course since Hindi was spoken by more than three hundred million people as a first language, while English was spoken by only 180,000 as a first language, many more people spoke Hindi than English.



of administration, meant to replace English within 15 years. This led to riots in non-Hindi speaking areas, the most violent of which occurred in Tamil Nadu in May 1963. Speakers of other languages felt at a disadvantage speaking Hindi and finally, in 1967, the government passed a law making Hindi and English joint official languages (Hohenthal 2003).

This background explains why English is more prevalent among people who speak languages distant to Hindi. In fact, in some states, more people speak English than Hindi. Over time, as I show below, the relationship between linguistic distance and English prevalence became institutionalized through English education. This theory has ambiguous predictions for whether native Hindi speakers learn English. On one hand, they do not need a second language to communicate within India; if they choose to learn a second language, it could be for other reasons. On the other hand, if they choose to learn a second language to communicate within India, English would allow them to interact with more additional people than any other language.

### **3.1 Medium of instruction in Indian schools**

When the British began to colonize India, they did not plan to provide mass education. They set up schools and colleges in large cities that taught entirely in English, meant to foster an elite class to help govern the country (Nurullah and Naik 1947, Kamat 1985). By 1850, other institutions such as missionary societies and princely states had set up rural schools that taught in native languages. Finally, in 1854, the recommendations set forth in Sir Charles Wood's Despatch marked the British government's committed to educating the entire population (Dakin, Tiffen and Widdowson 1968). Education spread to lower classes and an increasing number of schools taught in native languages. University education in major cities, however, was still primarily in English. Even today, it is the main medium of post-tertiary instruction (Hohenthal 2003). In 1993, according to the Sixth All India Educational Survey, there were over 28 different media of instruction in primary schools (regardless of government or private funding) across urban areas in India. While Hindi is

the most common medium of instruction with 38% of primary schools, English is second with 9%. At the secondary school level, there are 37 languages taught as a first or second language across urban areas of India.

### 3.2 Measuring linguistic distance

The 1961 and 1991 Census of India provide data on the number of people in each state that speak each of 114 distinct mother tongues and how many of them learn each of these languages as a second language. In addition, at the district level, we know how many people speak each language as a mother tongue. I first calculate various measures of the distance from Hindi of a language. In order to obtain a measure at the district level, I calculate the weighted average of the distance from Hindi of all native languages spoken in a district where the weights are the district population share. For an alternate measure, I calculate the percent of speakers in a district who speak languages sufficiently far from Hindi.

As there is no universally accepted measure of language distance among linguists, I calculate three independent measures. My preferred measure was developed in consultation with an expert on Indo-European languages, Jay Jasanoff, the Diebold Professor of Indo-European Linguistics and Philology at Harvard University. This measure is based on drawing seven concentric circles of languages around Hindi as they get linguistically more different (see figure 3). I count the circles (and call them degrees) from Hindi to each language. Figure 4 provides a map of India demonstrating the distribution of the weighted average of this measure and figure 5 provides a map of the percent of people who speak languages at least 3 degrees away from Hindi. Note that much of the variation is across regions (indicated by thick black lines). However, since we might worry that this variation is correlated with other factors (e.g. geography, culture), I include region fixed effects and differential trends.

A second measure of linguistic distance is based on language family trees. The most widely used language trees are from the Ethnologue database, one of the most comprehensive listings of currently known languages. Many linguists rely on and contribute to the database.

Figure 6 provides an extract from the Ethnologue’s language tree that includes all languages found in the Indian census. I define the distance between two languages as the number of nodes between the languages. For example, Urdu is two degrees away from Hindi, while Marwari is four degrees away. In order to link the other language families with Hindi, I assume there is a node connecting the different language families.

Another measure is taken from a method called glottochronology, which is used to estimate the time of divergence between languages (Swadesh 1972). The method involves making a list of 210 core words, i.e. words that are the most resistant to change as languages evolve. Then, using expert judgments on whether these words across languages are cognates with each other, we can calculate the percent of words that are cognates between each pair of languages.<sup>9</sup> I use the percent of cognates shared between each language and Hindi from the Dyen et al. (1992) dataset of 95 Indo-European languages. Table 1 provides example core meanings in English, Bengali and Hindi as well as cognate judgments for each pair of words. Since the dataset does not provide the words or cognate judgments for non Indo-European languages, I assume these languages have only 5% of words in common with Hindi, since the lowest percent of cognates with Hindi among Indo-European languages is 14.6%.<sup>10</sup> Finally, for Indo-European languages in the 1991 Census of India language data which do not appear in Dyen et al.’s list, I use the percent of cognates with Hindi of the closest language in the tree that also appears in the list. Close matches exist for the 12 Indo-European languages that require this.

The correlation between these 3 measures is quite high. Across languages, the correlation between degrees and nodes between languages is 0.9283. The correlation between these two measures and the percent of cognates is -0.9358 and -0.9731 respectively. Panel A of table 2 provides summary statistics on these and other measures regarding language in India.

---

<sup>9</sup>The original version of this method also involved a formula that converted this percent of cognates into a time of divergence, which is currently out of favor among linguists. Nevertheless, the percent of cognates is still an acceptable measure of similarity between languages.

<sup>10</sup>This choice is not arbitrary - linguists use 5% as a significance level when determining whether two languages are related. If less than 5% of words are cognates, linguists assume that those that are represent noise and the languages are unrelated.

### 3.3 Identification

We cannot estimate the impact of the cost of learning English using most measures of these costs since they would be endogenous. Local governments can influence these costs based on their preferences. For example, if the government cares about access to global opportunities, it may both promote education in English and provide incentives for foreign direct investment. We would also worry about reverse causality, since these outsourcing firms often set up English training centers. The variation in the cost of learning English that comes from linguistic distance to Hindi, however, does not suffer from these problems. Linguistic distance to Hindi impacts the cost of learning English in a manner that is orthogonal to preferences of different local governments. In addition, government policies or English-language opportunities will not affect the linguistic distance of a language from Hindi. Large movements of people across district boundaries may influence the linguistic distance from Hindi of languages spoken in a district, but migration in India is still quite low. According to the 1987 National Sample Survey, only 12.3% of individuals in urban areas had migrated in the past five years, only 6.8% had moved from outside their current district in this time and only 2.4% had moved from outside their current state.

In this subsection, I demonstrate that linguistic distance from Hindi predicts measures of English prevalence, but is not strongly correlated with other measures of economic development before 1990. I first use data on second languages spoken by different ethnic groups within Indian states from the Census of India to show that linguistic distance from Hindi predicts the percent of multilinguals who choose to study English

$$E_{lk} = \alpha_0 + \beta' D_l + \alpha_1' X_{lk} + \theta_f + \gamma_g + \epsilon_{lk}$$

where  $E_{lfk}$  is the percent of native speakers of language  $l$  of language family  $f$  in state  $k$  in region  $g$  who choose to learn English,  $D_l$  is a vector of measures capturing the linguistic distance of language  $l$  from Hindi, and  $X_{lk}$  is a vector of control variables at the language

and state level. To control for other characteristics of ethnic groups or particular regions of India, I include language family and region fixed effects. The regions include north, northeast, east, south, west and central India. The vector  $X_{lk}$  includes indicator variables for whether language  $l$  is Hindi and for whether language  $l$  is the most spoken language in the state. Finally, I control for a quadratic polynomial in the share of speakers of language  $l$  who reside in state  $k$  and who reside anywhere in India. I weight observations by the number of native speakers in the state and cluster at the state level.

The results using the 1991 Census of India are presented in table 3. Column 1 assumes a linear relationship implying that one degree increases the percent of multilinguals who learn English by 7.7 percent. Column 2 estimates that a 10% point increase in how many people speak languages more than two degrees away from Hindi increases the percent of multilingual English speakers by 2.16 percentage points. Finally, column 3 demonstrates that while the relationship between linguistic distance and English prevalence does not appear to be linear, the assumption of monotonicity is not problematic and this is true even excluding language family fixed effects.<sup>11</sup> The omitted group in this regression is Urdu speakers; Urdu is very close to Hindi (a distance of 0), but considered a separate language. From the F-statistics shown at the bottom of the table, it is clear that linguistic distance from Hindi does predict variation in the proportion of native speakers who learn English. In addition, I reject the hypothesis that having any distance between your mother tongue and Hindi has the same effect on English learning by testing the equality of all linguistic distance fixed effects (but allowing them to be different from Urdu).<sup>12</sup>

I next explore the relationship between linguistic distance and the percent of all native speakers who are multilingual (see columns 4-6 in table 3). While the linear measure of

---

<sup>11</sup>Note that linguistic distance of 5 degrees is dropped when including fixed effects for each number of degrees away from Hindi since the only Indo-European, but not Indo-Aryan language in the data is English and I omit observations for English. The results are robust to modifying the measure of linguistic distance to omit this concentric circle (i.e. give all non-Indo-Aryan languages a distance of 5 degrees). In addition, the fixed effect for linguistic distance of 6 is also dropped since it consists of all non-Indo-European languages and I include language family fixed effects.

<sup>12</sup>These results are also robust to using the percent of all native speakers who choose to learn English.

linguistic distance predicts multilingualism, the fixed effects for unit distances specification demonstrates that the effect is clearly not rising in linguistic distance. In fact, the results seems to be driven by a greater percent of multilingual individuals at a linguistic distance of 2 units relative to Urdu speakers. I also show that linguistic distance negatively predicts the percent of all multilingual native speakers who choose to learn Hindi but not English (see columns 7-9), as my theory would predict.<sup>13</sup>

Table 4 presents similar results from regressions using data from the 1961 Census of India. These regressions differ slightly due to changes in the data collection, but the results are very similar. For example, since Hindi and Urdu are categorized together, the omitted group is languages one degree away from Hindi. The last columns also use a different dependent variable, i.e. the percent of multilingual individuals who learn Hindi (even if they also learn English). These results discredit an important concern with this identification strategy. The identification strategy would be invalid if, for example, the ethnic groups speaking languages further from Hindi were more forward-looking or outward-oriented in the 1980s and anticipated the trade benefits to learning English. However, the tendency for these ethnic groups to learn English was evident in 1961, much before anyone could anticipate the trade liberalization of the early 1990s and the enormous returns to speaking English. This supports the use of linguistic distance to highlight exogenous variation in the costs of learning English. In addition, the number of English bilinguals among ethnic groups in different states in 1961 is highly correlated with the number of English learners in 1991.

Using a similar specification, I also explore how linguistic distance from Hindi of languages spoken in a district predicts whether schools teach in English with the following regression

$$M_{ij} = \alpha_0 + \beta' D_j + \alpha_1 P_j + \alpha_2' Z_j + \gamma_i + \epsilon_{ij} \quad (1)$$

---

<sup>13</sup>All results in this table are robust to including state fixed effects instead of region fixed effects and running unweighted regressions including only languages spoken by at least 100 people in the state (the median number of people represented by an observation). The results are similar when I separate out men and women; in fact, linguistic distance has a slightly larger effect on how many women learn English. In addition, the results are similar if I cluster by native language instead of state, to account for correlated error terms between speakers of a particular language across India.

where  $M_{ij}$  is a measure of language instruction at school level  $i$  (primary, upper primary, secondary or higher secondary) in area  $j$  (either a state or a district, depending on the outcome variable),  $D_j$  is a vector of measures capturing the linguistic distance from Hindi of languages spoken,  $P_j$  is child population and  $Z_j$  is a vector of control variables. The vector  $Z_j$  includes average household wage income, average income for individuals who have completed secondary school, the percent of adults who have regular wage or salaried jobs, the distance to the closest of the 10 biggest cities, and the percent of households that have electricity, all measured in 1987. In addition, I include school level fixed effects and control for the percent of people who: have graduated from college, have completed secondary school, are literate, are Muslim, are native English speakers or regularly use a train. To ensure that these results are not driven by large Hindi-speaking populations in the "Hindi Belt" states who are particularly poor due in large part to corruption and government inefficiency, I control for the percent of people who are native Hindi speakers and include an indicator variable for the following states: Bihar, Uttar Pradesh (and Uttaranchal), Madhya Pradesh (and Chhattisgarh), Haryana, Punjab, Rajasthan, Himachal Pradesh, Jharkhand, Chandigarh and Delhi. These variables come from a number of sources as discussed in the data section below. Summary statistics on the outcome variables, from the Sixth All India Educational Survey, are provided in panels A and B of table 2.<sup>14</sup>

The results show that linguistic distance from Hindi predicts the percent of schools that teach in English or teach English as a second language at the state level (see columns 1-6 of table 5). In columns 3 and 6, I re-estimate equation (1) using the percent of native speakers at each distance from Hindi in the state. The p-value from the F-test at the bottom of the table indicates that linguistic distance does predict the teaching of English but the individual measures are not significant. An increase in one degree in the distance from Hindi of the average speaker in a state would increase the percent of schools teaching in English by about 4.6 percentage points and the percent of schools teaching English by 8.5 percentage points.

---

<sup>14</sup>The percent of schools teaching in the mother tongue can be greater than 1 due to noise in the data.

Finally, I examine the percent of schools at the district level that teach in the regional mother tongue (see columns 7-9). Since English is not the regional mother tongue anywhere in India, the percent of schools that do not teach in the mother tongue is an upper bound for those that teach in English. The results demonstrate that linguistic distance from Hindi negatively predicts the percent of schools that teach in the mother tongue at the district level. An increase in one concentric circle in the distance measure reduces the probability that a school teaches in the region's mother tongue by 3.4 percentage points.

I next study whether linguistic distance from Hindi is correlated with any other educational and economic characteristics of these districts by using various outcome variables in specification (1). The results indicate that linguistic distance from Hindi is not strongly correlated with other characteristics of the educational system in 1993 (see table 6). Neither measure of linguistic distance predicts the number of schools in 1993 or the number of higher secondary schools (grades 11 and 12) offering courses in different subjects, all normalized by the population (in 10,000s) of children aged 5-18.<sup>15</sup> Similar results for other economic variables in 1987 can be seen in table 7. Most of the variables are uncorrelated with language distance except for the percent of the population that have college degrees. The correlation with the percent of graduates is negative, but the magnitude is very small. A increase in linguistic distance to Hindi of one degree reduces the population with college degrees by 0.3%. The distance to the 10 biggest cities in India is also positive correlated with linguistic distance probably because a disproportionate number of these cities are in Hindi-speaking areas. To save space, I have omitted the results using my alternate measure of linguistic distance (the percent of speakers of languages 3 or more degrees away from Hindi) which do not differ from the results shown. I have also omitted similar regressions for the population

---

<sup>15</sup>In order to save space, I have omitted the results using the percent of speakers at each distance from Hindi which generally support these results, but are more difficult to interpret. These results are also robust to using the absolute number of schools without normalizing. However, there does appear to be a significant negative correlation between linguistic distance and total enrollment in the arts and a positive, but not significant correlation between linguistic distance and total enrollment in the sciences (results not shown). Nevertheless, while we do see a correlation of linguistic distance with enrollment in certain subjects, there does not appear to be a strong relationship between linguistic distance and the availability of schools teaching those subjects.



growth from 1987 to 1991, average wage for educated workers, the percent who have completed secondary school and the percent of households with electricity all of which are all uncorrelated with linguistic distance.

Thus, while linguistic distance does predict how many people learn English and how many schools teach English or in English, it does not strongly predict the availability of educational institutions or other economic variables, reaffirming the validity of my identification strategy.

## 4 Theoretical model

In this section, I describe a simple model to provide intuition on how the effects of globalization vary across districts due to differences in the cost of learning English. Consider an economy made up of two separate districts that differ only in the cost of learning English. I specify a schooling model and production processes, and solve the model without and then with a globally traded good, i.e. before and after trade liberalization. IT serves as one example, but we can also think of this traded good as representing all traded goods and services that require English speaking workers. Final goods travel freely between these districts, but the movement of workers is negligible. Individuals choose to work as unskilled labor or obtain instantaneous, though costly, education in either English or the local language, Hindi, and work as skilled labor. Unskilled workers cannot learn English, but everyone, including the English-educated workers, speaks Hindi. English and Hindi skilled workers are equally productive in the production of other goods, but the globally traded good requires only English-speaking workers. Production of the globally traded good also uses a second factor that is fixed in the short run. We can think of this fixed factor as representing infrastructure, such as electricity or telecommunications services, that is slow to change or entrepreneurs who are immobile. Production of all final goods is perfectly competitive.

While the model provides predictions on returns to education and the demand for education, I abstract from a central question in the trade literature. I do not attempt to explain

why India exports a skilled labor-intensive good contrary to Heckscher-Ohlin predictions. A number of papers have explored theoretical modifications to match this fact. For example, Feenstra and Hanson (1996, 1997) focus on the role on global outsourcing and Kremer and Maskin (2006)) emphasize complementarities between people of different skill levels within and across countries. Nevertheless, as this paper does not provide any insights for this related question, I assume the price of the traded good is sufficiently high.

After trade liberalization, the globally traded good is produced in both districts; this simply increases the demand for skilled workers, increasing the return to education and school enrollment. How these changes differ between the two districts is more informative. Since the elasticity of English speakers is greater in the district with lower costs of learning English, this district will produce more IT since it can do so at a lower wage and will experience greater growth in education. The wage for English speakers rises by less, but the wage for Hindi workers will rise by more. Thus, the average return to education is ambiguous.

## 4.1 Schooling Decisions

Individuals live for one period and can choose to work as unskilled labor or get instantaneous education in English or Hindi and work as skilled labor. There are  $P$  people born each period. Individuals differ in a parameter  $c_i$  which governs the cost of education and is distributed uniformly between zero and one in each district. Studying in Hindi costs  $(t + c_i)w$  where  $t$  is a fixed cost of schooling ( $0 < t < 1$ ) and  $w$  is the unskilled wage. Studying in English costs  $\mu_j c_i w$  where  $\mu_j$  is a district-specific parameter,  $j \in \{LC, HC\}$ , and  $1 < t + 1 < \mu_{LC} < \mu_{HC}$ . LC denotes the district with a lower cost of learning English, while HC is the high cost district. Note that this cost structure is not symmetric; I chose this structure to ensure that English and Hindi skilled workers exist in equilibrium both before and after liberalization. Each individual's schooling decision is to maximize lifetime income

with respect to  $h = \text{Hindi schooling} \in \{0, 1\}$  and  $e = \text{English schooling} \in \{0, 1\}$

$$\max_{h, e} \left\{ \begin{array}{ll} -\mu_j c_i w + \max \{q_E, q_H\} & \text{if } e=1, h=0 \\ -tw - c_i w + q_H & \text{if } e=0, h=1 \\ w & \text{if } e=0, h=0 \end{array} \right\}$$

where  $q_H$  and  $q_E$  are the wages for Hindi and English skilled workers, respectively. Since Hindi and English workers are equally productive in the Y sector,  $q_E \geq q_H$ . Assume that, in equilibrium,  $q_H - tw > w > q_H - tw - w$  and  $q_E > w > q_E - \mu_j w$ ; otherwise either no one would get any schooling or everyone would get schooling. Solving the individual's decision problem is straightforward. There are two cases. In one case, people with low values of  $c_i$  get English schooling, and those above do not get schooling. I ignore this case since people still learn Hindi in India. In the other case, people with low values of  $c_i$  get English schooling, those in the middle get Hindi schooling and those with higher values of  $c_i$  remain unskilled. Figure 6 demonstrates this case when  $q_E = q_H$ . The labor supply functions are

$$\begin{aligned} S_E &= PQ_{HE} = P \frac{wt + \max \{q_E, q_H\} - q_H}{w (\mu_j - 1)} \\ S_H &= P(Q_H - Q_{HE}) = P \left( \frac{q_H - w - tw}{w} - \frac{wt + \max \{q_E, q_H\} - q_H}{w (\mu_j - 1)} \right) \\ S_L &= P(1 - Q_H) = P \left( 1 - \frac{q_H - w - tw}{w} \right) \end{aligned}$$

## 4.2 Production of Y

There is one final good, Y, consumed in both districts and produced using

$$Y = \min \left\{ \frac{L_Y}{\alpha_L}, \frac{H_Y + E_Y}{\alpha_H} \right\}$$

where  $L_Y$  is quantity of unskilled labor used,  $H_Y$  is the quantity of local language skilled labor used, and  $E_Y$  is the quantity of English speaking skilled labor used and  $\alpha_L > \alpha_H$  set

the productivity levels of skilled workers relative to unskilled workers in the production of good Y. To produce an amount  $Y$ , the labor demand functions are

$$\begin{aligned} D_{LY} &= L_Y = \alpha_L Y \\ D_{HEY} &= H_Y + E_Y = \alpha_H Y \end{aligned}$$

Note that firms can perfectly substitute Hindi skilled workers for English skilled workers, so if the wage for English-speaking workers is greater, production of good Y will hire only Hindi skilled workers. The amount of Y produced will be determined by the availability of labor. Prior to trade liberalization, English speakers must work in the Y industry and therefore earn  $q_E = q_H$ . If  $q_E > q_H$ ,  $E_Y = 0$ . The zero profit condition leads to

$$p = w\alpha_L + q_H\alpha_H \quad (2)$$

When the economy is open to trade, firms can set up in either district to produce a globally traded good X and take the price of X,  $p_X$ , as given. Production of good X is Cobb-Douglas and uses English-speaking skilled workers and a fixed factor  $F$ :

$$X = F^\beta E_X^{1-\beta} \quad (3)$$

where  $E_X$  is the amount of English skilled labor used. The endowment of F in a district is exogenous and cannot respond, at least in the short run. Since F has no outside return in this model, the industry will use all F available. To determine how much E is demanded:

$$\begin{aligned} \max_E profit &= \max_E p_X F^\beta E_X^{1-\beta} - q_E E_X - r_F F \\ D_{EX} &= F \left( \frac{q_E}{p_X (1-\beta)} \right)^{-\frac{1}{\beta}} \end{aligned}$$

where  $r_F$  is the return to the fixed factor F and  $0 < \beta < 1$ . The zero profit condition is

$$p_X X = r_F F + q_E F \left( \frac{q_E}{p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} \quad (4)$$

### 4.3 Characterizing the equilibrium

In equilibrium, the labor market for unskilled workers and skilled workers must clear. For unskilled workers this condition is simple

$$D_L = S_L \Rightarrow \alpha_L Y = P \left( 1 - \frac{q_S - w - tw}{w} \right) \quad (5)$$

Labor market clearing for skilled workers depends on whether the demand for English speakers exceeds the "natural" supply of English skilled workers. Recall from figure 6 that even when  $q_E = q_H$ , there will be some supply of English speakers. If the amount of F is sufficiently small such that in equilibrium, the demand for English workers is less than this natural supply, then we must have that  $q_E = q_H$  (otherwise firms could increase profits by reducing what they pay English workers since there is an excess supply). Call this case A. The labor market clearing condition for skilled labor is

$$D_{HEY} + D_{EX} = S_H + S_E \Rightarrow \alpha_H Y + F \left( \frac{q_H}{p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} = P \left( \frac{q_H - w - tw}{w} \right) \quad (6)$$

If F is large enough that the demand for English workers exceeds the natural supply of English workers, then  $q_E > q_H$  and no English workers are hired in the Y industry. Call this case B. The labor market clearing conditions are

$$D_{HEY} = S_H \Rightarrow \alpha_H Y = P \left( \frac{q_H - w - tw}{w} - \frac{wt + q_E - q_H}{w (\mu_j - 1)} \right) \quad (7)$$

$$D_{EX} = S_E \Rightarrow F \left( \frac{q_E}{p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} = P \frac{wt + q_E - q_H}{w (\mu_j - 1)} \quad (8)$$

I set good Y to be the numeraire with a price equal to 1; since final goods can move freely across borders, this price applies to both districts. These labor market clearing conditions plus the two zero profit conditions, equations (2) and (4), and the production function for good X close the model. Denote the equilibrium values of the endogenous variables ( $w$ ,  $q_H$ ,  $q_E$  ( $= q_H$  in case A),  $r$ ,  $Y$  and  $X$ ), with an asterisk. In addition, it will be useful to have defined two additional terms. Define the weighted average return to skill,  $\hat{q}$  as

$$\hat{q} = \frac{q_H H + q_E E}{w(H + E)}$$

and the total number of educated people,  $ED$

$$ED = H + E = P \left( \frac{q_H}{w} - 1 - t \right)$$

The equilibrium without any trade is a special case of case A when  $F = 0$ , described in Proposition 1 below. Since the demand for English skilled workers rises after trade liberalization, the wage for English workers has to rise. Now that fewer English speakers are working in the  $Y$  industry, the wage for Hindi skilled workers has to rise as well since the ratio of skilled to unskilled workers in  $Y$  production has to remain constant. School enrollment responds to these higher wages. To compare how these changes differ in districts with different levels of  $\mu_j$ , we have to first solve the equilibrium in both case A and B.

**Proposition 1** *Case A. If, in equilibrium for a small neighborhood around  $\mu_j$ ,  $q_E^* = q_H^*$*

*then a) for all variables  $Z \in \{w^*, q_H^*, r_F^*, Y^*, X^*, \hat{q}^*, ED^*\}$*

$$\frac{dZ}{d(\mu_j - 1)} = 0$$

*b)*

$$\frac{dE^*}{d(\mu_j - 1)} < 0 \quad \text{and} \quad \frac{dH^*}{d(\mu_j - 1)} > 0$$

c) and the following condition must hold

$$F \left( \frac{q_H^*}{p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} \leq P \frac{t}{\mu_j - 1} \quad (9)$$

**Proof.** See appendix. ■

If the two districts LC and HC are both in this case, they will have identical wages, production of X and returns to education. They will also have identical total education, but the low cost district will have a higher proportion of English speakers. I next solve case B.

**Proposition 2** *Case B. If, in equilibrium for a small neighborhood around  $\mu_j$ ,  $q_E^* > q_H^*$*

*then a) for all variables  $Z \in \{w^*, q_E^*, Y^*, H^*\}$*

$$\frac{dZ}{d(\mu_j - 1)} > 0$$

*and for all variables  $Z \in \{q_H^*, r_F^*, X^*, E^*, ED^*\}$*

$$\frac{dZ}{d(\mu_j - 1)} < 0$$

b) and the following condition must hold

$$F \left( \frac{q_H^*}{p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} > P \frac{t}{\mu_j - 1}$$

**Proof.** See appendix. ■

The effect of a small change in  $\mu_j$  on the return to education,  $\hat{q}^*$  is ambiguous. In the next subsection I provide some intuition and calibrate the model to determine which district has a higher average return to education.

Intuitively, a district is no longer in case A when the demand for English workers is greater than the natural supply. Therefore, the district with a higher cost of learning English, will leave case A at a lower value of  $F$  since it will exhaust its smaller natural amount of English

speakers sooner. I can now rewrite equation (9) in terms of exogenous variables.

**Proposition 3**  $q_E^* = q_H^*$  holds if and only if

$$F \leq P \frac{t}{\mu_j - 1} \left[ \frac{\alpha_L}{\alpha_H p_X (1 - \beta)} \frac{1}{\left( \frac{1}{\alpha_L} - \frac{\left( \frac{1}{\alpha_H} + \frac{1}{\alpha_L} \right)}{\left[ (2 + t) \left( 1 + \frac{\alpha_H}{\alpha_L} \right) + \frac{\alpha_L}{\alpha_H} \right] + \frac{t}{\mu_j - 1} \right)} \right]^{\frac{1}{\beta}} = \bar{F}(\mu_j) \quad (10)$$

**Proof.** See appendix. ■

District HC will leave case A at a lower value of F since  $\frac{d\bar{F}(\mu_j)}{d(\mu_j - 1)} < 0$ .

#### 4.4 Average return to education in case B

The intuition concerning the difference in the return to education between the two districts is standard. English-skilled workers are less elastic in district HC; therefore the wage must rise by more to get additional English-skilled workers, but in equilibrium, it cannot rise by enough to get the same number of English speakers as in district LC. Since fewer skilled workers are taken out of  $Y$  production,  $q_H^*$  increases by less in district HC. Thus, the weighted average return to education depends on the relative magnitudes of these two wages and the proportion of workers who study in English and Hindi, which in turn depend on the parameters. Using calibrated values of  $\alpha_L$ ,  $\alpha_H$ ,  $t$ ,  $\mu_C$  and  $\mu_T$ , I can show that when production of  $X$  is sufficiently intensive in the fixed factor,  $F$ ,  $\hat{q}^*$  is greater in district HC for all reasonable values of  $F$ . When production of  $X$  relies less on the fixed factor,  $\hat{q}^*$  is greater in district HC for small values of  $F$  but smaller in district LC for larger values.

I first calibrate  $\alpha_L$ ,  $\alpha_H$ , and  $t$  to match returns to high school and college and the percent of high school and college graduates in urban areas of India. According to the 1991 Census of India, 21.2% of people in urban areas of India have at least completed secondary school. Wage regressions using the NSS data described below demonstrate that high school increases wages by about 50% while college increases wages by about 100% in 1987. Using the fact that only 6.7% of people in urban areas have a college degree, I calculate a weighted average return



to education of 66% in a given year. Assuming a 10% interest rate, this is approximately an 84.2% lifetime return. I set  $\alpha_L = 1$ ,  $\alpha_H = 0.25$ , and  $t = 0.65$ . I then calibrate  $\mu_C = 9$  and  $\mu_T = 6$  to match the percent of people who learn English in districts that speak languages closer to (8.7%) and farther from Hindi (13.5%), respectively. Setting  $p_X = 1$  and  $P = 1$ , I am left with two free parameters,  $\beta$  and  $F$ .

Figure 7 demonstrates what happens to  $\hat{q}^*$  and  $\frac{d\hat{q}^*}{d(\mu_j-1)}$  in districts HC and LC as  $F$  rises when  $\beta = 0.3$ .  $\frac{d\hat{q}^*}{d(\mu_j-1)}$  is plotted against the left axis, while  $\hat{q}^*$  is plotted against the right axis. In addition, since  $F$  is difficult to interpret in real-world terms, I plot income from  $X$  production as a percent of total income with respect to  $F$  in figure 8. The return to education is always greater in district HC from small values of  $F$  and even until  $F$  is large enough that the  $X$  contributes around 90% of GDP. At this point  $\frac{d\hat{q}^*}{d(\mu_j-1)}$  is upward sloping in both districts implying that for all values of  $F$ , the return to education in district HC is greater than in district LC. I repeat this exercise in figure 9 using  $\beta = 0.18$ . Here,  $\hat{q}^*$  is greater in district HC only when  $F$  is between 1.05 and 6.8, corresponding to a percent of IT in GDP of 15% to 36.5%. At an even lower value of  $\beta = 0.1$ ,  $\hat{q}^*$  is greater in district C only when the percent of IT in GDP is between 15 and 20%.

This model has 3 main predictions. First, the district with lower costs of learning English produces more  $X$ . Second, the low cost district will employ more educated workers and have more educational attainment. Finally, the English wage will be lower but the Hindi wage will be greater. Since the data does not identify whether individuals speak English, I calculated the weighted average return to education. The prediction regarding this average return to education is ambiguous and depends on the  $F$ -intensity of  $X$  production. If  $F$  is important in  $X$  production, then the weighted average return to education is greater in the district with a higher cost of learning English. If  $F$  is less important, then for low values of  $F$ , we have the same prediction, but for higher values of  $F$ , the prediction is reversed.

## 5 Empirical methodology

### 5.1 Information technology and returns to education

To test the first prediction, I estimate the impact of linguistic distance from Hindi on geographic variation in IT firm location and growth. The impact on IT employment growth is one benefit of promoting the use of a global language, since IT firms hire mostly well-educated, English-speaking individuals. I estimate the following equation

$$T_{jt} = \alpha_0 + \beta' D_j + \alpha'_1 Z_j + \alpha'_2 Z_j^2 + \gamma_t + \gamma_g + \nu_{jt}$$

where  $T_j$  (for technology) is a measure of IT presence in district  $j$  in year  $t$ ,  $D_j$  is a vector of measures capturing the linguistic distance of languages spoken in district  $j$  from Hindi and  $Z_j^2$  is a vector of control variables. The measures of IT presence used are described below and include the existence of any IT headquarters or branches and the number of years IT firms have been present in the district.  $Z_j$  is as above (see equation (1)). The vector  $Z_j^2$  includes the natural log of district population, the number of engineering colleges in the district, the distance to the closest airport, and the percent of non-migrant adults with an engineering degree. I drop the ten most populous districts as of 1991 and cluster the standard errors by district. I also run this regression at the firm-level to study year of firm establishment.

To test the second and third predictions, I study returns to education via employment opportunities and wages. Unfortunately, the data does not distinguish between English-medium education and local language instruction, but since people who speak English almost always learned it in school, I focus on returns to education. I estimate the following regression

$$\begin{aligned} J_n = & \alpha_0 + \beta'_1 D_j \cdot I(t = 1999) + \beta'_2 D_j \cdot I(t = 1999) \cdot HS_n + \beta'_3 D_j \cdot I(t = 1999) \cdot C_n \\ & + \alpha'_1 D_j \cdot HS_n + \alpha'_2 D_j \cdot C_n + \alpha_3 I(t = 1999) \cdot HS_n + \alpha_4 I(t = 1999) \cdot C_n \\ & + \alpha_5 HS_n + \alpha_6 C_n + \alpha'_7 Y_n^1 + \alpha'_8 Y_j^2 \cdot I(t = 1999) + \theta_j + \gamma_t + \gamma_{gt} + \mu_n \end{aligned}$$

where  $J_n$  is a measure of employment of individual  $n$  in district  $j$  at year  $t \in \{1987, 1999\}$ ,  $HS_n$  and  $C_n$  are indicators for whether individual  $n$  has completed high school or college, respectively and  $Y_n^1$  and  $Y_j^2$  contain individual and district characteristics. For  $J_n$ , I use whether the individual has a regular wage or salaried job conditional on being in the labor force (as opposed to being self-employed, a casual laborer or seeking work) and the natural log of the individual's wage earnings per week. The vector  $Y_n^1$  includes the individual's age, age squared, gender, marital status (a dummy for being married), and whether the individual has ever moved. At the district level,  $Y_j^2$  includes the percent of native English speakers, the percent of native Hindi speakers and whether the state is in the Hindi belt. I also include district fixed effects and region fixed effects interacted with time, cluster the standard errors by district and weight the observations. I run this regression by gender and age and study the skilled wage premium across industries to establish the role played by trade in services.

## 5.2 Linguistic distance from Hindi and school enrollment

To further test the second prediction, I estimate the impact of linguistic distance from Hindi on school enrollment growth using

$$\begin{aligned} \log(S_{ijt}) - \log(S_{ijt-1}) &= \alpha_0 + \beta' D_j \cdot I(t = 2002) + \alpha_1 \log(S_{ijt-1}) + \alpha_2' P_{jt} \\ &+ \alpha_3 I(t = 2002) + \alpha_4' Z_j \cdot I(t = 2002) + \theta_j + \gamma_t + \gamma_{gt} + \gamma_{it} + \mu_{ijt} \end{aligned} \quad (11)$$

where  $S_{ijt}$  is a measure of enrollment at grade level  $i$  in district  $j$  in region  $g$  at time  $t \in \{1987, 1993, 2002\}$ ,  $D_j$  is a vector of measures capturing the linguistic distance of languages spoken in district  $j$  from Hindi,  $I(\cdot)$  is an indicator function, and  $P_{jt}$  is a vector of population at time  $t$  and time  $t - 1$ . The parameter  $\theta_j$  is a district fixed effect allowing each district to have its own trend in enrollment and  $\gamma_{it}$  is a cohort effect, controlling for grade-year trends. The vector  $Z_j$  includes the same controls as above. I use only data from the urban areas of all districts and correct the standard errors for intracluster correlation at the district level.

I also include fixed effects for region interacted with time. This strategy - using preexisting differences to predict changes in school enrollment, allowing for individual district effects and region-specific trends and controlling for many other preexisting differences - rules out a number of other explanations for these results.

Using individual level data, I confirm these results using the following regression:

$$\begin{aligned}
H_c = & \alpha_0 + \beta'_1 D_j \cdot I(t = 1999) + \beta'_2 D_j \cdot I(t = 1999) \cdot I(A_c \in [11, 15]) \\
& + \beta'_3 D_j \cdot I(t = 1999) \cdot I(A_c \in [16, 20]) + \alpha'_1 D_j \cdot I(A_c \in [11, 15]) \\
& + \alpha'_2 D_j \cdot I(A_c \in [16, 20]) + \alpha_3 I(t = 1999) \cdot I(A_c \in [11, 15]) \\
& + \alpha_4 I(t = 1999) \cdot I(A_c \in [16, 20]) + \alpha'_5 Y_c^1 + \alpha'_6 Y_j^2 + \theta_j + \gamma_t + \gamma_{gt} + \mu_c
\end{aligned} \tag{12}$$

where  $H_c$  denotes whether child  $c$  (between ages 6 and 20) in district  $j$  is attending school at year  $t \in \{1987, 1999\}$ ,  $A_c$  denotes the child's age, and  $Y_c^1$  and  $Y_j^2$  are vectors of characteristics about the child and the district respectively. I also include district fixed effects and region fixed effects interacted with time. I weight the observations and cluster the standard errors by district. The vector  $Y_c^1$  consists of fixed effects for age, gender, household religion, education categories of the household head; an indicator for having migrated and household wage income.  $Y_j^2$  is as above. I also run the same regression with a measure of educational achievement as the outcome. This variable takes on values 0 for children who are not literate, 1 for those who are literate without any formal education, 2 for those who have completed pre-primary schooling, 3 for primary, 4 for middle, 5 for secondary and 6 for college graduates. Due to the ordinal nature of this variable, I run this regression with an ordered probit model as well as linear probability.

While linguistic distance to Hindi is one measure of the cost of learning English, we might want to estimate the impact of other measures. For example, I estimate all specifications using the percent of schools in the district that teach in the mother tongue. However since this is endogenous, I instrument with measures of linguistic distance to Hindi.

## 6 Sources of data

### 6.1 Information technology

I collected and coded data on IT firms from the National Association of Software and Service Companies (NASSCOM) directories published in 1995, 1998, 1999-2000, 2002 and 2003. These directories are based on surveys that contain self-reported data from individual firms on the location of the headquarters, the number and location of branches, and the number of employees. According to NASSCOM, the sample accounts for 95% of the revenue in the industry in most years (Mehta 1995, 1998, 1999; Karnik 2002). To estimate the impact of linguistic distance from Hindi on the growth in IT, I use a number of different measures of IT presence. My primary measure is an indicator for the existence of IT headquarters or branches in a district, but I also examine the year of firm establishment, the number of branches and headquarters in a district, number of employees, total revenue, total exports and total capital subscribed. Since the data only contains this information at the headquarters level, I either assign them to the headquarters or divide evenly across branches. Summary statistics of IT presence by year of data can be found in table 8. Note that revenue and exports data were not available in the 1998 directory.

### 6.2 Employment, education in 1987 and district-level controls

Data on employment and returns to education are from the Indian National Sample Surveys (NSS) conducted in 1987-1988 and 1999-2000.<sup>16</sup> The NSS provides individual-level information on wages paid in cash and in-kind as well as employment status, industry and occupation codes and weights for each observation. Employment status includes working in household enterprises (i.e. self-employed), as a helper in a household enterprise, as a regular salaried wage employee and as casual wage labor. In addition, the data reveals whether individuals are seeking work, attending school or attending to domestic duties.

---

<sup>16</sup>An NSS survey was conducted in 1993-1994, but district identifiers are not available.

Matching industry codes across the two time periods allows examining employment growth and returns to education by industry in agriculture, manufacturing, wholesale/retail/repair, hotel & restaurant services, transport services, communications (post and courier), financial intermediate/insurance/real estate and other services (education, health care, civil services). Table 9 provides district averages of employment and wages for the two years in the sample.

This data was also used to calculate enrollment levels from 1987 and district-level control variables. I construct district-level measures of grade school enrollment at the primary, upper primary and secondary levels. The NSS also contains household-level information on household structure, demographics, employment, education, expenditures, migration and assets, from which I calculate district averages of all control variables mentioned above such as household wage income, the percent of working-age adults who are engineers, the percent Muslim, the percent who regularly travel by train and the percent of households that have electricity. In addition, I use this individual data as a second check of the education results, although I prefer the SAIES data for these education results since the latest year is more recent (2002 as opposed to 1999-2000) and as the SAIES is a census, it provides a more complete measure of district-wide enrollment.

Using latitude and longitude data, I calculate the distance from each district to the closest of the 10 biggest cities in India and to the closest airport operated by the Airports Authority of India. As a measure of engineering college presence I count the number of Indian Institutes of Technology and Regional Engineering Colleges (now called the National Institutes of Information Technology) in each district. All of them were established prior to 1990, although some were not given REC/NIIT status until the later 1990s. Summary statistics for all of these variables can also be found in panel C of table 2.

### **6.3 Enrollment, achievement and language of instruction**

Data on grade school enrollment comes from the Sixth and Seventh All India Educational Surveys (AIES), conducted by the National Council of Educational Research and Training

(NCERT) which began on September 30, 1993 and September 30, 2002 respectively. The surveys collect data at the school level on enrollment, school facilities, languages taught, courses available, teacher qualifications, incentive programs and other aspects of education. I focus on urban enrollment at each grade in a district. Note that the data consist of the actual number of students enrolled in each grade, not enrollment rates. In order to account for this, I control for population and population growth among 5-19 year-olds from the 1991 and 2001 Census of India. I cannot calculate a more accurate measure of the relevant population for each grade since there is no definitive age at which a child should be in a particular grade, given that children start school at different ages and often fail to be promoted. Enrollment in each grade increased dramatically between 1993 and 2002, as can be seen in table 10. Total enrollment in all grades increased by 32%, but the increases are greater in the highest grades. I pool all grades and also separate them by school level (grades 1-5 constitute primary, 6-8 upper primary, 9-12 secondary).

## **7 Impact of linguistic distance on returns to education**

### **7.1 Geographic variation in the growth of information technology**

The results indicate strong positive effects of linguistic distance from Hindi on IT presence (see tables 11 to 13). I find that this cost of learning English, when measured either as the weighted average or the percent of speakers sufficiently removed from Hindi, predicts whether any IT firm establishes a headquarters or a branch in a district (see columns 1 and 2 of table 11). One degree away from Hindi of the average speaker in a district results in a 6% increase in the probability of having any IT presence. However, many districts may be very unlikely to receive IT firms for a number of reasons. While these reasons should be orthogonal to linguistic distance to Hindi, it is possible they are not, so I use firm-level data to focus on districts that see any IT firms between 1995 and 2003 (columns 3 and 4 of table 11). These results confirm that areas that are linguistically further removed from Hindi saw the

establishment of IT firms earlier, by approximately 3 years per degree of linguistic distance.

Linguistic distance from Hindi also explains geographic variation in the number of headquarters and branches and the number of employees when divided evenly by branch (see columns 1-4 of table 12). However, it does not predict IT firm employment when assigned to the firm headquarters or firm performance, when measured either way (columns 5-12 of table 12). There are many possible explanations. First, many employees and much of a firm's production may not be located at the headquarters. A firm may set up in a district to which the founder has personal ties but produce all its revenue at a different branch. Either allocation method would then be incorrect and bias the results. Second, the entire effect of linguistic distance on IT could be at the extensive margin of where firms establish, not on the intensive margin. Another possible explanation is that the relationship between linguistic distance and firm performance could be nonlinear or nonparametric and not captured by these reduced form measures of linguistic distance from Hindi. This is suggested when I instrument for the percent of schools that teach in the regional mother tongue with the percent of people in the district that speak languages at each distance away from Hindi (see table 13). The F-statistics are sufficiently strong and the coefficients suggest an impact of the correct sign, but the standard errors are too large to judge significance.

## **7.2 Returns to education**

Testing the second and third predictions, I find evidence of greater growth in employment of educated workers but smaller growth in wage premiums in districts with lower costs of English as predicted by the model. I first demonstrate that the college premium for the probability of regular employment rose faster in districts with lower costs of learning English from 1987 to 1999 (see table 14). These effects are driven by increases in the employment for young adults (below the age of 30) rather than older adults. These results are sensible given that many of the new firms in services that have risen due to trade, such as in IT, hire predominantly young adults. The coefficients are also larger for women than men, which is



also sensible since IT firms employ more women relative to traditional Indian firms. The male-female ratio among those working was 80:20 according to the 1987 NSS, but 77:23 in software firms and 35:65 in business processing firms (NASSCOM 2004).

At the same time, confirming the third prediction, I find that skilled wage premiums rose by less in districts with lower costs of English, particularly for secondary school graduates (see table 15). These results are driven by wages for older adults and the magnitudes are relatively small. The wage premium for high school graduates rises by 3% less over 12 years per degree in linguistic distance from Hindi, relative to a premium of 54% for high school graduates in 1987. Further exploring these wage results by industry, I find that the fall in wage premium for educated workers in these districts seems to be driven by wages in manufacturing, hotels and restaurants, transportation and communications (see tables 16 and 17). I find no evidence of differential changes in wage premium by linguistic distance from Hindi in agriculture and wholesale, retail and repair. I also find an increase in overall wages in transportation and communications. Studying other industries, I find that high school and college wage premiums rose more in districts linguistically further from Hindi only in the business services industry, which includes financial institutions, insurance, real estate, computer related activities, research and development and other business activities. In addition, I find no wage effects on other services which includes public service, education, health, sanitary and community services.

These results are clearly consistent with an increase in trade-related jobs in districts with lower costs of learning English since most trade in services requiring English would be in business services. While trade in manufacturing could increase as well following the trade liberalization in the 1990s, these results indicate that wages did not rise in the manufacturing industry, perhaps because tariff cuts in the 1990s resulted in import competition in manufacturing. Export growth in India since the 1990s was driven by services. The increase in wages in the transportation and communications industries are also consistent with this story since they are important to trade-related services.

## 8 Impact of linguistic distance on education

Further testing the second prediction, I find that educational attainment rises more in districts with lower costs of learning English. I first demonstrate this result using district-level data, estimating specification (11) (see tables 18 and 19). Both measures of linguistic distance show an increase in school enrollment, although the results are stronger at different levels. At the primary and upper primary levels, the coefficients for girls are larger than boys, but the increase is comparable at the secondary level. This is partly because enrollment for girls starts from a lower level than boys and the outcome is a percent improvement.

I find similar results, when using the percent of each type of school (primary and upper primary) that teach in the regional mother tongue as the measure of cost of learning English and instrumenting with the linguistic distance (see table 20). I instrument both with the weighted average measure of linguistic distance as well as the percent of speakers each degree away from Hindi as in table 13. The F-statistics from the excluded variables in the first stage are shown at the bottom of the table; these instruments clearly have sufficient predictive power. The second set of instruments allows for a more flexible relationship between distance of mother tongues from Hindi and the prevalence of English, increasing the predictive power of the first stage and significantly reducing the standard errors (at least for primary schools).

The magnitudes of these results are economically significant. Using the estimates from column 1 of table 18, an increase in 1 degree from Hindi of the average speaker's mother tongue (44% of 1 standard deviation) would increase enrollment growth by 11% over the 9-year period. Using estimates from column 1 in table 19, an increase of 1 standard deviation in how many people speak languages far from Hindi (44%) would increase enrollment by 12% over the 9-year period. Finally, a 1 standard deviation increase in how many primary schools teach in the mother tongue (22%) would reduce enrollment by 43% over 9 years (averaging the effect using different instrument sets). For upper primary schools, an increase in 1 standard deviation (25 percentage points) would reduce enrollment by 39% over 9 years.

As an additional test of these results, I examine individual data from another source.

Note that the time period in this exercise is 1987 to 1999-2000, so we would expect to find smaller effects. Estimating equation (12) on the NSS data, I find that the main interaction of linguistic distance and post is not significant (in fact, the coefficients are often negative), but the interactions with older age groups is quite significant (see table 21). For boys, the most precisely estimated increase in attendance is at ages 16 - 20 while the attendance for 11 - 15 year-old girls seems to respond the most. I find similar effects using the alternate dependent variable of educational achievement (see table 22).<sup>17</sup>

## 9 Conclusion

In this paper, I studied the effect of promoting a global language on education and employment during a period of trade liberalization and globalization, by exploiting exogenous variation in the cost of learning English across districts in India. This exogenous variation was determined by historical linguistic diversity in India and the imposition of Hindi as the official language of the country in the 1950s. I first showed that linguistic distance from Hindi does predict whether individuals choose to learn English as a second language and that linguistic distance from Hindi of languages spoken in a district does predict how many schools teach English. I showed that one clear benefit of promoting a global language is access to global job opportunities such as those in information technology. IT firms were established in districts that are linguistically further from Hindi with a greater probability and earlier than in other districts. I next demonstrated that the high school and college premium in the probability of having a regular salaried job rose faster, but the wage premium rose by less for individuals in these districts. Lastly, I showed that districts further removed from Hindi experienced greater increases in school enrollment growth.

There are two avenues through which new job opportunities created through trade liberalization may have increased school enrollment. In this paper, I demonstrated that returns

---

<sup>17</sup>I get similar results using an ordered probit model. I also repeat the exercise from table 20 using this specification and find the presence of fewer English-instruction schools reduces enrollment growth particularly for older children (results not shown).

to education rose faster in districts where residents speak languages further from Hindi. Another channel could be through increased family income. It is unlikely that this is driving all the results presented here since the greater job opportunities were concentrated among young adults. Thus, this should increase enrollment of children at lower grades by more while the results indicate bigger effects at older ages. In future work, I plan to further distinguish between these two effects by exploiting the household structure of the NSS data and comparing households that are more affected by the increase in labor market opportunities with other households that would only be affected by general income growth in the district. Unfortunately, the data is not currently available past 1999. I also hope to find additional data sources which may contain language knowledge for the individual to determine more directly the consequences of speaking English.

In addition, the interaction of linguistic distance to Hindi and trade liberalization may have impacted fertility and marriage rates. Anecdotal evidence suggests that women work in call centers and other business services firms between finishing their education and getting married. This might then increase the age of first marriage and it might increase the probability that women continue to work past marriage, potentially impacting fertility rates. Finally, in the long run, there might be an impact of future job opportunities on child health. If parents think their daughters are more likely to work for a few years before getting married (since some business processing firms such as call centers and medical transcription firms prefer to hire women) they may invest more in their daughters' health as well as education.

Thus, I demonstrated how the impact of trade liberalization varies across regions with different elasticities of skilled labor - areas with greater labor supply response experience greater employment of skilled workers and human capital accumulation, but smaller increases in wage inequality as measured by skilled wage premiums.

## 10 References

Angrist, Joshua, Aimee Chin and Ricardo Godoy (2006). "Is Spanish-Only Schooling Responsible for the Puerto Rican Language Gap?," *NBER Working Paper* 12005, National Bureau of Economic Research, Cambridge, MA.

Angrist, Joshua and Lavy, Victor (1997). "The Effect of a Change in Language of Instruction on the Returns to Schooling in Morocco," *Journal of Labor Economics*, University of Chicago Press, vol. 15(1), pages S48-76, January.

Attanasio, O., Goldberg P., and N. Pavcnik (2004). "Trade Reforms and Wage Inequality in Colombia," *Journal of Development Economics* 74, 331-366.

Attanasio, O. and M. Szekely (2000): "Household Saving in East Asia and Latin America: Inequality Demographics and All That", in B. Pleskovic and N. Stern (eds.), Annual World Bank Conference on Development Economics 2000. Washington, DC: World Bank.

"Busy signals: Too many chiefs, not enough Indians," *The Economist*, September 8, 2005.

"Can India Fly? A Special Report," *The Economist*, June 3-9, 2006.

Clingingsmith, David (2006). "Bilingualism, Language Shift and Economic Development in India, 1931-1961." (mimeo) Harvard University.

Cragg, M.I. and M. Epelbaum (1996). "Why Has Wage Dispersion Grown in Mexico? Is It Incidence of Reforms or Growing Demand for Skills?" *Journal of Development Economics* 51(1), 99-116.

Dakin, Julian, Brian Tiffen, and H.G. Widdowson (1968). *Language in Education: The Problem in Commonwealth Africa and the Indo-Pakistan Sub-continent*. Oxford: Oxford University Press.

Dyen, Isidore, Joseph Kruskal and Paul Black (1997). FILE IE-DATA1. Available at <http://www.ntu.edu.au/education/langs/ielex/HEADPAGE.html>.

Edmonds, Eric, Nina Pavcnik and Petia Topalova (2006). "Trade Policy, Child Labor, and Schooling: Evidence from Indian Districts." (mimeo) Dartmouth College.

Feenstra, R.C. and G. Hanson (1996). "Foreign Investment, Outsourcing and Relative Wages." In R.C. Feenstra, G.M. Grossman and D.A. Irwin, eds., *The Political Economy of Trade Policy: Papers in Honor of Jagdish Bhagwati*, MIT Press, 89-127.

Feenstra, R.C. and G. Hanson (1997). "Foreign Direct Investment and Relative Wages: Evidence from Mexico's Maquiladoras." *Journal of International Economics*, 42(3), 371-393.

Feliciano, Z. (1993). "Workers and Trade Liberalization: The Impact of Trade Reforms in Mexico on Wages and Employment." (mimeo) Harvard University.

Goldberg, Pinelopi K. and Nina Pavcnik (2004). "Trade, Inequality, and Poverty: What Do We Know? Evidence from Recent Trade Liberalization Episodes in Developing Countries," *Brookings Trade Forum*, Washington, DC: Brookings Institution Press: 223-269.

Hanson, Gordon H. and Ann Harrison (1999). "Trade, Technology and Wage Inequality in Mexico." *Industrial and Labor Relations Review* 52(2), 271-288.

Hohenthal, Annika (2003). "English in India; Loyalty and Attitudes," *Language in India*, **3**, May 5, 2003.

Kamat, A.R., 1985. *Education and Social Change in India*. Bombay: Somaiya Publications Private Limited.

Karnik, Kiran, ed. (2002). *Indian IT Software and Services Directory 2002*. National Association of Software and Service Companies, New Delhi.

Kremer, Michael and Eric Maskin (2006). "Globalization and Inequality." (mimeo) Harvard University.

Lang, Kevin and Erez Siniver (2006). "The Return to English in a Non-English Speaking Country: Russian Immigrants and Native Israelis in Israel," *NBER Working Paper* 12464, National Bureau of Economic Research, Cambridge, MA.

Lindert, Peter H. and Jeffrey G. Williamson (2001). "Does Globalization Make the World More Unequal?" *NBER Working Paper* No. 8228, National Bureau of Economic Research, Cambridge, MA.

Mehta, Dewang, ed. (1995). *Indian Software Directory 1995-1996*. National Association of Software and Service Companies, New Delhi.

Mehta, Dewang, ed. (1998). *Indian Software Directory 1998*. National Association of Software and Service Companies, New Delhi.

Mehta, Dewang, ed. (1999). *Indian IT Software and Services Directory 1999-2000*. National Association of Software and Service Companies, New Delhi.

Munshi, Kaivan & Mark Rosenzweig (2006). "Traditional Institutions Meet the Modern World: Caste, Gender, and Schooling Choice in a Globalizing Economy," *American Economic Review*, American Economic Association, vol. 96(4), pages 1225-1252, September.

NASSCOM, 2004. *Strategic Review 2004*. National Association of Software and Service Companies, New Delhi, 185-194.

Nurullah, Syed and J.P. Naik, 1949. *A Student's History of Education in India, 1800-1947*. Bombay: Macmillan and Company Limited.

Panagariya, Arvind (2003). "The WTO Trade Policy Review of India, 1998." *International Trade* 0309012, EconWPA.

Panagariya, Arvind (2004). "India's Trade Reform: Progress, Impact and Future Strategy" *International Trade* 0403004, EconWPA.

Robbins, Donald (1995a). "Earnings Dispersion in Chile after Trade Liberalization." Harvard Institute for International Development, Cambridge, MA.

Robbins, Donald (1995b). "Trade, Trade Liberalization, and Inequality in Latin America and East Asia: Synthesis of Seven Country Studies." Harvard Institute for International Development, Cambridge, MA.

Robbins, Donald (1996a). "Stolper-Samuelson (Lost) in the Tropics: Trade Liberalization and Wages in Colombia 1976-94." Harvard Institute for International Development, Cambridge, MA.

Robbins, Donald (1996b). "HOS Hits Facts: Facts Win. Evidence on Trade and Wages in the Developing World." Harvard Institute for International Development, Cambridge, MA.

Robbins, Donald, and Thomas Gindling (1997). "Educational Expansion, Trade Liberalisation, and Distribution in Costa Rica." In Albert Berry, ed., *Poverty, Economic Reform and Income Distribution in Latin America*. Boulder, Colo.: Lynne Rienner Publishers.

Robbins, Donald, Martin Gonzales, and Alicia Menendez (1995). "Wage Dispersion in Argentina, 1976-93: Trade Liberalization amidst Inflation, Stabilization, and Overvaluation." Harvard Institute for International Development, Cambridge, MA.

Sanchez-Paramo, C. and N. Schady (2003): "Off and Running? Technology, Trade, and the Rising Demand for Skilled Workers in Latin America," World Bank Policy Research

Working Paper 3015. Washington, DC: World Bank.

Swadesh, Morris (1972). "What is glottochronology?" In M. Swadesh, The origin and diversification of languages. London: Routledge & Kegan Paul: 281-284.

United Nations Development Programme, Human Development Report 2004: Cultural Liberty in Today's Diverse World, New York: Oxford University Press, 2004.

Wei, Shang-Jin and Yi Wu (2001). "Globalization and Inequality: Evidence from Within China." *NBER Working Paper* No. 8611, National Bureau of Economic Research, Cambridge, MA.

Wood, Adrian (1997). "Openness and Wage Inequality in Developing Countries: The Latin American Challenge to East Asian Conventional Wisdom." *World Bank Economic Review*, 11(1), 33-57.

## 11 Appendix A

### 11.1 Proof of Proposition 1

a) From equations (5) and (6), we can show that

$$\frac{\alpha_H}{\alpha_L} P \left( 1 - \frac{q_H^* - w^* - tw^*}{w^*} \right) + F \left( \frac{q_H^*}{p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} = P \left( \frac{q_H^* - w^* - tw^*}{w^*} \right)$$

Substituting  $q_H^* = \frac{1}{\alpha_H} - \frac{\alpha_L}{\alpha_H} w^*$  into this expression implicitly solves for  $w^*$ :

$$(2 + t) \left( 1 + \frac{\alpha_H}{\alpha_L} \right) + \frac{\alpha_L}{\alpha_H} = \frac{1}{w^*} \left( \frac{1}{\alpha_H} + \frac{1}{\alpha_L} \right) - \frac{F}{P} \left( \frac{1 - \alpha_L w^*}{\alpha_H p_X (1 - \beta)} \right)^{-\frac{1}{\beta}} \quad (13)$$

Note that this expression does not depend on  $\mu_j$ . Thus,

$$\frac{dw^*}{d(\mu_j - 1)} = 0$$

The variables,  $q_H^*$ ,  $r_F^*$ ,  $Y^*$ ,  $X^*$ ,  $\hat{q}^*$ ,  $ED^*$  can be written as functions of  $w^*$  which also do not depend on  $\mu_j$ .

b) These follow from the expressions,  $E = P \frac{t}{\mu_j - 1}$  and  $H = P \left( \frac{1}{w^* \alpha_H} - \frac{\alpha_L}{\alpha_H} - 1 - t - \frac{t}{\mu_j - 1} \right)$ .

c) Since  $q_E^* = q_H^*$ , we know that the supply of English skilled workers does not depend on the wages. To be in this equilibrium, the demand for English skilled labor from  $X$  production

must be less than or equal to the supply; English speakers not working in the  $X$  industry can work in  $Y$  production and earn the same wage.

$$F\left(\frac{q_E^*}{p_X(1-\beta)}\right)^{-\frac{1}{\beta}} = F\left(\frac{q_H^*}{p_X(1-\beta)}\right)^{-\frac{1}{\beta}} \leq P \frac{t}{\mu_j - 1}$$

## 11.2 Proof of Proposition 2

a) I first solve for the equilibrium in case B. From equations (5) and (7), we have

$$\frac{a_L}{\alpha_H} P \left( \frac{q_H^* - w^* - tw^*}{w^*} - \frac{w^*t + q_E^* - q_H^*}{w^*(\mu_j - 1)} \right) = P \left( 1 - \frac{q_H^* - w^* - tw^*}{w^*} \right)$$

Substituting  $q_H^* = \frac{1}{\alpha_H} - \frac{a_L}{\alpha_H} w^*$  into this expression and solving for  $q_E^*$  gives us

$$q_E^* = \frac{1}{\alpha_H} + (\mu_j - 1) \left( \frac{1}{\alpha_L} + \frac{1}{\alpha_H} \right) - w^* \left[ \frac{a_L}{\alpha_H} + t + (\mu_j - 1) \left[ \frac{a_L}{\alpha_H} + \left( \frac{a_H}{\alpha_L} + 1 \right) (2 + t) \right] \right] = A - w^* B$$

Define

$$A = \frac{1}{\alpha_H} + (\mu_j - 1) \left( \frac{1}{\alpha_L} + \frac{1}{\alpha_H} \right) \quad \text{and} \quad B = \left[ \frac{a_L}{\alpha_H} + t + (\mu_j - 1) \left[ \frac{a_L}{\alpha_H} + \left( \frac{a_H}{\alpha_L} + 1 \right) (2 + t) \right] \right]$$

Plugging these expressions for  $q_H$  and  $q_E$  into equation (8), we get

$$0 = \left( \frac{P}{F} \right)^{-\beta} \left( \frac{1}{w^*} \left( \frac{1}{\alpha_L} + \frac{1}{\alpha_H} \right) - \left[ \frac{a_L}{\alpha_H} + \left( \frac{a_H}{\alpha_L} + 1 \right) (2 + t) \right] \right)^{-\beta} - \frac{A - w^* B}{p_X(1-\beta)} = G(w; \mu_j) \quad (14)$$

Thus,

$$\frac{dw}{d(\mu_j - 1)} = - \frac{\frac{dG}{d(\mu_j - 1)}}{\frac{dG}{dw}} > 0$$

Now we can write the other variables in terms of  $w$  and differentiate with respect to  $\mu_j - 1$ .

$$\frac{dq_E^*}{d(\mu_j - 1)} = \left[ \left( \frac{1}{\alpha_L} + \frac{1}{\alpha_H} \right) - w^* \left[ \frac{a_L}{\alpha_H} + \left( \frac{a_H}{\alpha_L} + 1 \right) (2 + t) \right] \right] \left[ 1 - \frac{\frac{B}{p_X(1-\beta)}}{\frac{dG}{dw}} \right] > 0$$



$$\begin{aligned}
\frac{dY^*}{d(\mu_j - 1)} &= \frac{P}{a_L \alpha_H w^{*2}} \frac{dw^*}{d(\mu_j - 1)} > 0 \\
\frac{dH^*}{d(\mu_j - 1)} &= \frac{P}{a_L w^{*2}} \frac{dw^*}{d(\mu_j - 1)} > 0 \\
\frac{dq_H^*}{d(\mu_j - 1)} &= -\frac{a_L}{\alpha_H} \frac{dw^*}{d(\mu_j - 1)} < 0 \\
\frac{dr_F^*}{d(\mu_j - 1)} &= -\left(\frac{1}{p_X(1-\beta)}\right)^{-\frac{1}{\beta}} q_E^{*\frac{-1}{\beta}} \frac{dq_E^*}{d(\mu_j - 1)} < 0 \\
\frac{dX^*}{d(\mu_j - 1)} &= F\left(1 - \frac{1}{\beta}\right) \left(\frac{1}{p_X(1-\beta)}\right)^{\frac{\beta-1}{\beta}} q_E^{*\frac{-1}{\beta}} \frac{dq_E^*}{d(\mu_j - 1)} < 0 \\
\frac{dE^*}{d(\mu_j - 1)} &= -\frac{P}{w^{*2}} \left(\frac{1}{\alpha_L} + \frac{1}{\alpha_H}\right) \frac{dw^*}{d(\mu_j - 1)} < 0 \\
\frac{d(ED^*)}{d(\mu_j - 1)} &= -\frac{P}{a_H w^{*2}} \frac{dw^*}{d(\mu_j - 1)} < 0
\end{aligned}$$

b) To be in this equilibrium, the demand for English skilled labor from  $X$  production must be equal to the supply; if there was excess supply,  $q_E^*$  would fall to increase firm profits and if there was excess demand,  $q_E^*$  would rise to attract additional English workers.

$$F\left(\frac{q_H}{p_X(1-\beta)}\right)^{-\frac{1}{\beta}} > F\left(\frac{q_E}{p_X(1-\beta)}\right)^{-\frac{1}{\beta}} = P \frac{wt + q_E - q_H}{w(\mu_j - 1)} > P \frac{t}{(\mu_j - 1)}$$

### 11.3 Proof of Proposition 3

a) First, I prove that if  $q_E^* = q_H^*$ , condition (10) holds. From Lemma 1, part c, we know

$$\frac{F}{P} \left(\frac{q_H^*}{p_X(1-\beta)}\right)^{-\frac{1}{\beta}} = \frac{F}{P} \left(\frac{1 - \alpha_L w^*}{\alpha_H p_X(1-\beta)}\right)^{-\frac{1}{\beta}} \leq \frac{t}{\mu_j - 1}$$

From the proof in Lemma 1 part a, equation (13), we can write

$$\frac{F}{P} \left(\frac{1 - \alpha_L w^*}{\alpha_H p_X(1-\beta)}\right)^{-\frac{1}{\beta}} = \frac{1}{w^*} \left(\frac{1}{\alpha_H} + \frac{1}{\alpha_L}\right) - (2+t) \left(1 + \frac{\alpha_H}{a_L}\right) - \frac{\alpha_L}{\alpha_H} \leq \frac{t}{\mu_j - 1}$$

Solving for  $w^*$ , we get

$$w^* \geq \frac{\left(\frac{1}{\alpha_H} + \frac{1}{\alpha_L}\right)}{(2+t) \left(1 + \frac{\alpha_H}{\alpha_L}\right) + \frac{\alpha_L}{\alpha_H} + \frac{t}{\mu_j - 1}}$$

Rewriting (10) and plugging in this inequality for  $w^*$ , we get

$$\begin{aligned} F &\leq P \frac{t}{\mu_j - 1} \left( \frac{1 - \alpha_L w^*}{\alpha_H p_X (1 - \beta)} \right)^{\frac{1}{\beta}} \leq \\ &P \frac{t}{\mu_j - 1} \left[ \frac{\alpha_L}{\alpha_H} \frac{1}{p_X (1 - \beta)} \left( \frac{1}{\alpha_L} - \frac{\left(\frac{1}{\alpha_H} + \frac{1}{\alpha_L}\right)}{\left[(2+t) \left(1 + \frac{\alpha_H}{\alpha_L}\right) + \frac{\alpha_L}{\alpha_H}\right] + \frac{t}{\mu_j - 1}} \right) \right]^{\frac{1}{\beta}} \end{aligned}$$

b) Next I prove that if  $F \leq \bar{F}(\mu_j)$ , then  $q_E^* = q_H^*$  by contradiction. We know that  $q_E^* \geq q_H^*$  since English skilled workers can take jobs as Hindi skilled workers. Suppose  $q_E^* > q_H^*$ . From condition (8), we know that

$$\begin{aligned} F &= \left( \frac{q_E}{p_X (1 - \beta)} \right)^{\frac{1}{\beta}} P \left( \frac{t}{\mu_j - 1} + \frac{1}{(\mu_j - 1)} \frac{q_E - q_H}{w} \right) \\ &> \left( \frac{q_E}{p_X (1 - \beta)} \right)^{\frac{1}{\beta}} P \left( \frac{t}{\mu_j - 1} \right) > P \frac{t}{\mu_j - 1} \left( \frac{1}{p_X (1 - \beta)} \right)^{\frac{1}{\beta}} \left( \frac{1 - \alpha_L w^*}{\alpha_H} \right)^{\frac{1}{\beta}} \end{aligned}$$

where the first inequality is due to  $q_E^* - q_H^* > 0$  and the second is due to  $q_E^* > q_H^* = \frac{1 - \alpha_L w^*}{\alpha_H}$ .

Putting this together with  $F \leq \bar{F}(\mu_j)$  and rearranging terms, we can show that

$$\frac{1}{w^*} \left( \frac{1}{\alpha_L} + \frac{1}{\alpha_H} \right) < \left[ (2+t) \left( 1 + \frac{\alpha_H}{\alpha_L} \right) + \frac{\alpha_L}{\alpha_H} \right] + \frac{t}{\mu_j - 1}$$

However, this contradicts what we know from the proof of Lemma 2, part a, equation (14)

$$\frac{1}{w^*} \left( \frac{1}{\alpha_L} + \frac{1}{\alpha_H} \right) - \left[ \frac{\alpha_L}{\alpha_H} + \left( \frac{\alpha_H}{\alpha_L} + 1 \right) (2+t) \right] = \frac{F}{P} \left( \frac{A - w^* B}{p_X (1 - \beta)} \right)^{-1/\beta} = \frac{wt + q_E - q_H}{w (\mu_j - 1)} > \frac{t}{(\mu_j - 1)}$$

where the second equality is from the fact that in this equilibrium,  $D_{EX} = S_E$  and the inequality is from  $q_E^* - q_H^* > 0$ . Thus,  $q_E^* = q_H^*$ .

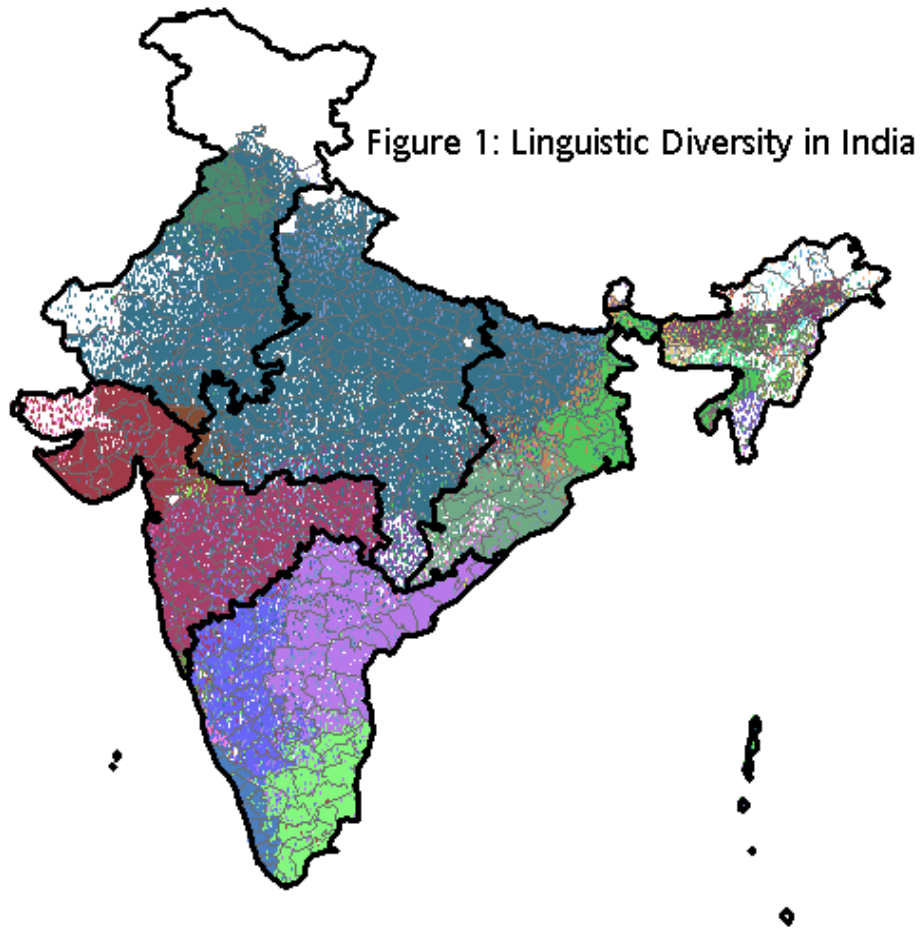


Figure 1: Linguistic Diversity in India

Note: A single dot represents 5000 individuals and different colors represent different languages.

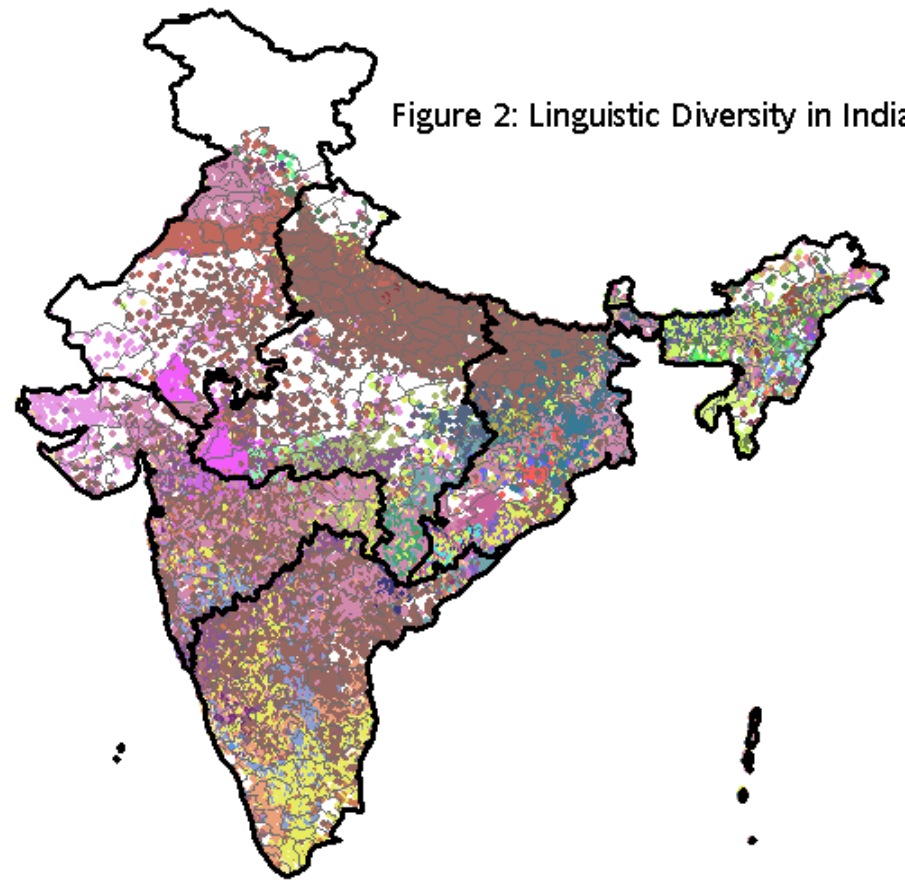


Figure 2: Linguistic Diversity in India

Note: A single dot represents 3000 individuals and different colors represent different languages, excluding the most spoken language in the district.

Figure 3: Chart of Language Distances

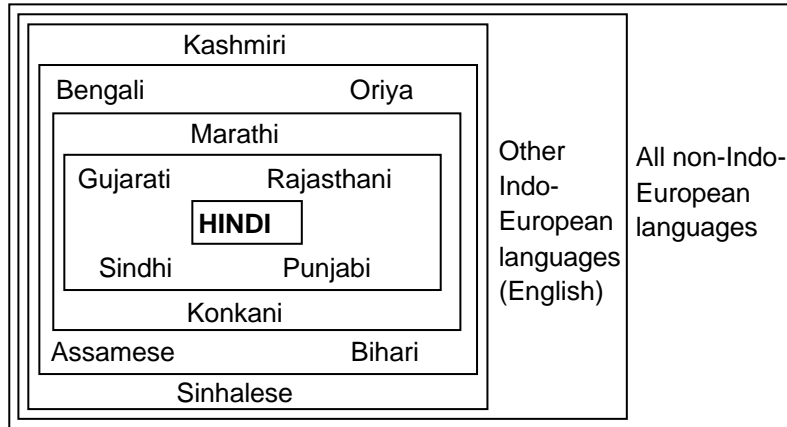


Figure 4: Map of Linguistic Distance in India

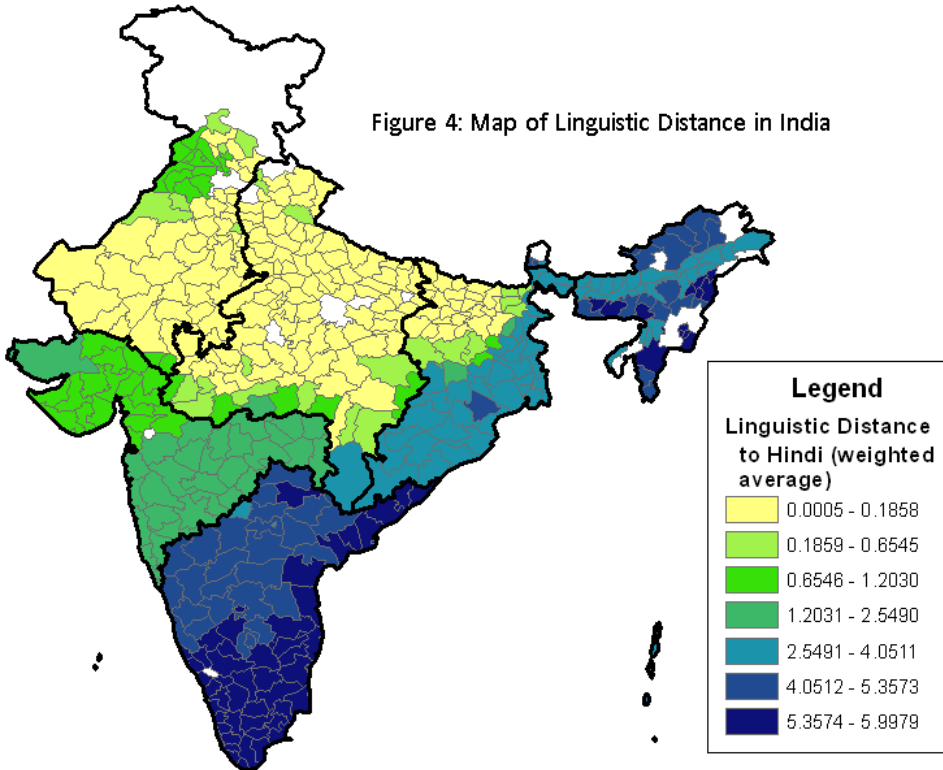
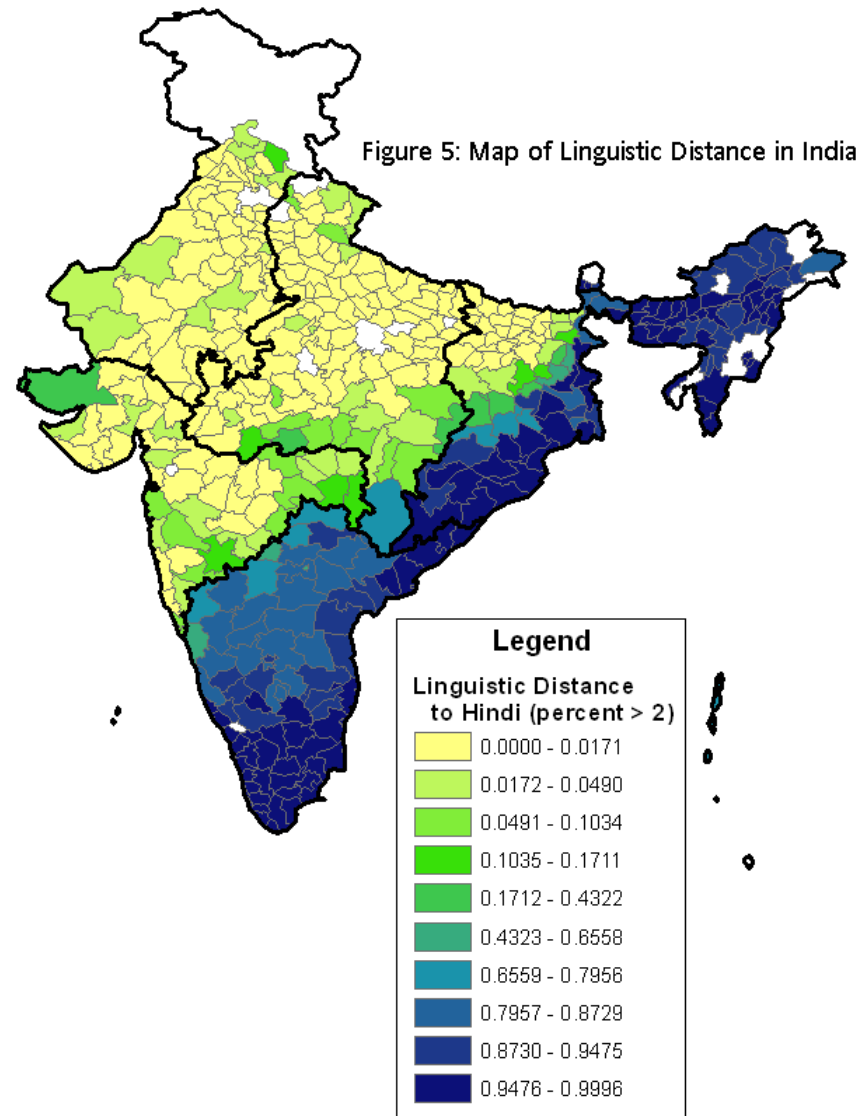
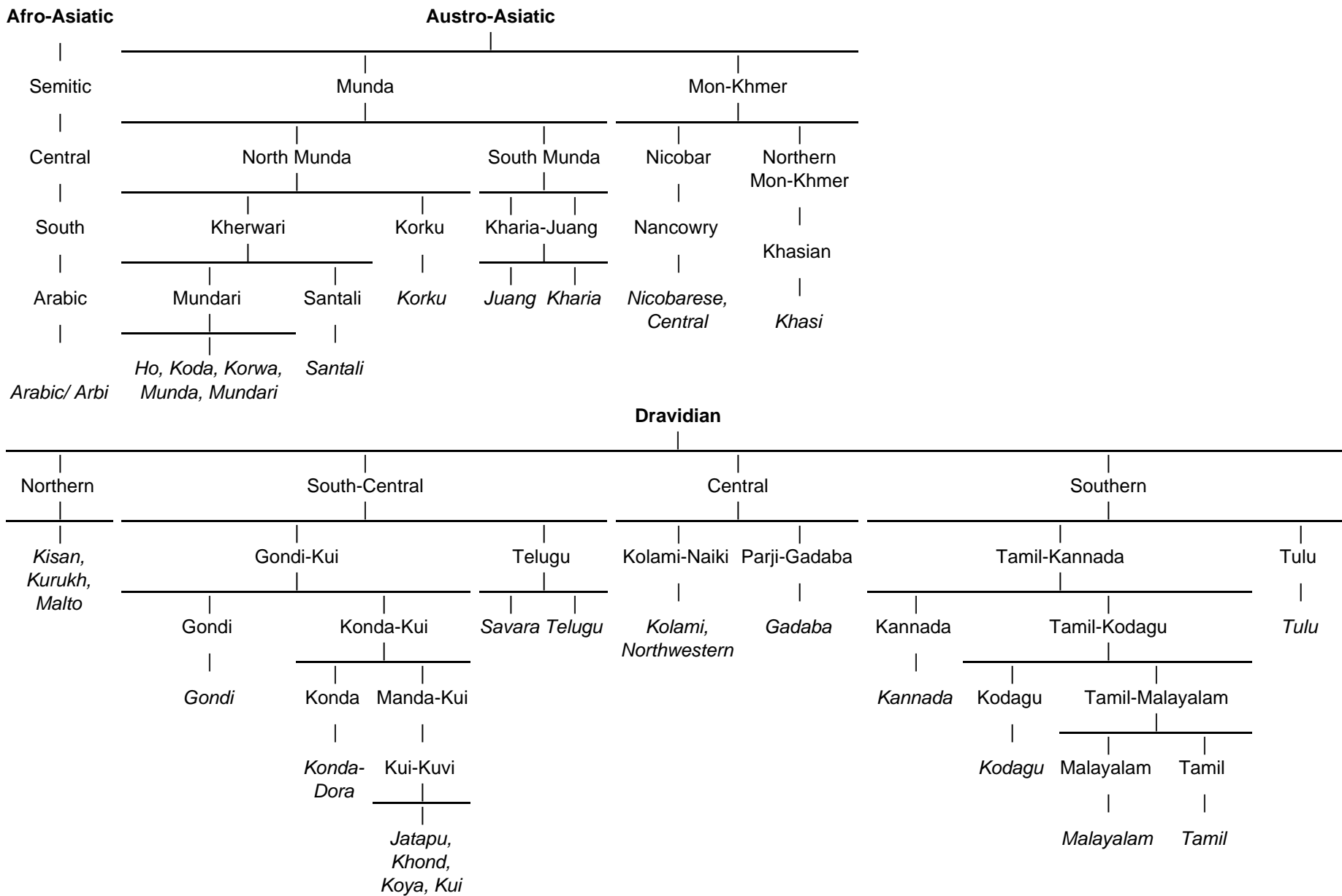


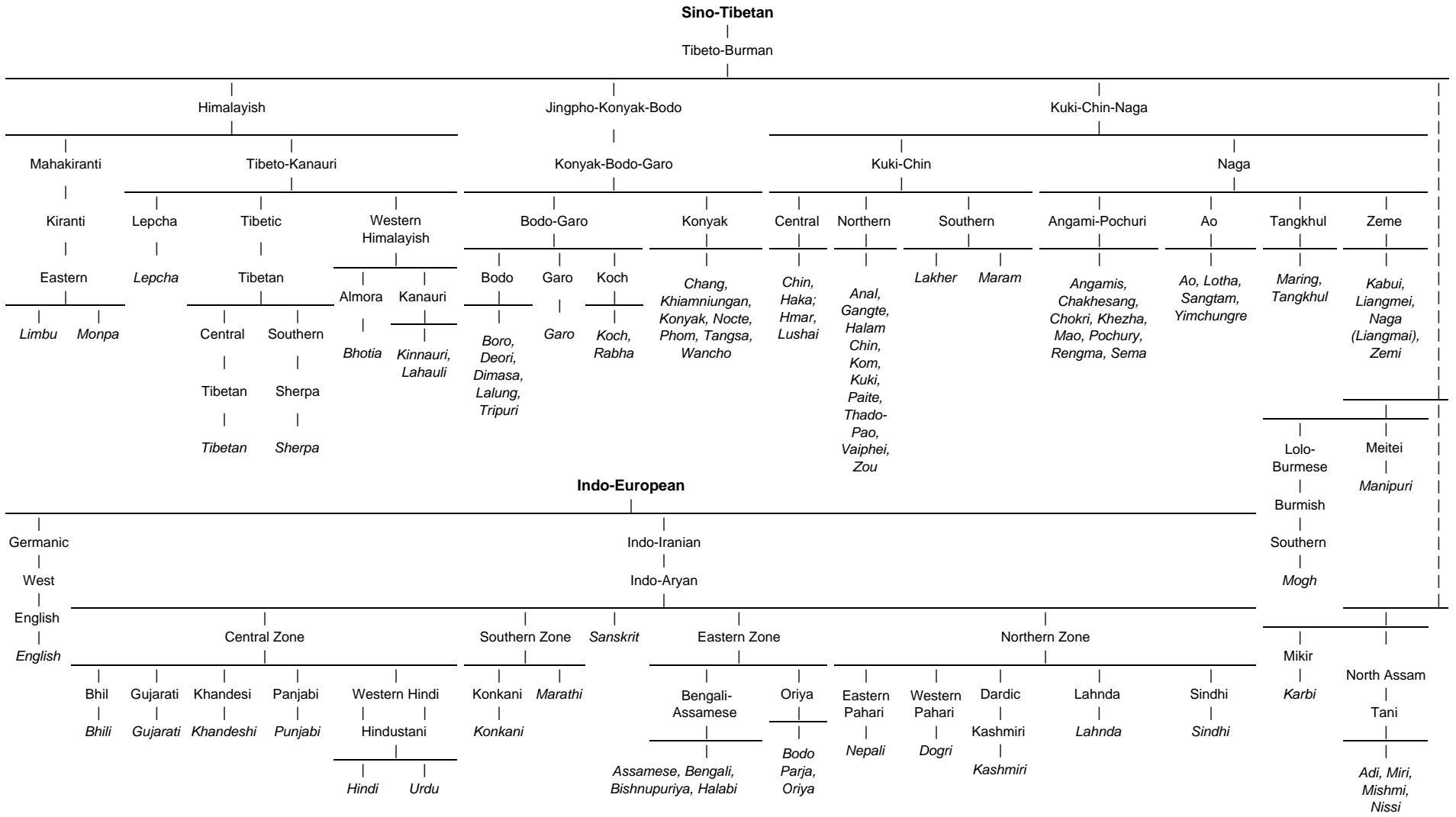
Figure 5: Map of Linguistic Distance in India



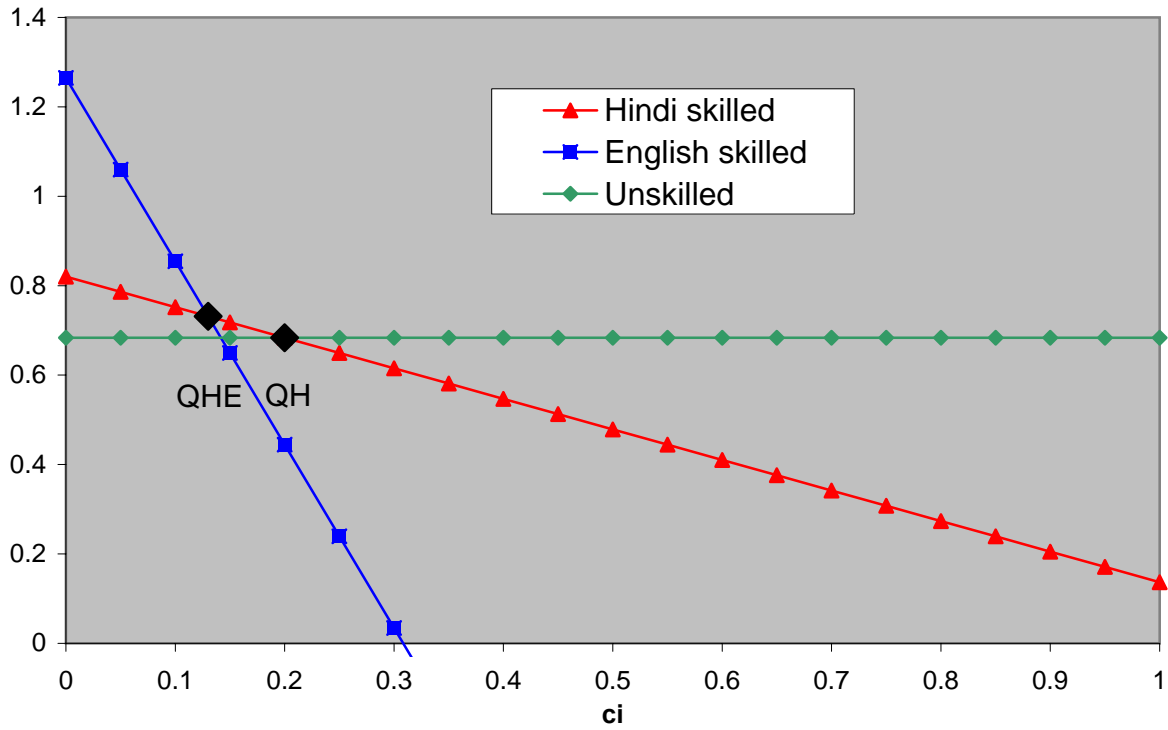
**Figure 6: Language trees showing the relative position of languages in India**



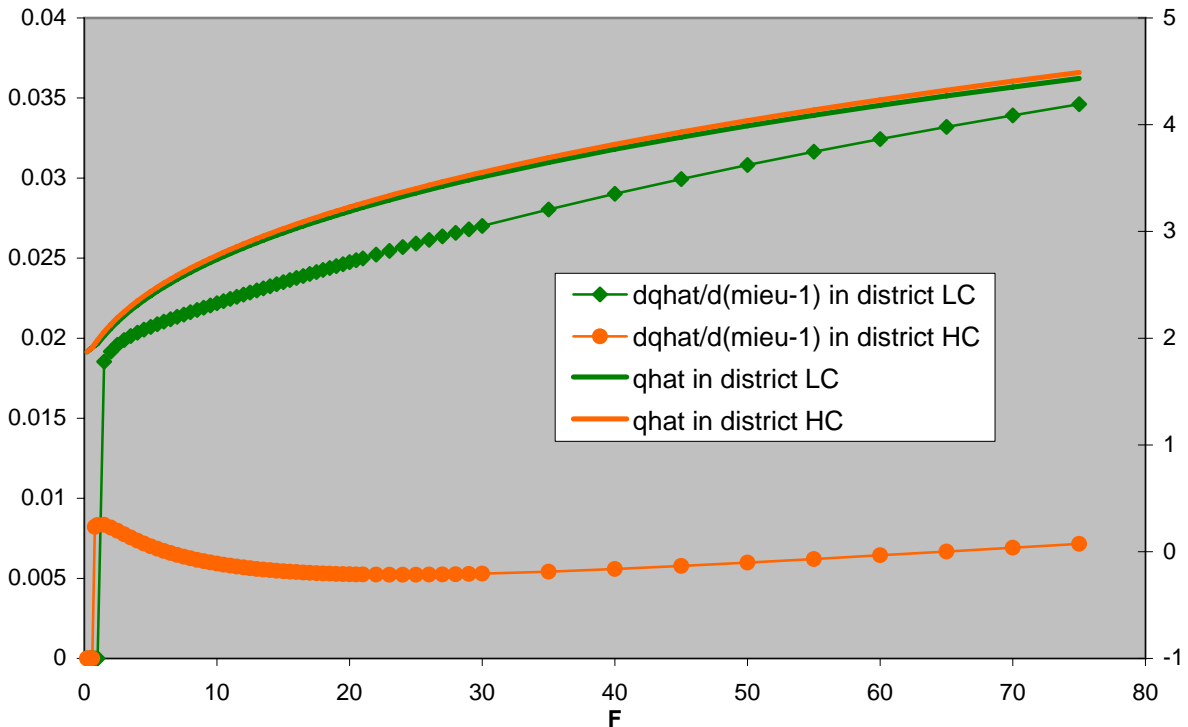
**Figure 6 (continued): Language trees showing the relative position of languages in India**



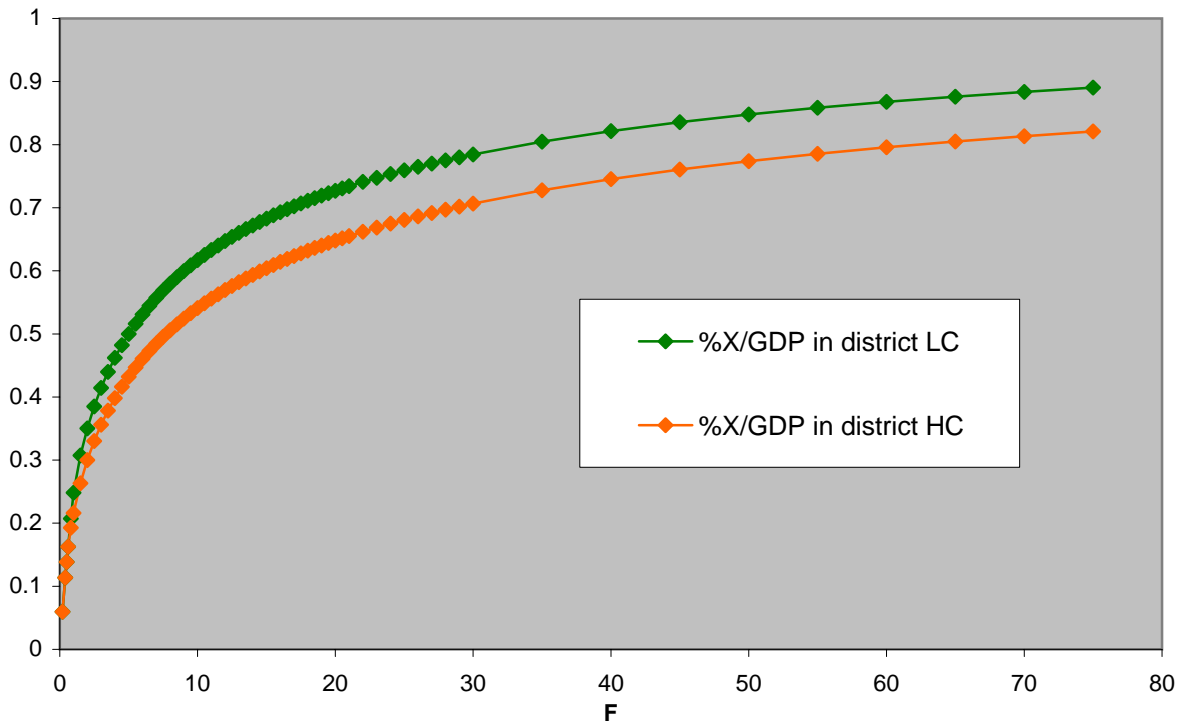
**Figure 7: Schooling Decisions**



**Figure 8: Weighted average return to wages and the differential with respect to  $\mu$ ,  $\beta=0.3$**



**Figure 9: X production as a share of GDP in both districts,  $\beta=0.3$**



**Figure 10: Weighted average return to wages and the differential with respect to  $\mu$ ,  $\beta=0.18$**

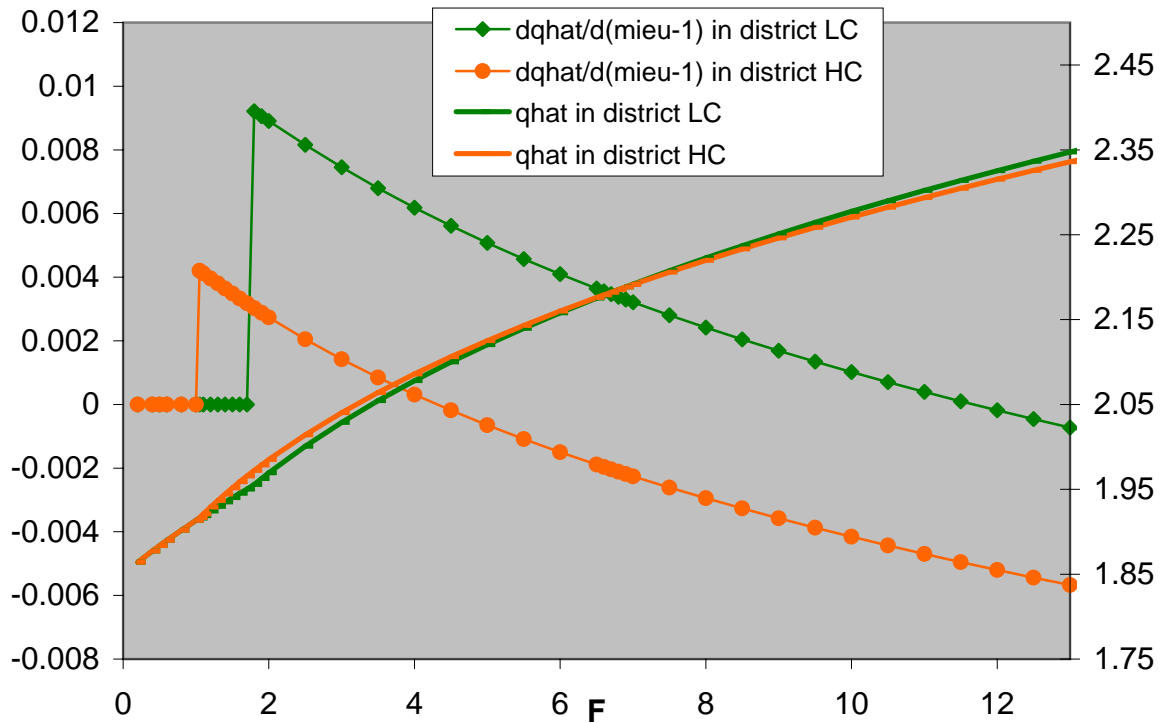




Table 1: Translations and cognate judgments for sample words in English, Bengali and Hindi

English	Meanings		Cognate Judgment		
	Hindi	Bengali	Hindi - Bengali	Hindi - English	Bengali - English
ALL	SEB, SARA	SOB	Yes	No	No
AND	OR	AR, EBON	Yes	No	No
ANIMAL	JANVER	JANOAR, JONTU	Yes, doubtful	No	No
BAD	KHERAB, BURA	KHARAP	Yes	No	No
CLOUD	BADEL	MEG	No	No	No
EYE	AKH	COK	No	Yes	No
FEATHER	PER	PALOK	No	Yes, doubtful	No
FIVE	PAC	PAC	Yes	Yes	Yes
FOOT	PER	PA	Yes	Yes	Yes
FOUR	CAR	CAR	Yes	Yes	Yes
FRUIT	PHEL	PHOL	Yes	No	No
GOOD	ECCHA	BHALO	No	No	No
GRASS	GHAS	GHAS	Yes	No	No
HOW	KESA	KEMON	Yes	Yes	Yes
I	ME	AMI	Yes	Yes	Yes
IN	ENDER, -ME	ONDOR	Yes	Yes	Yes
MOTHER	MA	MA	Yes	Yes	Yes
NAME	NAM	NAM	Yes	Yes	Yes
NOSE	NAK	NAK	Yes	Yes	Yes
OTHER	DUSRA	ONNO	No	No	Yes
STAR	TARA	TARA	Yes	Yes	Yes
TO COME	ANA	ASA	No	No	No
TO FREEZE	JEMNA	JOMAT+BADHANO	Yes	No	No
WITH	SATH	SATH, SONGE	Yes	No	No
<b>Percent Cognates</b>			<b>64.10%</b>	<b>14.60%</b>	<b>14.20%</b>

Table 2: Summary Statistics

Variable abbreviated name	Notes	Num Obs	Mean	St. Dev.	Min.	Max.
<b>Panel A (at the district level)</b>						
Linguistic distance (weighted average)	Degree measure of distance from native languages to Hindi	390	2.158	2.281	0.001	5.998
Percent with linguistic distance > 2	Percent of people who speak languages at distance > 2	390	0.373	0.443	0.000	1.000
Linguistic distance measure 2	Node measure of distance from native languages to Hindi	390	5.063	4.778	0.021	13.994
Linguistic distance measure 3	Cognate measure of distance from native languages to Hindi	390	64.814	35.222	5.031	99.991
Native Hindi speakers	Percent of people who speak Hindi (as a native language)	390	0.424	0.433	0.000	0.996
Speakers at distance 0	Percent of people who speak languages at distance 0	390	0.465	0.443	0.000	1.000
Speakers at distance 1	Percent of people who speak languages at distance 1	390	0.095	0.265	0.000	0.991
Speakers at distance 2	Percent of people who speak languages at distance 2	390	0.067	0.212	0.000	0.958
Speakers at distance 3	Percent of people who speak languages at distance 3	390	0.094	0.251	0.000	0.985
Speakers at distance 4	Percent of people who speak languages at distance 4	390	0.013	0.069	0.000	0.766
Native English Speakers	Percent of people who speak languages at distance 5	390	0.000	0.000	0.000	0.007
Speakers at distance 6	Percent of people who speak languages at distance 6	390	0.265	0.391	0.000	1.000
Primary schools teaching in mother tongue	Percent of urban primary schools that teach in the mother tongue	408	0.889	0.222	0	1.08
Upper primary schools teaching in mother tongue	Percent of urban upper primary schools that teach in the mother tongue	408	0.840	0.245	0	1.02
<b>Panel B (at the state level, only urban areas)</b>						
Primary schools teaching English	Percent of primary schools that teach English	32	0.263	0.164	0.023	0.664
Upper primary schools teaching English	% upper primary schools that teach English	32	0.310	0.070	0.233	0.511
Secondary schools teaching English	% secondary schools that teach English	32	0.337	0.104	0.182	0.615
Primary schools teaching in English	% primary schools with English instruction	32	0.222	0.237	0	1.00
Upper primary schools teaching in English	% upper primary schools with English instruction	32	0.321	0.268	0.053	1.00
Secondary schools teaching in English	% secondary schools with English instruction	32	0.380	0.285	0.047	1.00
Higher secondary schools teaching in English	% higher secondary schools with English instruction	31	0.470	0.295	0.051	1.00
<b>Panel C (at the district level, only urban areas)</b>						
Household wage income	Average weekly total wage income in 1000s of Rupees	397	0.183	0.178	0	1.763
Educated wage	Average wage of individuals with at least high school education	396	0.233	0.197	0	2.714
Salaried	Percent of adults with a regular wage or salaried job	397	0.188	0.082	0	0.568
Graduate	Percent of people with a college degree	397	0.051	0.043	0	0.5
Secondary	Percent of people with a high school degree	397	0.125	0.057	0	0.3333333
Literate	Percent of people who are literate	397	0.625	0.137	0.1111111	0.969697
Muslim	Percent of people who are Muslim	397	0.168	0.193	0	1
Train	Percent of people who recently made a journey by train	395	0.077	0.094	0	0.686
Electricity	% households that use electricity as their main energy source	395	0.695	0.203	0	1
Hindi belt	Districts that are in the Hindi belt	409	0.477	0.500	0	1
Child population 1991	District population of 5-19 year olds in 1991	379	181846	278514	631	2952148
Child population 2001	District population of 5-19 year olds in 2001	379	181846	278514	631	2952148
Child population growth	Population growth rate for 5-18 year olds	379	0.318	0.347	-0.422	2.695
Distance to closest big city	Distance to closest of the 10 biggest cities in India	401	31.495	18.582	0.206	120.363
Distance to closest airport	Distance to closest airport operated by Airport Authority of India	401	7.797	4.695	0.443	24.761
Number of engineering colleges	Number of IITs and NITs	409	0.064	0.244	0	1
Engineers	% engineers among non-migrant 26-65-year-olds	395	0.002	0.008	0	0.065

Table 3: Impact of Linguistic Distance on % of Native Speakers who Learn English

Dependent Variable:	% of Multilinguals who Learn English			% of Native Speakers who are Multilingual		% of Multilinguals who Learn Hindi, but not English			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Linguistic distance (0 - 6)	0.077 *** (0.025)			0.030 ** (0.015)			-0.060 ** (0.026)		
Linguistic distance > 2		0.216 *** (0.055)			0.007 (0.036)			-0.229 *** (0.079)	
Linguistic distance = 1			-0.095 (0.063)			0.111 ** (0.047)			0.055 (0.042)
Linguistic distance = 2			-0.004 (0.070)			0.194 *** (0.038)			-0.005 (0.036)
Linguistic distance = 3			0.171 ** (0.068)			0.111 *** (0.036)			-0.271 *** (0.071)
Linguistic distance = 4			0.203 ** (0.084)			0.117 (0.098)			0.085 (0.079)
Language = Hindi	5.950 (4.212)	3.455 (3.664)	3.469 (3.699)	2.993 (3.135)	0.435 (2.932)	4.577 (3.268)			
Most spoken language in state	0.153 (0.123)	0.215 ** (0.105)	0.200 * (0.109)	0.082 (0.101)	0.136 (0.110)	0.026 (0.100)	0.052 (0.097)	0.029 (0.087)	0.073 (0.089)
Share of native speakers in state	0.156 (0.377)	0.172 (0.362)	0.069 (0.377)	-1.037 *** (0.286)	-1.137 *** (0.320)	-1.036 *** (0.273)	-0.264 (0.272)	-0.336 (0.255)	-0.286 (0.253)
Share of native speakers in state squared	0.263 (0.330)	0.186 (0.354)	0.345 (0.356)	0.612 *** (0.202)	0.673 *** (0.227)	0.659 *** (0.187)	0.072 (0.201)	0.177 (0.196)	0.071 (0.179)
Share of native speakers in India	6.057 ** (3.039)	4.068 (2.632)	3.616 (2.907)	1.127 (2.815)	-0.700 (2.645)	2.319 (2.970)	-2.340 (2.324)	-1.408 (2.468)	0.800 (2.612)
Share of native speakers in India squared	-49.112 (32.910)	-29.255 (28.681)	-28.646 (29.371)	-21.218 (25.812)	-1.394 (24.211)	-33.546 (27.084)	26.191 (21.351)	16.070 (23.104)	0.237 (23.525)
Lang. Family = Afro-Asiatic	-0.227 ** (0.096)	-0.053 (0.047)	0.110 (0.079)	-0.120 (0.076)	-0.020 (0.073)	0.110 (0.082)	0.240 (0.146)	0.127 (0.084)	-0.028 (0.058)
Lang. Family = Austro-Asiatic	-0.365 *** (0.131)	-0.161 ** (0.080)	-0.001 (0.089)	-0.175 ** (0.074)	-0.079 (0.068)	0.033 (0.078)	0.335 (0.212)	0.169 (0.158)	-0.027 (0.119)
Lang. Family = Dravidian	-0.262 ** (0.118)	-0.075 * (0.045)	0.109 *** (0.038)	-0.196 *** (0.067)	-0.073 ** (0.031)	0.034 (0.035)	0.274 * (0.160)	0.168 * (0.099)	-0.050 (0.047)
Lang. Family = Sino-Tibetan	-0.055 (0.097)	0.162 ** (0.075)	0.311 *** (0.118)	-0.165 *** (0.062)	-0.085 (0.058)	0.034 (0.073)	0.112 (0.099)	-0.066 (0.071)	-0.259 *** (0.094)
No. of Obs. (weighted)	8.4E+08	8.4E+08	8.4E+08	8.4E+08	8.4E+08	8.4E+08	5E+08	5E+08	5E+08
No. of Obs.	1466	1466	1466	1741	1741	1741	1435	1435	1435
R-squared	0.914	0.918	0.921	0.798	0.791	0.812	0.823	0.832	0.847
p-value of language distance measures			0.00			0.00			0.00

Robust standard errors, clustered by state are shown in parentheses. All columns include region fixed effects.

Table 4: Impact of Linguistic Distance on % of Native Speakers who Learn English in 1961

Dependent Variable:	% of Multilinguals who Learn English			% of Native Speakers who are Multilingual			% of Multilinguals who Learn Hindi		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Linguistic distance (0 - 6)	0.184 *** (0.044)			0.021 (0.013)			-0.067 * (0.039)		
Linguistic distance > 2		0.360 *** (0.083)			0.032 (0.024)			-0.254 ** (0.113)	
Linguistic distance = 2			0.156 * (0.091)			0.039 * (0.021)			-0.022 (0.049)
Linguistic distance = 3			0.476 *** (0.120)			0.063 * (0.035)			-0.348 ** (0.150)
Linguistic distance = 4			0.269 ** (0.130)			-0.005 (0.068)			-0.116 (0.144)
Language = Hindi	14.180 ** (5.975)	9.763 (6.227)	9.893 (6.344)	3.297 * (1.749)	2.825 * (1.698)	2.715 (2.173)			
Most spoken language in state	-0.096 (0.167)	-0.036 (0.139)	-0.130 (0.154)	-0.026 (0.033)	-0.015 (0.033)	-0.039 (0.034)	0.278 ** (0.133)	0.275 ** (0.117)	0.291 ** (0.113)
Share of native speakers in state	-1.725 *** (0.544)	-1.707 *** (0.611)	-1.650 *** (0.575)	-0.861 *** (0.143)	-0.872 *** (0.144)	-0.859 *** (0.148)	-0.010 (0.406)	-0.105 (0.403)	-0.068 (0.391)
Share of native speakers in state squared	2.452 *** (0.461)	2.358 *** (0.544)	2.424 *** (0.490)	0.464 *** (0.100)	0.461 *** (0.099)	0.479 *** (0.105)	-0.196 (0.259)	-0.073 (0.281)	-0.127 (0.263)
Share of native speakers in India	11.625 *** (4.325)	9.113 ** (4.238)	8.003 * (4.607)	3.380 *** (1.303)	3.139 ** (1.311)	2.732 * (1.623)	-9.618 ** (4.428)	-8.928 ** (4.508)	-7.309 (4.796)
Share of native speakers in India squared	-100.656 ** (40.787)	-72.397 * (42.457)	-70.334 (43.606)	-24.677 ** (12.107)	-21.756 * (11.920)	-20.178 (15.173)	105.347 ** (44.299)	97.164 ** (45.191)	84.708 * (48.930)
Lang. Family = Austro-Asiatic	-0.705 *** (0.177)	-0.184 ** (0.086)	0.228 ** (0.101)	-0.079 (0.060)	-0.018 (0.062)	0.026 (0.080)	0.380 (0.234)	0.199 (0.149)	-0.095 (0.130)
Lang. Family = Dravidian	-0.645 *** (0.153)	-0.158 ** (0.067)	0.298 *** (0.097)	-0.089 * (0.051)	-0.030 (0.034)	0.027 (0.033)	0.176 (0.201)	0.083 (0.125)	-0.196 *** (0.073)
Lang. Family = Sino-Tibetan	-0.649 *** (0.162)	-0.107 (0.129)	0.295 * (0.177)	0.006 (0.046)	0.070 *** (0.026)	0.111 *** (0.037)	0.206 (0.155)	0.007 (0.091)	-0.292 * (0.151)
Number of Observations	4.2E+08	4.2E+08	4.2E+08	4.2E+08	4.2E+08	4.2E+08	2.4E+08	2.4E+08	2.4E+08
R-squared	0.793	0.794	0.806	0.827	0.826	0.829	0.772	0.781	0.786
p-value of language distance measures			0.02			0.00			0.11

Robust standard errors, clustered by region are shown in parentheses. All columns include region fixed effects.

Table 5: Impact of Linguistic Distance on Percent of Schools that Teach English

Dependent Variable:	% of Schools in State Teaching In English			% of Schools in State Teaching English			% of Schools in District Teaching in Mother Tongue		
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Linguistic distance (0 - 6)	0.046 *** (0.010)			0.085 *** (0.025)			-0.034 *** (0.010)		
Linguistic distance > 2		0.157 *** (0.051)			0.406 *** (0.144)			-0.121 *** (0.042)	
Linguistic distance of 1			0.074 (0.862)			-0.677 (1.780)			0.235 (0.228)
Linguistic distance of 2			0.512 (0.984)			-0.383 (1.920)			0.282 (0.253)
Linguistic distance of 3			0.099 (0.847)			-0.955 (1.728)			0.282 (0.242)
Linguistic distance of 4			0.016 (0.894)			0.347 (1.836)			-0.506 * (0.283)
Linguistic distance of 6			0.303 (0.896)			-0.426 (1.855)			0.106 (0.236)
Native English speakers	27.639 *** (10.124)	33.392 *** (11.498)	-60.580 (45.040)	60.242 ** (27.205)	91.428 *** (26.720)	-0.405 (90.986)	-70.574 *** (24.972)	-76.546 *** (23.444)	-61.636 ** (26.770)
Native Hindi speakers	0.159 *** (0.039)	0.153 *** (0.050)	0.141 (0.853)	0.303 ** (0.137)	0.307 * (0.160)	-0.687 (1.770)	-0.003 (0.057)	0.031 (0.057)	0.293 (0.211)
Hindi belt states	-0.010 (0.064)	-0.048 (0.063)	0.015 (0.057)	-0.057 (0.155)	-0.065 (0.151)	0.004 (0.112)	-0.054 (0.062)	-0.039 (0.061)	-0.033 (0.074)
Child population in 1991	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 * (0.000)	0.000 (0.000)	0.000 *** (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
Household wage income	0.000 * (0.000)	0.000 (0.000)	0.000 *** (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	-0.044 (0.042)	-0.037 (0.041)	-0.022 (0.046)
Educated wage	0.000 (0.000)	-0.001 (0.000)	0.000 (0.000)	-0.001 ** (0.000)	-0.001 ** (0.001)	0.000 (0.001)	0.032 (0.143)	0.039 (0.147)	-0.012 (0.141)
Salaried	0.880 * (0.451)	0.499 (0.503)	1.691 *** (0.445)	2.336 ** (1.168)	1.587 (1.031)	2.095 ** (0.862)	-0.397 ** (0.193)	-0.383 ** (0.187)	-0.388 ** (0.191)
Distance to closest big city	0.001 (0.001)	0.001 (0.001)	-0.001 (0.001)	-0.001 (0.003)	0.000 (0.002)	0.001 (0.002)	-0.003 * (0.002)	-0.004 * (0.002)	-0.003 * (0.002)
Graduate	1.386 ** (0.660)	0.984 (0.838)	1.434 ** (0.560)	0.046 (2.241)	-0.695 (2.619)	2.837 ** (1.226)	0.172 (0.371)	0.313 (0.394)	0.138 (0.336)
Secondary	-0.918 * (0.558)	-0.857 (0.753)	-1.529 *** (0.575)	-0.468 (1.310)	-0.219 (1.458)	-0.940 (1.489)	0.147 (0.266)	0.121 (0.262)	0.133 (0.236)
Literate	0.173 (0.408)	0.438 (0.484)	0.160 (0.506)	-0.339 (0.538)	0.032 (0.656)	0.099 (1.025)	0.109 (0.193)	0.120 (0.191)	0.030 (0.179)
Muslim	-0.099 (0.084)	-0.046 (0.096)	-0.055 (0.084)	0.030 (0.156)	0.098 (0.180)	0.216 (0.192)	0.015 (0.068)	0.009 (0.068)	0.030 (0.053)
Train	-0.051 (0.228)	-0.294 (0.216)	-0.353 (0.458)	0.434 (0.591)	0.055 (0.553)	0.538 (0.907)	0.071 (0.107)	0.109 (0.095)	0.025 (0.107)
Electricity	0.173 (0.120)	0.296 (0.196)	-0.025 (0.152)	0.459 (0.346)	0.805 (0.499)	-0.483 (0.334)	-0.073 (0.105)	-0.111 (0.117)	-0.003 (0.088)
Number of Observations	90	90	90	119	119	119	732	732	732
R-squared	0.605	0.538	0.651	0.591	0.570	0.718	0.312	0.294	0.373
p-value of language distance measures			0.000			0.000			0.001

Robust standard errors, clustered by state, are shown in parentheses. All columns include fixed effects for school level.

Table 6: Impact of Linguistic Distance on Number of Schools

Dependent Variable:	Number of Schools		Number of HS Schools Offering ...					
			Arts		Commerce		Science	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Linguistic distance (weighted average)	0.096 (0.278)		-0.092 (0.065)		-0.001 (0.065)		0.071 (0.079)	
Linguistic distance > 2		0.526 (0.890)		-0.195 (0.236)		-0.100 (0.245)		0.349 (0.259)
Native English speakers	954 ** (449)	978 ** (452)	135 (136)	123 (132)	84 (92)	80 (88)	200 ** (99)	216 ** (95)
Native Hindi speakers	2.715 ** (1.276)	2.681 ** (1.217)	0.435 (0.303)	0.569 ** (0.257)	0.301 (0.467)	0.271 (0.461)	0.635 (0.438)	0.597 (0.412)
Hindi belt states	-1.306 (1.006)	-1.271 (0.971)	0.165 (0.372)	0.255 (0.358)	-0.130 (0.413)	-0.167 (0.416)	0.546 (0.488)	0.556 (0.482)
Child population in 1991	0.000 *** (0.000)	0.000 *** (0.000)	0.000 * (0.000)	0.000 * (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 * (0.000)	0.000 ** (0.000)
Household wage income	-0.779 (0.897)	-0.796 (0.896)	-1.038 *** (0.403)	-1.019 ** (0.403)	-0.446 (0.298)	-0.447 (0.297)	-0.917 * (0.528)	-0.930 * (0.529)
Educated wage	3.423 ** (1.493)	3.377 ** (1.494)	0.310 (0.435)	0.310 (0.441)	-0.201 (0.256)	-0.187 (0.256)	-0.146 (0.442)	-0.175 (0.436)
Salaried	1.202 (4.288)	1.142 (4.303)	-0.490 (0.885)	-0.459 (0.921)	-0.604 (0.687)	-0.596 (0.695)	0.779 (0.899)	0.738 (0.893)
Distance to closest big city	0.024 (0.018)	0.023 (0.018)	0.001 (0.006)	0.000 (0.006)	-0.008 ** (0.004)	-0.007 ** (0.004)	0.001 (0.006)	0.000 (0.006)
Graduate	14.715 * (8.213)	14.226 * (8.023)	3.554 ** (1.728)	3.889 ** (1.740)	0.778 (1.406)	0.825 (1.370)	3.341 * (1.875)	2.999 (1.836)
Secondary	1.170 (4.897)	1.231 (4.927)	3.233 (2.112)	3.179 (2.161)	0.096 (1.142)	0.094 (1.141)	-0.084 (2.159)	-0.040 (2.153)
Literate	-50.154 ** (21.698)	-49.970 ** (22.536)	-12.713 ** (5.803)	-12.043 ** (5.747)	-6.381 ** (3.028)	-6.638 ** (3.114)	-13.978 ** (6.128)	-13.946 ** (6.010)
Muslim	-1.643 (1.297)	-1.604 (1.295)	-0.529 (0.428)	-0.540 (0.429)	-0.809 *** (0.276)	-0.818 *** (0.273)	-0.886 * (0.512)	-0.860 * (0.518)
Train	-3.797 * (2.212)	-3.888 * (2.220)	-1.422 * (0.820)	-1.307 (0.815)	-0.403 (0.412)	-0.411 (0.420)	-0.998 (0.819)	-1.068 (0.797)
Electricity	1.571 (1.217)	1.747 (1.297)	0.859 (0.559)	0.818 (0.581)	1.448 *** (0.259)	1.407 *** (0.262)	1.566 *** (0.521)	1.679 *** (0.499)
Number of Observations	1464	1464	366	366	366	366	366	366
R-squared	0.630	0.630	0.288	0.281	0.243	0.244	0.264	0.264

Robust standard errors, clustered by state, are shown in parentheses. All columns include fixed effects for school level.

Table 7: Impact of Linguistic Distance on Economic Variables

Dependent Variable:	Child Population in 1991 (in 10,000s)	Child population growth (1991-2001)	Household wage income	Percent Salaried	Percent with College Degree	Percent Literate	Percent Engineers	Percent Muslim	Distance to closest big city	Distance to Closest Airport
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Linguistic distance (weighted average)	0.388 (1.176)	-0.007 (0.016)	-0.005 (0.003)	-0.001 (0.003)	-0.003 ** (0.001)	0.011 (0.007)	0.000 (0.000)	0.004 (0.010)	3.001 ** (1.465)	-0.172 (0.293)
Native English speakers	28621 *** (8752)	-21.235 (21.388)	0.631 (7.702)	-0.726 (6.102)	0.778 (1.679)	11.693 (7.687)	2.770 ** (1.230)	-16.659 (10.453)	1279.655 (3004.179)	52.591 (587.689)
Native Hindi speakers	-2.497 (4.660)	0.007 (0.053)	-0.053 (0.039)	0.000 (0.023)	0.023 *** (0.006)	-0.013 (0.034)	-0.004 * (0.002)	0.059 (0.067)	8.350 (8.713)	-1.627 (2.063)
Hindi belt states	-1.592 (3.678)	0.043 (0.033)	0.065 ** (0.029)	-0.016 (0.022)	-0.005 (0.009)	-0.039 (0.037)	0.002 (0.003)	-0.081 (0.053)	1.755 (6.334)	1.016 (2.314)
Child population in 1991		0.000 (0.000)	0.000 (0.000)	0.000 *** (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	-0.227 *** (0.058)	-0.044 *** (0.012)
Household wage income	2.075 (5.148)	0.044 (0.044)		0.039 (0.043)	-0.008 (0.009)	0.034 * (0.020)	0.002 (0.003)	0.003 (0.042)	-4.189 (4.891)	-0.601 (1.174)
Educated wage	2.484 (6.321)	0.012 (0.066)	0.617 *** (0.094)	0.104 *** (0.038)	0.022 * (0.013)	0.034 (0.031)	0.000 (0.002)	0.045 (0.051)	7.364 (7.981)	0.713 (1.362)
Salaried	49.844 *** (16.229)	0.756 *** (0.236)	0.191 (0.162)		0.058 * (0.030)	0.095 (0.097)	0.006 (0.011)	-0.443 *** (0.146)	20.561 (20.862)	7.306 ** (3.232)
Distance to closest big city	-0.377 *** (0.092)	0.001 (0.001)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 * (0.000)	-0.001 (0.001)		-0.022 (0.021)
Graduate	34.056 (37.894)	0.164 (0.679)	-0.189 (0.228)	0.262 ** (0.116)		0.575 *** (0.158)	0.046 ** (0.019)	-0.175 (0.266)	-15.437 (39.529)	-15.395 ** (6.984)
Secondary	6.334 (26.307)	0.367 (0.348)	0.317 (0.208)	0.129 (0.101)	0.023 (0.038)	0.739 *** (0.130)	0.007 (0.009)	-0.111 (0.166)	21.331 (32.663)	-11.752 ** (5.702)
Literate	-16.185 (10.776)	-0.283 (0.208)	0.106 (0.071)	0.061 (0.064)	0.082 *** (0.021)		-0.001 (0.004)	-0.486 *** (0.172)	10.213 (18.021)	-1.476 (3.631)
Muslim	8.655 * (4.441)	0.026 (0.083)	0.003 (0.038)	-0.083 *** (0.026)	-0.007 (0.011)	-0.142 *** (0.039)	0.000 (0.003)		-11.681 (7.744)	-1.393 (1.788)
Train	36.541 * (20.471)	-0.053 (0.137)	0.093 (0.121)	-0.057 * (0.031)	-0.001 (0.015)	0.031 (0.030)	-0.005 ** (0.003)	0.077 (0.114)	-13.839 (10.620)	1.978 (3.139)
Electricity	0.675 (4.491)	0.002 (0.076)	-0.017 (0.034)	0.050 ** (0.024)	0.016 * (0.008)	0.181 *** (0.038)	-0.006 (0.003)	0.131 ** (0.064)	-9.980 (12.743)	-3.680 * (2.111)
Number of Observations	367	367	367	367	367	367	366	367	367	367
R-squared	0.476	0.159	0.398	0.396	0.292	0.675	0.104	0.241	0.251	0.163

Robust standard errors, clustered by state, are shown in parentheses.

Table 8: Summary Statistics on IT Presence Across Districts

Year of Data	Average Across Districts					
	Number of Districts with IT	Number of HQ or Branches	Employees	Revenue (in millions Rs.)	Exports (in millions Rs.)	Subscribed Capital (in millions Rs.)
1995	47	1.36	120	56.71	39.92	21.32
1998	47	2.25	279			54.82
1999	54	2.48	348	277.42	232.31	48.54
2002	76	3.24	559	470.37	344.63	118.55
2003	72	2.54	604	587.34	508.96	111.59



Table 9: Summary Statistics on Employment, only urban areas

Variable name	1987					1999				
	Num Obs	Mean	St. Dev.	Min.	Max.	Num Obs	Mean	St. Dev.	Min.	Max.
<b>Salaried Employment</b>	68567	0.412	0.492	0.000	1.000	72993	0.389	0.488	0.000	1.000
if male	56654	0.423	0.494	0.000	1.000	60008	0.393	0.488	0.000	1.000
if female	11913	0.360	0.480	0.000	1.000	12985	0.374	0.484	0.000	1.000
if under age 30	24818	0.333	0.471	0.000	1.000	23880	0.318	0.466	0.000	1.000
if over age 29	43749	0.457	0.498	0.000	1.000	49113	0.424	0.494	0.000	1.000
<b>Log (wage)</b>	35931	5.130	0.976	-4.605	11.292	37390	6.479	0.937	2.639	11.527
if male	29776	5.229	0.902	-4.605	11.292	30657	6.558	0.870	2.708	11.527
if female	6155	4.649	1.160	-1.715	11.097	6733	6.122	1.132	2.639	10.401
if under age 30	11549	4.796	0.889	-1.204	11.252	11248	6.043	0.776	2.708	9.210
if over age 29	24382	5.288	0.976	-4.605	11.292	26142	6.667	0.939	2.639	11.527
in manufacturing	8428	5.045	0.937	-1.715	11.292	7577	6.345	0.822	2.708	9.741
in hotels and restaurants	675	4.781	0.778	0.811	11.252	833	6.190	0.679	3.829	8.448
in agriculture	2222	4.080	0.856	-0.511	9.687	1995	5.454	0.768	3.219	11.527
in wholesale, retail and repair	2320	4.767	0.782	-0.357	10.922	4311	6.117	0.758	2.708	9.306
in business	1630	5.827	0.809	1.253	11.002	1763	7.113	0.914	3.555	10.401
in other services	13687	5.365	0.951	0.049	11.275	12656	6.835	0.943	2.639	9.580
in transportation	3119	5.245	0.791	0.588	9.393	3146	6.550	0.786	3.401	9.770
in communications	353	5.608	0.717	0.952	6.839	459	6.913	0.837	4.248	8.854

Table 10: Summary Statistics on Grade School Enrollment, only urban areas

Group	Mean in 1993	Standard Deviation in 1993	Mean in 2002	Standard Deviation in 2002	Class size / Class size in 1st grade 1993	Class size / Class size in 1st grade 2002	% Growth since 1993
Boys in grade 1	7618	11435	8739	12778			15%
Boys in grade 2	6354	10130	7663	11819	83%	88%	21%
Boys in grade 3	6146	10068	7416	11662	81%	85%	21%
Boys in grade 4	5777	9592	7058	11123	76%	81%	22%
Boys in grade 5	5816	9979	7316	11602	76%	84%	26%
Boys in grade 6	6066	9789	7330	12284	80%	84%	21%
Boys in grade 7	5425	8489	6917	10744	71%	79%	28%
Boys in grade 8	5247	8179	6907	10471	69%	79%	32%
Boys in grade 9	5301	7416	6609	9231	70%	76%	25%
Boys in grade 10	4269	5785	5719	7311	56%	65%	34%
Boys in grade 11	2544	3528	5006	6133	33%	57%	97%
Boys in grade 12	2366	3329	4503	5205	31%	52%	90%
Girls in grade 1	6707	10639	7775	11601			16%
Girls in grade 2	5668	9492	6871	10757	85%	88%	21%
Girls in grade 3	5423	9325	6623	10574	81%	85%	22%
Girls in grade 4	5059	8761	6282	10007	75%	81%	24%
Girls in grade 5	4963	8807	6441	10247	74%	83%	30%
Girls in grade 6	5011	8600	6381	10765	75%	82%	27%
Girls in grade 7	4529	7613	6064	9729	68%	78%	34%
Girls in grade 8	4196	6943	5903	9212	63%	76%	41%
Girls in grade 9	3819	6216	5378	8107	57%	69%	41%
Girls in grade 10	3041	4768	4666	6485	45%	60%	53%
Girls in grade 11	1714	2926	3800	5421	26%	49%	122%
Girls in grade 12	1529	2707	3368	4706	23%	43%	120%
						Overall:	32%

Table 11: Impact of Linguistic Distance on Growth of IT presence

Dependent variable: Linguistic distance measure:	Any HQ or branch (district level)		Year Firm Established (firm level)	
	Weighted Average	Percent speakers at distance > 2	Weighted Average	Percent speakers at distance > 2
	(1)	(2)	(3)	(4)
Linguistic distance	0.0673 *** (0.0257)	0.3407 *** (0.1216)	-2.7051 *** (0.8505)	-15.0796 *** (5.4308)
Log (population)	0.0346 ** (0.0172)	0.0247 (0.0172)	1.1363 * (0.6315)	1.1485 * (0.6548)
Number of IITs and NIITs	0.1851 *** (0.0652)	0.1816 *** (0.0659)	-1.0899 (0.9763)	-0.7981 (0.9571)
Distance to closest airport	-0.0110 *** (0.0025)	-0.0108 *** (0.0026)	0.0464 (0.0527)	-0.0897 (0.0622)
Engineers	1.9767 (1.6256)	2.0175 (1.6227)	-96.9205 (48.3120)	** -59.0885 (43.5064)
Household wage income	0.0468 (0.0819)	0.0492 (0.0861)	8.4556 (3.4638)	** 8.2529 (3.9692)
Educated wage	0.0387 (0.0935)	0.0263 (0.1027)	-2.2287 (5.1815)	-5.0521 (5.4987)
Salaried	-0.2449 (0.2033)	-0.2217 (0.2017)	-13.3473 (7.6883)	* -15.3950 (7.7061)
Distance to closest big city	0.0006 (0.0009)	0.0008 (0.0010)	0.0948 (0.0258)	*** 0.0728 (0.0254)
Graduate	2.5263 *** (0.4566)	2.4622 *** (0.4660)	4.0136 (20.0222)	3.0525 (19.2820)
Secondary	2.0722 *** (0.3068)	2.1199 *** (0.3029)	6.4765 (12.9487)	4.8880 (12.6780)
Literate	-0.5737 *** (0.1338)	-0.6149 *** (0.1374)	-2.1973 (10.3447)	-3.7686 (10.9752)
Muslim	0.0800 (0.0597)	0.0691 (0.0592)	-28.5462 (11.0981)	** -30.3656 (11.8281)
Train	-0.1232 (0.1385)	-0.1257 (0.1410)	-4.4506 (3.4270)	-0.4865 (2.9307)
Electricity	-0.0741 (0.0813)	-0.0505 (0.0831)	-0.4878 (8.0346)	0.8848 (8.5110)
Native English speakers	166.1686 * (100.6208)	167.3251 * (99.0629)	714.8604 (479.5840)	147.5548 (373.1754)
Native Hindi speakers	-0.0437 (0.0771)	-0.0901 (0.0697)	-1.5676 (4.4488)	-0.3602 (3.8328)
Hindi belt states	0.1363 * (0.0748)	0.2734 ** (0.1063)	-8.2277 ** (4.0398)	** -14.7351 (3.5102)
Observations	1845	1845	2407	2407
R-squared	0.38	0.38	0.15	0.15

include region fixed effects. Regressions also include the percent of people in the district who have college degrees or secondary school degrees, are literate or Muslim, ride a train and the percent of households with electricity. Columns 1 & 2 drop observations for the ten most populous cities in India (as of 1987).

Table 12: Impact of Linguistic Distance on Growth of IT presence

Dependent variable:	Number of HQ and Branches		Number of Employees per Branch		Total Number of Employees		Total Revenue		Total Exports	
	Weighted Average	Percent speakers at distance > 2	Weighted Average	Percent speakers at distance > 2	Weighted Average	Percent speakers at distance > 2	Weighted Average	Percent speakers at distance > 2	Weighted Average	Percent speakers at distance > 2
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Linguistic distance	0.185 (0.137)	1.415 ** (0.593)	57.39 * (33.11)	416.83 ** (194.70)	34.02 (34.68)	167.94 (165.69)	9.07 (25.87)	19.62 (102.49)	10.23 (20.61)	25.05 (81.51)
Log (population)	0.026 (0.108)	0.001 (0.102)	-16.71 (22.94)	-24.64 (22.96)	-23.41 (23.34)	-28.42 (24.11)	-29.44 (25.44)	-30.89 (27.20)	-22.54 (20.87)	-24.16 (22.58)
Number of IITs and NIITs	0.102 (0.430)	0.082 (0.431)	-60.90 (78.31)	-66.74 (75.23)	-101.94 (77.82)	-103.59 (78.04)	-93.20 (85.14)	-93.06 (84.93)	-77.63 (71.84)	-77.53 (71.66)
Distance to closest airport	-0.063 *** (0.019)	-0.063 *** (0.019)	-10.63 ** (4.39)	-10.73 ** (4.34)	-8.86 * (4.80)	-8.77 * (4.70)	-6.61 (5.10)	-6.54 (4.99)	-5.50 (4.27)	-5.42 (4.17)
Engineers	7.609 (14.775)	7.901 (14.618)	1845.37 (2463.75)	1927.77 (2431.60)	1994.55 (3967.27)	2013.48 (3989.20)	-2271.49 (2432.18)	-2276.38 (2441.57)	-1896.11 (2037.42)	-1900.54 (2045.73)
Household wage income	0.403 (0.779)	0.423 (0.786)	16.11 (140.48)	21.79 (140.93)	-76.55 (103.63)	-75.49 (102.84)	-111.42 (112.38)	-111.88 (112.62)	-79.70 (87.28)	-80.14 (87.49)
Educated wage	-0.383 (0.543)	-0.435 (0.558)	-62.85 (97.81)	-78.19 (98.53)	-53.23 (81.26)	-59.34 (79.65)	-70.27 (85.87)	-70.93 (86.09)	-55.03 (68.87)	-55.88 (69.08)
Salaried	3.314 (2.945)	3.323 (2.934)	1029.48 (778.53)	1034.91 (781.94)	1177.09 (894.81)	1189.33 (899.93)	1310.57 (990.15)	1316.55 (999.09)	1026.53 (816.12)	1032.93 (824.64)
Distance to closest big city	-0.008 (0.012)	-0.008 (0.012)	-4.17 (3.04)	-4.00 (3.02)	-4.69 (3.43)	-4.58 (3.40)	-4.81 (3.81)	-4.78 (3.76)	-3.72 (3.16)	-3.69 (3.11)
Graduate	9.793 ** (4.678)	9.716 ** (4.611)	1865.31 (1287.19)	1837.12 (1261.99)	1364.59 (1461.80)	1331.28 (1431.19)	1332.52 (1631.39)	1318.86 (1602.34)	1180.22 (1367.41)	1165.43 (1343.67)
Secondary	9.734 *** (2.882)	9.842 *** (2.872)	1705.17 ** (687.30)	1739.56 ** (691.05)	1116.80 (748.29)	1141.17 (753.76)	1124.10 (814.82)	1132.49 (827.02)	894.16 (677.74)	903.46 (689.24)
Literate	-2.267 (1.478)	-2.495 * (1.488)	-570.89 (446.34)	-636.40 (448.03)	-491.57 (521.86)	-511.37 (520.88)	-606.60 (605.92)	-605.97 (605.57)	-506.02 (511.76)	-506.01 (511.42)
Muslim	0.813 ** (0.382)	0.790 ** (0.374)	172.03 * (96.08)	164.55 * (92.41)	135.53 (103.55)	129.97 (99.98)	109.64 (116.86)	107.80 (113.74)	95.95 (98.02)	93.91 (95.36)
Train	-1.017 (0.725)	-1.038 (0.722)	-265.26 (179.48)	-270.96 (179.14)	-280.67 (199.99)	-281.82 (199.71)	-223.24 (219.68)	-222.86 (219.20)	-181.36 (183.83)	-181.02 (183.44)
Electricity	-1.395 ** (0.649)	-1.284 ** (0.649)	-346.25 ** (164.78)	-313.82 * (163.09)	-298.48 (185.59)	-286.98 (184.61)	-263.10 (203.46)	-262.52 (203.33)	-203.39 (168.61)	-202.45 (168.23)
Native English speakers	451 (486)	448 (481)	51230 (67236)	50598 (65621)	-7008 (31271)	-6366 (30732)	-16492 (28936)	-16018 (28469)	-14941 (23267)	-14443 (22777)
Native Hindi speakers	0.258 (0.342)	0.216 (0.339)	32.35 (64.47)	15.67 (62.42)	97.45 (68.47)	73.19 (65.48)	52.09 (58.33)	41.07 (54.56)	39.75 (47.06)	27.84 (41.66)
Hindi belt states	0.012 (0.380)	0.702 (0.466)	103.60 (108.01)	303.12 * (178.49)	47.91 (97.90)	114.36 (156.99)	-26.98 (58.88)	-25.76 (85.72)	-15.95 (45.88)	-12.67 (67.09)
Observations	1845	1845	1845	1845	1845	1845	1476	1476	1476	1476
R-squared	0.14	0.14	0.08	0.08	0.04	0.04	0.04	0.04	0.04	0.04

Robust standard errors, clustered by district are in parentheses. All regressions include year of data fixed effects and region fixed effects. The measure of linguistic distance in this table is the weighted average of all languages spoken; results are similar using other measures. All regressions drop observations for the ten most populous cities in India (as of 1987).

Table 13: Impact of Linguistic Distance on Growth of IT presence

Dependent variable:	Any HQ or Branch	Years of IT presence	Number of HQ and Branches	Number of Employees per Branch	Total Number of Employees	Total Revenue	Total Exports	Total Subscribed Capital
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Teaching in mother tongue	-1.035 ** (0.447)	-21.4 *** (7.5)	-2.8 ** (1.1)	-440.7 ** (201.2)	-302.0 (198.3)	-84.5 (147.9)	-61.7 (118.9)	-95.9 (60.5)
Log (population)	0.081 *** (0.030)	1.5 ** (0.6)	0.2 (0.1)	-0.6 (25.4)	-11.6 (24.3)	-24.3 (27.0)	-19.2 (22.5)	2.4 (5.7)
Number of IITs and NIITs	0.186 *** (0.071)	3.5 * (1.9)	0.1 (0.4)	-58.0 (80.7)	-100.3 (76.8)	-91.6 (83.7)	-76.1 (70.5)	23.4 (28.4)
Distance to closest airport	-0.010 *** (0.003)	-0.2 *** (0.1)	-0.1 *** (0.0)	-9.9 ** (4.5)	-8.4 * (4.7)	-6.5 (4.9)	-5.4 (4.1)	-1.2 * (0.7)
Engineers	2.468 (1.819)	58.6 (56.8)	7.2 (15.6)	1702.9 (2754.8)	1921.7 (4485.6)	-2752.4 (2830.7)	-2297.7 (2377.0)	-611.3 (441.7)
Household wage income	0.020 (0.083)	1.9 (2.4)	0.3 (0.7)	2.1 (132.4)	-85.6 (104.7)	-113.1 (113.3)	-81.3 (88.1)	-18.7 (23.6)
Educated wage	0.075 (0.179)	1.0 (3.9)	-0.3 (0.7)	-44.1 (122.4)	-40.6 (97.0)	-62.2 (86.5)	-48.6 (69.0)	-9.6 (21.7)
Salaried	-0.517 (0.325)	-6.2 (6.7)	2.7 (3.0)	952.9 (789.6)	1119.0 (907.7)	1312.2 (1013.6)	1032.8 (839.0)	204.2 (146.4)
Distance to closest big city	-0.001 (0.001)	0.0 (0.0)	0.0 (0.0)	-4.7 (3.1)	-5.1 (3.4)	-4.8 (3.7)	-3.7 (3.1)	-0.8 (0.6)
Native English speakers	148.27 (92.07)	2505 (2557)	414 (465)	47807 (62960)	-9894 (31516)	-14661 (27452)	-13050 (21610)	-1999 (5892)
Native Hindi speakers	-0.095 (0.078)	-1.1 (1.6)	0.1 (0.4)	-37.4 (81.9)	58.3 (70.4)	39.5 (58.7)	24.5 (43.4)	22.3 (18.7)
Hindi belt states	-0.022 (0.097)	-0.4 (2.1)	-0.4 (0.4)	2.1 (86.3)	-15.0 (67.5)	-40.2 (54.0)	-29.8 (43.4)	-16.6 (14.8)
Observations	1840	1840	1840	1840	1840	1472	1472	1840
R-squared	0.22	0.19	0.12	0.07	0.04	0.04	0.04	0.04
F-statistic of Instruments	31.6	31.6	31.6	31.6	31.6	25.2	25.2	31.6

Robust standard errors, clustered by district are in parentheses. All regressions include year of data fixed effects and region fixed effects, the percent of college graduates, secondary school graduates, literates, Muslims, people who travel by train and households with electricity. The measure of English cost is primary schools that teach in the mother tongue as the measure of costs of English; results are similar using upper primary schools. As instruments I use the percent of speakers at each distance away from Hindi. All regressions drop observations for the ten most populous cities in India (as of 1987).

Table 14: Impact of District Linguistic Distance on Returns to Education by Age Group and Gender

Dependent variable: Linguistic distance measure:	Salaried Employment									
	Weighted average					Percent speakers at distance > 2				
	All	Men	Women	Age < 30	Age > 29	All	Men	Women	Age < 30	Age > 29
(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	
Linguistic distance * post	-0.002 (0.013)	-0.001 (0.010)	-0.019 (0.032)	-0.005 (0.014)	0.000 (0.014)	-0.003 (0.061)	0.010 (0.050)	-0.071 (0.155)	-0.026 (0.064)	0.005 (0.070)
High school * Linguistic distance * post	-0.001 (0.005)	-0.001 (0.005)	0.004 (0.011)	0.005 (0.008)	-0.010 (0.005)	** 0.008 (0.026)	0.001 (0.026)	0.081 (0.054)	0.044 (0.041)	-0.031 (0.027)
College * Linguistic distance * post	0.016 (0.005)	*** 0.014 (0.006)	** 0.021 (0.013)	* 0.024 (0.009)	*** 0.008 (0.006)	0.093 (0.026)	*** 0.077 (0.028)	*** 0.146 (0.065)	** 0.114 (0.043)	*** 0.061 (0.030)
High school	0.172 (0.017)	*** 0.141 (0.015)	*** 0.411 (0.050)	*** 0.062 (0.017)	*** 0.237 (0.021)	*** 0.188 (0.017)	*** 0.150 (0.015)	*** 0.443 (0.043)	*** 0.080 (0.015)	*** 0.254 (0.021)
College	0.288 (0.020)	*** 0.245 (0.018)	*** 0.456 (0.049)	*** 0.180 (0.023)	*** 0.342 (0.022)	*** 0.300 (0.020)	*** 0.250 (0.019)	*** 0.491 (0.041)	*** 0.195 (0.023)	*** 0.356 (0.023)
Post	-0.028 (0.044)	-0.037 (0.042)	0.115 (0.107)	0.011 (0.051)	-0.052 (0.049)	-0.037 (0.055)	-0.052 (0.048)	0.111 (0.136)	0.012 (0.059)	-0.058 (0.063)
High school * post	-0.019 (0.017)	-0.013 (0.016)	-0.043 (0.048)	-0.015 (0.027)	-0.019 (0.019)	-0.026 (0.015)	* -0.015 (0.015)	-0.074 (0.034)	** -0.020 (0.022)	-0.030 (0.017)
College * post	-0.068 (0.018)	*** -0.061 (0.019)	*** -0.072 (0.042)	* -0.092 (0.028)	*** -0.060 (0.018)	*** -0.066 (0.017)	*** -0.058 (0.018)	*** -0.079 (0.037)	** -0.082 (0.026)	*** -0.065 (0.018)
High school * Linguistic distance	0.016 (0.004)	*** 0.016 (0.004)	*** -0.009 (0.010)	0.013 (0.006)	** 0.021 (0.004)	*** 0.062 (0.026)	** 0.073 (0.024)	*** -0.115 (0.056)	** 0.039 (0.032)	0.083 (0.029)
College * Linguistic distance	0.009 (0.005)	* 0.011 (0.004)	** -0.010 (0.013)	0.002 (0.007)	0.015 (0.005)	*** 0.020 (0.029)	0.048 (0.025)	* -0.134 (0.074)	* -0.025 (0.037)	0.052 (0.033)
Native English speakers * post	19.80 (4.05)	*** 21.11 (4.30)	*** 15.05 (4.87)	*** 25.68 (3.14)	*** 16.25 (5.76)	*** 19.67 (3.99)	*** 20.88 (4.18)	*** 16.24 (5.00)	*** 25.76 (3.15)	*** 15.98 (5.74)
Native Hindi speakers * post	-0.038 (0.037)	-0.034 (0.033)	-0.140 (0.098)	0.009 (0.049)	-0.061 (0.042)	-0.037 (0.032)	-0.032 (0.029)	-0.129 (0.088)	0.008 (0.045)	-0.060 (0.037)
Hindi belt states * post	0.019 (0.032)	0.019 (0.034)	0.028 (0.079)	-0.035 (0.041)	0.042 (0.036)	0.027 (0.044)	0.032 (0.041)	0.025 (0.110)	-0.035 (0.047)	0.049 (0.052)
Observations	1.25E+08	1.03E+08	2.17E+07	4.35E+07	8.10E+07	1.25E+08	1.03E+08	2.17E+07	4.35E+07	8.10E+07
R-squared	0.16	0.15	0.28	0.13	0.19	0.16	0.15	0.28	0.13	0.18

Robust standard errors, clustered by district, in parentheses, also controlling for age, age squared, married, male, whether the individual has ever moved, region fixed effects interacted with post.

Table 15: Impact of District Linguistic Distance on Returns to Education by Age Group

Dependent variable: Linguistic distance measure:	Log (Wages)									
	Weighted average					Percent speakers at distance > 2				
	All	Men	Women	Age < 30	Age > 29	All	Men	Women	Age < 30	Age > 29
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Linguistic distance * post	0.000 (0.017)	0.013 (0.018)	-0.041 (0.046)	-0.028 (0.025)	0.016 (0.020)	-0.010 (0.088)	0.064 (0.093)	-0.214 (0.253)	-0.065 (0.148)	0.017 (0.097)
High school * Linguistic distance * post	-0.033 *** (0.010)	-0.030 *** (0.010)	-0.052 * (0.031)	-0.002 (0.014)	-0.046 *** (0.012)	-0.122 ** (0.055)	-0.116 ** (0.055)	-0.250 (0.158)	-0.002 (0.070)	-0.168 *** (0.062)
College * Linguistic distance * post	-0.017 * (0.010)	-0.013 (0.010)	-0.028 (0.024)	-0.020 (0.021)	-0.022 ** (0.011)	-0.032 (0.052)	-0.021 (0.054)	-0.078 (0.121)	-0.067 (0.101)	-0.054 (0.060)
High school	0.553 *** (0.027)	0.489 *** (0.027)	0.999 *** (0.074)	0.457 *** (0.039)	0.594 *** (0.029)	0.586 *** (0.026)	0.506 *** (0.026)	1.048 *** (0.061)	0.460 *** (0.033)	0.647 *** (0.029)
College	0.996 *** (0.025)	0.904 *** (0.021)	1.456 *** (0.066)	0.919 *** (0.037)	1.022 *** (0.029)	1.033 *** (0.021)	0.926 *** (0.019)	1.510 *** (0.062)	0.939 *** (0.032)	1.073 *** (0.025)
Post	1.113 *** (0.073)	1.094 *** (0.074)	1.255 *** (0.187)	1.160 *** (0.126)	1.082 *** (0.082)	1.129 *** (0.084)	1.086 *** (0.089)	1.337 *** (0.237)	1.146 *** (0.153)	1.124 *** (0.093)
High school * post	0.031 (0.035)	0.037 (0.036)	0.044 (0.118)	-0.171 *** (0.046)	0.118 *** (0.042)	-0.002 (0.035)	0.010 (0.034)	0.006 (0.104)	-0.175 *** (0.039)	0.074 * (0.041)
College * post	0.109 *** (0.030)	0.116 *** (0.035)	0.028 (0.075)	-0.070 (0.067)	0.172 *** (0.036)	0.085 *** (0.031)	0.095 *** (0.034)	-0.006 (0.071)	-0.090 (0.056)	0.146 *** (0.038)
High school * Linguistic distance	0.041 *** (0.009)	0.029 *** (0.008)	0.056 *** (0.021)	0.007 (0.013)	0.059 *** (0.009)	0.174 *** (0.045)	0.133 *** (0.043)	0.252 ** (0.101)	0.036 (0.062)	0.226 *** (0.049)
College * Linguistic distance	0.039 *** (0.008)	0.027 *** (0.008)	0.047 ** (0.018)	0.019 (0.012)	0.049 *** (0.010)	0.133 *** (0.045)	0.101 ** (0.043)	0.155 * (0.092)	0.062 (0.059)	0.158 *** (0.055)
Native English speakers * post	0.774 (4.589)	2.740 (4.486)	-5.116 (9.515)	5.968 (5.358)	-1.845 (5.985)	0.252 (4.663)	2.211 (4.479)	-5.403 (10.164)	5.581 (5.522)	-2.325 (6.251)
Native Hindi speakers * post	-0.074 (0.067)	-0.036 (0.070)	-0.174 (0.165)	-0.008 (0.117)	-0.078 (0.065)	-0.064 (0.063)	-0.033 (0.065)	-0.128 (0.155)	0.030 (0.112)	-0.089 (0.060)
Hindi belt states * post	0.050 (0.049)	0.042 (0.052)	0.031 (0.152)	-0.048 (0.123)	0.070 (0.054)	0.035 (0.056)	0.058 (0.062)	-0.076 (0.190)	-0.067 (0.141)	0.053 (0.066)
Observations	6.72E+07	5.57E+07	1.16E+07	2.16E+07	4.57E+07	6.72E+07	5.57E+07	1.16E+07	2.16E+07	4.57E+07
R-squared	0.67	0.67	0.69	0.65	0.67	0.67	0.67	0.68	0.65	0.67

Robust standard errors, clustered by district, in parentheses, also controlling for age, age squared, married, male, whether the individual has ever moved, region fixed effects interacted with post and a dummy variable for whether the individual is self-employed.

Table 16: Impact of District Linguistic Distance on Returns to Education By Industry

Dependent variable: Log (Wages)	Manufacturing		Hotels and Restaurants		Agriculture, Hunting, Forestry and Fishing		Wholesale, Retail and Repair	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Linguistic distance * post	0.045 (0.042)	0.166 (0.172)	0.092 (0.071)	0.552 (0.396)	0.009 (0.067)	-0.080 (0.303)	-0.003 (0.034)	-0.135 (0.212)
High school * Linguistic distance * post	-0.039 * (0.022)	-0.227 ** (0.115)	0.023 (0.041)	0.238 (0.214)	-0.079 (0.061)	-0.314 (0.330)	0.022 (0.019)	0.079 (0.095)
College * Linguistic distance * post	0.027 (0.026)	0.124 (0.128)	-0.147 * (0.079)	-0.825 * (0.421)	-0.050 (0.095)	-0.504 (0.542)	-0.032 (0.042)	-0.095 (0.240)
High school	0.377 *** (0.042)	0.405 *** (0.036)	0.363 *** (0.129)	0.322 *** (0.122)	0.347 * (0.202)	0.323 (0.201)	0.322 *** (0.057)	0.297 *** (0.050)
College	0.901 *** (0.058)	0.921 *** (0.045)	0.518 ** (0.234)	0.556 ** (0.222)	0.740 *** (0.267)	0.741 *** (0.274)	0.429 ** (0.173)	0.425 ** (0.173)
Post	0.933 *** (0.124)	0.913 *** (0.142)	0.925 * (0.551)	0.780 (0.601)	1.294 *** (0.282)	1.398 *** (0.330)	1.393 *** (0.197)	1.481 *** (0.215)
High school * post	-0.004 (0.070)	-0.013 (0.069)	0.092 (0.157)	0.060 (0.131)	0.222 (0.272)	0.107 (0.244)	-0.076 (0.072)	-0.049 (0.063)
College * post	-0.101 (0.097)	-0.082 (0.087)	0.552 ** (0.227)	0.471 ** (0.210)	0.286 (0.442)	0.400 (0.435)	0.323 * (0.173)	0.288 (0.176)
High school * Linguistic distance	0.052 *** (0.017)	0.262 *** (0.082)	-0.048 (0.038)	-0.231 (0.199)	0.008 (0.049)	0.099 (0.272)	0.002 (0.014)	0.064 (0.071)
College * Linguistic distance	0.024 (0.024)	0.089 (0.114)	-0.011 (0.071)	-0.160 (0.375)	0.089 (0.058)	0.519 (0.322)	0.044 (0.037)	0.255 (0.221)
Native English speakers * post	-12.669 (10.996)	-12.813 (11.532)	6.610 (20.558)	3.858 (21.168)	-71.522 *** (23.647)	-68.499 *** (24.495)	3.490 (8.838)	3.290 (8.500)
Native Hindi speakers * post	0.134 (0.101)	0.093 (0.091)	0.489 (0.718)	0.528 (0.725)	0.003 (0.257)	-0.031 (0.240)	-0.123 (0.131)	-0.138 (0.125)
Hindi belt states * post	0.015 (0.090)	0.070 (0.114)	0.059 (0.514)	0.175 (0.480)	-0.036 (0.185)	-0.103 (0.231)	0.031 (0.177)	-0.042 (0.192)
Observations	1.68E+07	1.68E+07	1337810	1337810	4110879	4110879	6373433	6373433
R-squared	0.71	0.71	0.74	0.74	0.73	0.73	0.67	0.67

Robust standard errors, clustered by district, in parentheses, also controlling for age, age squared, married, male, whether the individual has ever moved, region fixed effects interacted with post and a dummy variable for whether the individual is self-employed.. Odd-numbered columns use the weighted average measure of linguistic distance while even-number columns use the percent of distant speakers.



Table 17: Impact of District Linguistic Distance on Returns to Education By Industry

Dependent variable: Log (Wages)	Financing, Insurance, Real Estate, Computer Related Activities, R & D, Other Business Activities		Other Services (Public Service, Education, Health, Sanitary, Community Services)		Transportation (Land, Water, Air) and Related Services		Communications (Post, Courier, Telecommunications)	
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Linguistic distance * post	-0.152 ** (0.062)	-0.995 *** (0.338)	0.030 (0.030)	0.217 (0.184)	0.097 *** (0.035)	0.613 *** (0.164)	0.274 ** (0.115)
High school * Linguistic distance * post	0.084 * (0.046)	0.607 ** (0.241)	-0.004 (0.018)	0.004 (0.090)	-0.089 *** (0.030)	-0.375 *** (0.142)	-0.079 * (0.046)	-0.388 * (0.232)
College * Linguistic distance * post	0.081 ** (0.041)	0.586 ** (0.248)	0.000 (0.018)	-0.001 (0.092)	-0.053 * (0.030)	-0.227 (0.155)	0.007 (0.108)	0.213 (0.470)
High school	0.519 *** (0.092)	0.509 *** (0.081)	0.557 *** (0.027)	0.594 *** (0.026)	0.387 *** (0.057)	0.425 *** (0.052)	0.237 ** (0.115)	0.252 ** (0.099)
College	1.060 *** (0.091)	1.070 *** (0.103)	0.940 *** (0.044)	0.959 *** (0.041)	0.680 *** (0.056)	0.726 *** (0.059)	0.486 *** (0.114)	0.539 *** (0.092)
Post	1.101 *** (0.273)	1.283 *** (0.296)	1.180 *** (0.142)	1.084 *** (0.191)	0.851 *** (0.212)	0.685 *** (0.215)	-0.125 (0.423)	-0.205 (0.488)
High school * post	-0.112 (0.134)	-0.152 (0.114)	0.093 (0.062)	0.084 (0.061)	0.201 (0.125)	0.129 (0.112)	0.205 (0.160)	0.139 (0.131)
College * post	-0.085 (0.160)	-0.124 (0.174)	0.109 * (0.056)	0.110 ** (0.055)	0.217 ** (0.099)	0.184 ** (0.094)	0.129 (0.223)	0.015 (0.190)
High school * Linguistic distance	-0.049 (0.034)	-0.268 (0.175)	0.048 *** (0.009)	0.199 *** (0.047)	0.047 *** (0.015)	0.190 ** (0.078)	0.016 (0.043)	0.078 (0.195)
College * Linguistic distance	-0.042 * (0.023)	-0.267 * (0.144)	0.038 *** (0.013)	0.175 *** (0.062)	0.047 ** (0.018)	0.163 (0.103)	-0.003 (0.072)	-0.093 (0.300)
Native English speakers * post	-3.688 (11.750)	-3.306 (11.714)	-16.904 ** (8.180)	-17.961 ** (7.649)	3.525 (12.109)	1.093 (10.944)	-7.528 (17.365)	-12.955 (18.735)
Native Hindi speakers * post	-0.181 (0.190)	-0.134 (0.181)	-0.036 (0.092)	-0.048 (0.084)	-0.044 (0.139)	-0.060 (0.128)	0.306 (0.353)	0.080 (0.310)
Hindi belt states * post	0.403 ** (0.158)	0.211 (0.192)	0.042 (0.128)	0.149 (0.176)	0.147 (0.113)	0.346 *** (0.111)	1.193 *** (0.382)	1.540 *** (0.466)
Observations	2952146	2952146	2.07E+07	2.07E+07	6028767	6028767	740670	740670
R-squared	0.73	0.73	0.69	0.69	0.7	0.7	0.79	0.79

Robust standard errors, clustered by district, in parentheses, also controlling for age, age squared, married, male, whether the individual has ever moved, region fixed effects interacted with post and a dummy variable for whether the individual is self-employed.. Odd-numbered columns use the weighted average measure of linguistic distance while even-number columns use the percent of distant speakers.

Table 18: Impact of Weighted Average Linguistic Distance on Grade School Enrollment

	All Grades		Primary (Grades 1-5)		Upper Primary (Grades 6-8)		Secondary (Grades 9-12)		
	All	Girls	Boys	Girls	Boys	Girls	Boys	Girls	Boys
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Linguistic distance * post	0.111 *** (0.033)	0.110 *** (0.032)	0.116 *** (0.036)	0.041 * (0.025)	0.029 (0.026)	0.063 ** (0.026)	0.037 (0.024)	0.246 *** (0.065)	0.249 *** (0.077)
Log enrollment in pre year	0.040 * (0.024)	0.018 (0.025)	0.004 (0.029)	-0.118 (0.078)	-0.162 ** (0.080)	-0.103 (0.063)	-0.113 (0.072)	-0.100 *** (0.028)	-0.127 *** (0.034)
Log child population in pre year	-0.012 (0.048)	-0.037 (0.047)	0.002 (0.062)	0.044 (0.078)	0.109 (0.095)	-0.007 (0.050)	0.012 (0.061)	0.252 *** (0.070)	0.222 ** (0.110)
Log child population in post year	0.813 *** (0.112)	0.660 *** (0.130)	0.721 *** (0.132)	0.705 *** (0.154)	0.723 *** (0.135)	0.861 *** (0.098)	0.622 *** (0.119)	0.654 ** (0.257)	0.839 *** (0.251)
Household wage income * post	0.067 (0.054)	0.013 (0.051)	0.101 * (0.058)	-0.057 (0.054)	-0.030 (0.057)	0.015 (0.050)	0.020 (0.059)	0.131 (0.117)	0.207 * (0.118)
Educated wage * post	-0.062 (0.113)	-0.042 (0.114)	-0.074 (0.116)	-0.167 (0.131)	-0.172 (0.119)	-0.195 * (0.104)	-0.143 (0.100)	0.163 (0.283)	0.061 (0.275)
Salaried * post	-0.257 (0.277)	-0.164 (0.302)	-0.158 (0.277)	-0.303 (0.308)	-0.197 (0.301)	-0.274 (0.307)	0.177 (0.283)	-0.209 (0.487)	-0.098 (0.529)
Distance to closest big city * post	-0.003 *** (0.001)	-0.003 *** (0.001)	-0.003 *** (0.001)	-0.003 *** (0.001)	-0.002 *** (0.001)	-0.001 (0.001)	-0.002 * (0.001)	-0.004 * (0.002)	-0.003 (0.002)
Graduate * post	-0.771 (0.553)	-1.058 * (0.550)	-0.390 (0.665)	0.337 (0.502)	0.338 (0.490)	-0.208 (0.562)	-0.462 (0.574)	-2.170 ** (0.973)	-1.479 (1.161)
Secondary * post	-0.280 (0.462)	-0.051 (0.444)	-0.262 (0.506)	0.159 (0.380)	0.113 (0.381)	-0.217 (0.387)	-0.289 (0.371)	0.482 (0.848)	-0.182 (1.015)
Literate * post	0.082 (0.233)	-0.048 (0.249)	0.022 (0.231)	-0.401 (0.247)	-0.275 (0.231)	-0.456 ** (0.225)	-0.375 * (0.212)	0.043 (0.400)	0.286 (0.483)
Muslim * post	-0.006 (0.123)	0.037 (0.124)	-0.004 (0.123)	0.138 (0.149)	0.031 (0.142)	-0.011 (0.117)	-0.089 (0.093)	0.267 (0.187)	0.003 (0.202)
Train * post	0.016 (0.153)	0.097 (0.156)	-0.037 (0.163)	-0.122 (0.129)	-0.086 (0.139)	0.036 (0.144)	0.048 (0.126)	0.093 (0.280)	-0.015 (0.311)
Electricity * post	0.098 (0.141)	0.020 (0.148)	0.134 (0.132)	-0.027 (0.134)	0.003 (0.129)	0.202 (0.168)	0.232 * (0.134)	-0.152 (0.269)	-0.028 (0.271)
Native English speakers * post	-13.359 (16.720)	-20.196 (18.018)	-18.968 (17.100)	7.673 (15.821)	0.845 (16.591)	-14.346 (16.438)	-11.818 (13.682)	-55.245 (39.851)	-64.649 * (37.123)
Native Hindi speakers * post	0.336 *** (0.090)	0.426 *** (0.101)	0.299 *** (0.084)	0.327 *** (0.099)	0.256 *** (0.093)	0.434 *** (0.126)	0.265 *** (0.086)	0.593 *** (0.154)	0.430 *** (0.146)
Hindi belt states * post	0.148 (0.108)	0.119 (0.113)	0.162 (0.112)	-0.214 *** (0.075)	-0.242 *** (0.073)	-0.389 *** (0.095)	-0.295 *** (0.077)	0.629 *** (0.217)	0.821 *** (0.265)
Constant	-2.445 * (1.323)	-0.492 (1.557)	-1.020 (1.633)	1.232 (1.888)	0.850 (1.680)	-0.509 (1.315)	2.145 (1.533)	-2.483 (2.803)	-2.281 (2.737)
Number of Observations	17169	8509	8660	3605	3635	2112	2157	2792	2868
R-squared	0.888	0.902	0.884	0.978	0.978	0.983	0.982	0.872	0.847

Robust standard errors, clustered by district are in parentheses. The measure of linguistic distance used is the weighted average of all languages spoken. Also includes district, timeperiod, gender and grade level fixed effects, and both region fixed effects and grade level fixed effects interacted with post.

Table 19: Impact of Percent Distant Speakers on Grade School Enrollment

	All Grades		Primary (Grades 1-5)		Upper Primary (Grades 6-8)		Secondary (Grades 9-12)		
	All	Girls	Boys	Girls	Boys	Girls	Boys	Girls	Boys
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
Linguistic distance * post	0.278 (0.172)	0.255 (0.167)	0.269 (0.176)	0.242 * (0.139)	0.196 * (0.116)	0.428 *** (0.155)	0.288 ** (0.122)	0.251 (0.279)	0.294 (0.379)
Log enrollment in pre year	0.039 (0.024)	0.017 (0.025)	0.002 (0.029)	-0.117 (0.078)	-0.160 ** (0.079)	-0.105 * (0.064)	-0.113 (0.071)	-0.108 *** (0.028)	-0.135 *** (0.034)
Log child population in pre year	-0.006 (0.048)	-0.030 (0.047)	0.012 (0.063)	0.047 (0.079)	0.108 (0.095)	0.000 (0.050)	0.013 (0.061)	0.263 *** (0.071)	0.245 ** (0.111)
Log child population in post year	0.842 *** (0.110)	0.692 *** (0.127)	0.758 *** (0.132)	0.720 *** (0.152)	0.733 *** (0.134)	0.883 *** (0.094)	0.639 *** (0.119)	0.707 *** (0.261)	0.926 *** (0.257)
Household wage income * post	0.064 (0.056)	0.008 (0.054)	0.098 * (0.059)	-0.053 (0.052)	-0.027 (0.056)	0.019 (0.051)	0.025 (0.058)	0.109 (0.123)	0.192 (0.126)
Educated wage * post	-0.067 (0.124)	-0.046 (0.127)	-0.078 (0.128)	-0.175 (0.133)	-0.178 (0.119)	-0.204 * (0.104)	-0.152 (0.099)	0.155 (0.322)	0.063 (0.302)
Salaried * post	-0.231 (0.283)	-0.138 (0.308)	-0.133 (0.281)	-0.311 (0.308)	-0.207 (0.300)	-0.298 (0.309)	0.159 (0.286)	-0.079 (0.501)	-0.038 (0.538)
Distance to closest big city * post	-0.002 ** (0.001)	-0.002 ** (0.001)	-0.002 ** (0.001)	-0.002 ** (0.001)	-0.002 ** (0.001)	-0.001 (0.001)	-0.001 (0.001)	-0.002 (0.002)	-0.002 (0.002)
Graduate * post	-1.042 * (0.613)	-1.340 ** (0.591)	-0.680 (0.713)	0.293 (0.527)	0.303 (0.498)	-0.302 (0.597)	-0.475 (0.612)	-3.044 *** (1.067)	-2.093 * (1.268)
Secondary * post	-0.335 (0.475)	-0.108 (0.457)	-0.329 (0.525)	0.153 (0.374)	0.115 (0.377)	-0.223 (0.380)	-0.293 (0.363)	0.323 (0.918)	-0.438 (1.081)
Literate * post	0.128 (0.229)	0.006 (0.244)	0.075 (0.226)	-0.424 * (0.252)	-0.298 (0.234)	-0.503 ** (0.227)	-0.410 * (0.218)	0.262 (0.397)	0.482 (0.481)
Muslim * post	-0.035 (0.124)	0.007 (0.126)	-0.035 (0.124)	0.129 (0.147)	0.026 (0.140)	-0.022 (0.117)	-0.093 (0.092)	0.189 (0.192)	-0.071 (0.203)
Train * post	0.014 (0.149)	0.096 (0.153)	-0.037 (0.157)	-0.134 (0.131)	-0.093 (0.140)	0.018 (0.142)	0.035 (0.125)	0.116 (0.265)	0.013 (0.285)
Electricity * post	0.122 (0.147)	0.042 (0.153)	0.157 (0.139)	-0.002 (0.139)	0.023 (0.131)	0.242 (0.173)	0.259 * (0.139)	-0.152 (0.275)	0.003 (0.278)
Native English speakers * post	-7.622 (16.276)	-14.532 (17.709)	-12.683 (16.515)	10.936 (16.004)	3.356 (16.824)	-8.434 (15.806)	-8.519 (13.321)	-44.615 (40.043)	-53.783 (36.768)
Native Hindi speakers * post	0.209 *** (0.080)	0.295 *** (0.092)	0.163 ** (0.071)	0.306 *** (0.091)	0.245 *** (0.085)	0.408 *** (0.124)	0.261 *** (0.080)	0.242 ** (0.116)	0.100 (0.115)
Hindi belt states * post	0.187 (0.148)	0.145 (0.138)	0.185 (0.155)	-0.111 (0.097)	-0.152 * (0.084)	-0.190 (0.118)	-0.157 (0.107)	0.482 ** (0.198)	0.678 ** (0.319)
Constant	-2.822 ** (1.319)	-0.901 (1.525)	-1.517 (1.655)	1.042 (1.881)	0.735 (1.688)	-0.804 (1.269)	1.944 (1.528)	-3.102 (2.850)	-3.401 (2.852)
Number of Observations	17169	8509	8660	3605	3635	2112	2157	2792	2868
R-squared	0.888	0.902	0.884	0.978	0.978	0.983	0.982	0.870	0.846

Robust standard errors, clustered by district are in parentheses. The measure of linguistic distance used is the percent of speakers at a distance greater than 2. Also includes district, timeperiod, gender and grade level fixed effects, and both region fixed effects and grade level fixed effects interacted with post.

Table 20: IV Regression of Cost of English on Grade School Enrollment

Type of School:	Primary Schools						Upper Primary Schools					
Linguistic Distance Instrument:	Weighted Average			Percent Speakers at each distance			Weighted Average			Percent Speakers at each distance		
	All	Girls	Boys	All	Girls	Boys	All	Girls	Boys	All	Girls	Boys
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
Teaching in mother tongue * post	-2.87 (1.94)	-2.87 (1.99)	-3.18 (2.21)	-1.05 *** (0.29)	-0.81 *** (0.30)	-1.25 *** (0.33)	-2.09 ** (0.99)	-2.10 ** (1.07)	-2.23 ** (1.02)	-1.07 *** (0.26)	-0.90 *** (0.27)	-1.27 *** (0.29)
Log enrollment in pre year	0.04 * (0.02)	0.02 (0.02)	0.00 (0.03)	0.05 ** (0.02)	0.03 (0.02)	0.01 (0.03)	0.05 ** (0.02)	0.03 (0.02)	0.01 (0.03)	0.05 ** (0.02)	0.03 (0.02)	0.01 (0.03)
Log child population in pre year	-0.09 (0.08)	-0.10 (0.09)	-0.13 (0.13)	-0.03 (0.05)	-0.04 (0.04)	-0.04 (0.06)	-0.04 (0.06)	-0.05 (0.06)	-0.07 (0.08)	-0.02 (0.05)	-0.03 (0.05)	-0.03 (0.06)
Log child population in post year	0.60 ** (0.23)	0.36 (0.33)	0.46 * (0.28)	0.78 *** (0.12)	0.64 *** (0.14)	0.65 *** (0.15)	0.62 *** (0.17)	0.40 * (0.22)	0.49 *** (0.19)	0.74 *** (0.11)	0.60 *** (0.13)	0.61 *** (0.13)
Household wage income * post	0.01 (0.16)	-0.05 (0.16)	0.03 (0.18)	0.04 (0.08)	-0.02 (0.07)	0.07 (0.10)	0.04 (0.12)	-0.02 (0.12)	0.07 (0.13)	0.05 (0.08)	-0.01 (0.07)	0.08 (0.09)
Educated wage * post	0.08 (0.43)	0.14 (0.40)	0.03 (0.49)	0.02 (0.22)	0.05 (0.18)	-0.01 (0.25)	0.11 (0.32)	0.16 (0.30)	0.08 (0.35)	0.05 (0.21)	0.07 (0.18)	0.03 (0.24)
Salaried * post	-0.83 (0.70)	-0.79 (0.71)	-0.71 (0.75)	-0.37 (0.34)	-0.25 (0.35)	-0.29 (0.37)	-0.95 (0.61)	-0.88 (0.63)	-0.88 (0.63)	-0.54 (0.35)	-0.39 (0.36)	-0.51 (0.38)
Distance to closest big city * post	-0.01 * (0.00)	-0.01 * (0.00)	-0.01 (0.01)	0.00 *** (0.00)	0.00 *** (0.00)	0.00 *** (0.00)	-0.01 ** (0.00)	-0.01 ** (0.00)	-0.01 ** (0.00)	0.00 *** (0.00)	0.00 *** (0.00)	0.00 *** (0.00)
Graduate * post	0.53 (1.38)	0.16 (1.27)	1.16 (1.80)	-0.55 (0.63)	-1.00 * (0.58)	-0.03 (0.80)	-0.48 (0.78)	-0.78 (0.87)	-0.03 (0.87)	-0.82 (0.55)	-1.17 ** (0.57)	-0.36 (0.65)
Secondary * post	0.44 (0.84)	0.65 (0.84)	0.59 (0.96)	0.02 (0.53)	0.18 (0.50)	0.09 (0.61)	0.45 (0.65)	0.65 (0.68)	0.55 (0.71)	0.12 (0.48)	0.28 (0.47)	0.22 (0.53)
Literate * post	0.53 (0.52)	0.36 (0.52)	0.53 (0.60)	0.35 (0.26)	0.19 (0.26)	0.32 (0.27)	0.27 (0.37)	0.11 (0.38)	0.23 (0.39)	0.26 (0.26)	0.12 (0.27)	0.21 (0.27)
Muslim * post	0.13 (0.17)	0.15 (0.17)	0.16 (0.20)	0.05 (0.12)	0.06 (0.12)	0.06 (0.12)	0.08 (0.15)	0.10 (0.15)	0.09 (0.16)	0.04 (0.12)	0.06 (0.12)	0.05 (0.13)
Train * post	0.18 (0.30)	0.23 (0.29)	0.22 (0.37)	0.06 (0.16)	0.12 (0.16)	0.05 (0.18)	0.10 (0.26)	0.16 (0.27)	0.10 (0.29)	0.04 (0.17)	0.11 (0.17)	0.03 (0.20)
Electricity * post	-0.32 (0.35)	-0.39 (0.37)	-0.32 (0.39)	-0.11 (0.14)	-0.16 (0.14)	-0.09 (0.15)	-0.20 (0.23)	-0.28 (0.25)	-0.16 (0.24)	-0.10 (0.14)	-0.16 (0.14)	-0.07 (0.15)
Native English speakers * post	-127.5 (97.8)	-132.7 (99.8)	-153.1 (116.4)	-50.0 * (27.5)	-46.0 * (27.5)	-66.7 ** (30.8)	-179.8 (114.7)	-187.6 (121.4)	-201.0 * (121.8)	-94.7 * (52.8)	-86.7 * (48.8)	-118.9 * (61.1)
Native Hindi speakers * post	0.26 * (0.14)	0.36 ** (0.15)	0.22 (0.14)	0.19 ** (0.09)	0.27 *** (0.10)	0.15 * (0.09)	0.35 *** (0.13)	0.44 *** (0.14)	0.32 ** (0.13)	0.25 *** (0.09)	0.33 *** (0.10)	0.23 *** (0.08)
Hindi belt states * post	-0.19 (0.21)	-0.21 (0.21)	-0.21 (0.23)	-0.06 (0.11)	-0.07 (0.11)	-0.07 (0.11)	-0.09 (0.15)	-0.10 (0.16)	-0.11 (0.15)	-0.04 (0.10)	-0.05 (0.11)	-0.05 (0.10)
Constant	0.71 (3.10)	5.27 (4.20)	4.95 (4.08)	-1.94 (1.35)	-0.35 (1.55)	1.63 (1.81)	0.04 (1.96)	4.32 * (2.63)	3.85 (2.38)	-1.67 (1.20)	0.00 (1.47)	1.96 (1.54)
Number of Observations	17121	8485	8636	17121	8485	8636	17121	8485	8636	17121	8485	8636
R-squared	0.858	0.875	0.848	0.879	0.895	0.873	0.869	0.884	0.863	0.879	0.895	0.874
F-statistic of instruments	277.04	136.08	20.91	447.32	203.73	228.17	411.88	196.82	193.22	387.63	174.86	196.85

Robust standard errors, clustered by district are in parentheses. Also includes district, timeperiod, gender and grade level fixed effects, and both region fixed effects and grade level fixed effects interacted with post.

Table 21: Impact of District Linguistic Distance on Individual School Enrollment

Linguistic distance measure: Dep Variable: Attending	Weighted average			Percent speakers at distance > 2		
	All	Boys	Girls	All	Boys	Girls
	(1)	(2)	(3)	(4)	(5)	(6)
Linguistic distance * post	-0.004 (0.009)	-0.012 * (0.007)	0.003 (0.013)	-0.017 (0.041)	-0.045 (0.038)	0.005 (0.060)
Linguistic distance * post * ages 11 - 15	0.008 ** (0.003)	0.004 (0.004)	0.010 ** (0.005)	0.035 ** (0.017)	0.022 (0.018)	0.043 * (0.023)
Linguistic distance * post * ages 16 - 20	0.014 *** (0.004)	0.018 *** (0.005)	0.007 * (0.004)	0.069 *** (0.021)	0.087 *** (0.026)	0.036 (0.023)
Linguistic distance * ages 11 - 15	-0.010 *** (0.003)	-0.009 *** (0.003)	-0.009 ** (0.004)	-0.056 *** (0.014)	-0.055 *** (0.014)	-0.054 *** (0.020)
Linguistic distance * ages 16 - 20	-0.017 *** (0.003)	-0.016 *** (0.004)	-0.017 *** (0.004)	-0.082 *** (0.017)	-0.073 *** (0.020)	-0.084 *** (0.021)
Native English speakers * post	-6.912 *** (2.313)	-2.819 (2.032)	-11.485 *** (3.225)	-6.936 *** (2.324)	-2.761 (2.040)	-11.534 *** (3.293)
Native Hindi speakers * post	-0.021 (0.026)	-0.063 ** (0.027)	0.021 (0.040)	-0.023 (0.023)	-0.057 ** (0.025)	0.012 (0.036)
Hindi belt states * post	0.039 (0.024)	0.038 (0.026)	0.040 (0.036)	0.047 (0.029)	0.039 (0.031)	0.051 (0.043)
Ages 11 - 15 * post	0.029 *** (0.010)	0.025 ** (0.010)	0.040 ** (0.016)	0.033 *** (0.009)	0.025 *** (0.009)	0.046 *** (0.014)
Ages 16 - 20 * post	0.010 (0.014)	-0.013 (0.017)	0.041 *** (0.015)	0.017 (0.013)	-0.005 (0.016)	0.044 *** (0.015)
Whether individual moved	0.061 *** (0.008)	0.024 ** (0.011)	0.094 *** (0.007)	0.061 *** (0.008)	0.024 ** (0.011)	0.094 *** (0.007)
Male	0.067 *** (0.004)			0.067 *** (0.004)		
Number of HH members	-0.006 *** (0.001)	-0.005 *** (0.001)	-0.006 *** (0.001)	-0.006 *** (0.001)	-0.005 *** (0.001)	-0.006 *** (0.001)
Household wage income	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
Observations	1.19E+08	6.28E+07	5.63E+07	1.19E+08	6.28E+07	5.63E+07
R-squared	0.44	0.41	0.47	0.44	0.41	0.47

Robust standard errors, clustered by district, in parentheses. Also controlling for age dummies, household religion dummies, head of household education dummies and region interacted with post dummies.

Table 22: Impact of District Linguistic Distance on Individual Educational Achievement

Linguistic distance measure: Dep Variable: Edu. Level	Weighted average			Percent speakers at distance > 2		
	All	Boys	Girls	All	Boys	Girls
	(1)	(2)	(3)	(4)	(5)	(6)
Linguistic distance * post	-0.042 (0.025)	* -0.060 (0.026)	** -0.029 (0.032)	-0.055 (0.119)	-0.160 (0.133)	0.029 (0.148)
Linguistic distance * post * ages 11 - 15	0.042 (0.011)	*** 0.041 (0.013)	*** 0.039 (0.013)	*** 0.148 (0.057)	*** 0.163 (0.067)	** 0.120 (0.069)
Linguistic distance * post * ages 16 - 20	0.050 (0.016)	*** 0.050 (0.016)	*** 0.040 (0.022)	* 0.200 (0.076)	*** 0.213 (0.081)	*** 0.142 (0.105)
Linguistic distance * ages 11 - 15	0.012 (0.007)	0.008 (0.009)	0.017 (0.010)	* 0.038 (0.038)	0.023 (0.043)	0.059 (0.049)
Linguistic distance * ages 16 - 20	-0.021 (0.011)	* -0.004 (0.010)	-0.028 (0.018)	-0.141 (0.056)	** -0.053 (0.054)	-0.191 (0.088)
Native English speakers * post	-37.783 (9.035)	*** -20.039 (10.790)	* -58.485 (9.499)	*** -38.426 (9.146)	*** -20.559 (10.905)	* -59.204 (9.484)
Native Hindi speakers * post	0.030 (0.113)	-0.104 (0.127)	0.145 (0.132)	0.064 (0.101)	-0.060 (0.118)	0.178 (0.120)
Hindi belt states * post	-0.045 (0.097)	-0.067 (0.100)	-0.003 (0.130)	0.000 (0.093)	-0.060 (0.097)	0.063 (0.133)
Ages 11 - 15 * post	-0.011 (0.035)	-0.042 (0.043)	0.025 (0.044)	0.024 (0.035)	-0.012 (0.042)	0.065 (0.042)
Ages 16 - 20 * post	0.173 (0.050)	*** 0.047 (0.053)	0.323 (0.072)	*** 0.207 (0.045)	*** 0.077 (0.050)	0.354 (0.063)
Whether individual moved	0.250 (0.022)	*** 0.079 (0.030)	*** 0.389 (0.027)	*** 0.250 (0.022)	*** 0.078 (0.030)	*** 0.389 (0.027)
Male	0.180 (0.014)	*** 0.000 (0.000)	0.000 (0.000)	0.179 (0.014)	*** 0.000 (0.000)	0.000 (0.000)
Number of HH members	-0.015 (0.004)	*** -0.015 (0.004)	*** -0.011 (0.006)	* -0.015 (0.004)	*** -0.015 (0.004)	*** -0.011 (0.006)
Household wage income	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)	0.000 (0.000)
Observations	1.30E+08	6.84E+07	6.18E+07	1.30E+08	6.84E+07	6.18E+07
R-squared	0.47	0.49	0.46	0.47	0.49	0.46

Robust standard errors, clustered by district, in parentheses. Also controlling for age dummies, household religion dummies, head of household education dummies and region interacted with post dummies.