

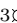









Statistical analysis of 1-grams flow among Indo-European languages

Name1 Surname ^{1,2}, Name2 Surname ², Name3 Surname ^{2,3}, Name4 Surname ²,
Name5 Surname ², Name6 Surname ², Name7 Surname ^{1,2,3}^{*}, with the Lorem Ipsum Consortium[¶]

- 1** Affiliation Dept/Program/Center, Institution Name, City, State, Country
2 Affiliation Dept/Program/Center, Institution Name, City, State, Country
3 Affiliation Dept/Program/Center, Institution Name, City, State, Country

 These authors contributed equally to this work.
 These authors also contributed equally to this work.
 Current Address: Dept/Program/Center, Institution Name, City, State, Country
 Deceased
 Membership list can be found in the Acknowledgments section.
^{*} correspondingauthor@institute.edu

Abstract

Author summary

Introduction

In recent years, the field of linguistics has been benefited from the development of more sophisticated computational tools, helping to process a greater amount of data in less time and allowing the study of linguistics from a statistical perspective. This statistical study began with the works of George Zipf [?], in them Zipf argues that if the words used in a text are ranked by their frequency of appearance, where the lower ranks belong to the most frequent words, then the frequency f of any word and its rank k are related by a power law of the form $f \propto 1/k$. The previous expression is known as Zipf's law, and it has not only been tested on language datasets, Morales et al. [?] have also proved on sports and games data, and Cristelli et al. [?] on the gross domestic product of several countries, wealth of American citizens, and population of cities.

Although Zip's law has opened several statistical studies in linguistics, nowadays few studies have been done about how within the vocabulary of a language, words from the language itself and from other languages are mixed. Currently in the spanish language, there are words from the english language that do not have a translation or that sometimes displace those that already exist in spanish. For example, for native spanish speakers it is common to hear the word marketing instead of mercadeo when dealing with economic or business issues; also the word online has replaced en línea, when referring to issues related to the internet, a word without translation in spanish.

This trending is not only affecting spanish, but also other languages that are being influenced by topics where english is the main and common language for communication. However, in different periods of time, the flow of words came from other languages. D'Amore [?] discusses with linguistic rigor the flow of words between english and spanish, showing historical and cultural causes that allowed such flow; in addition to mentioning the influence of arabic in spanish and french in english.

In this work, we use the Google Books N-gram [?] dataset of the most frequent words in books published in english, french, german, italian and spanish languages. With that dataset, we develop an algorithm that identifies the words of one language and that are being used by others. Once these words have been classified, we construct two models to quantify the influence that one language has had on another during the 20th century. The first model, we count the number of new words that a language received from another, while in the second method, we develop the concept of the use of one language in another, from quantifying the relative frequency of the words of a language that are being used in another language. In both we identify historical, social and cultural causes that are responsible for the flow of words and the constantly use of them.

Finally, we use the concept of rank diversity, that shows the number of words occupying a certain rank across the time. This study shows that regardless of who is the language the flow comes from or who language receives them, the lower ranks are always occupied by fewer words, and as the rank increases the diversity curve also increases as a logarithmic.

Metodology

For the development of this work, we used the Google Books Ngram data set. This dataset is made up of listing for each year and for each language of publication the most used “n-grams”. The “n-grams” are the words or set of words that make up the text of a book, where the number n indicates the number of words that make up the gram, being a 1-gram an individual word, a 2-gram a phrase composed of two 1-gram, 3-gram the set of three 1-gram and so on.

From this data set, the lists of the five thousand most used 1-grams each year between 1740 and 2009 were extracted for the English, French, German, Italian and Spanish languages. In each list, the words are ranked according to their frequency of appearance, where the most frequent words have the lowest ranks.

To determine the presence of one language in another, an algorithm was developed to find the words that are common between at least two languages, these must be the same in writing, letter by letter. These words were defined as **migrant words**.

A migrant word is associated with a **source language** and a **receiving language**, where the source language is the one where the word appeared for the first time within the five thousand most used words, while the receiving language is the one where the word is also present, being a different set from the source language. To determine the source language, we established that this will be the language where the word appeared for the first time within the five thousand most used words; in the case of a migrant word has appeared in the same year in two or more languages, the source is the one where the word has the lowest rank.

The previous criteria of looking for words with the same writing and later associating them with a source language, is not perfect. One of the most common errors was finding words with the same writing, but with different meanings, for example mayor in English refers to the representative of the government in a locality, while in Spanish, mayor is an adjective to indicate that something is bigger or older. Another recurring error was to not distinguish words with the same meaning but with different endings, for example, the word imagine is written imaginer in French and imaginari in Spanish. Finally, in some cases, the authentic source language is some other language for which there is no information in the data set, for example the word natural comes from Greek, but there is no data from Greek language in the google books n-grams data set, consequently this word was associated with English as its source language.

The above errors were detected by individually analyzing each of the migrant words and their corresponding source and receiving languages. One way to have cleaner data

is by consulting an expert in each language, who review the words and decide which ones were classified properly, however, this is not practical since if there were more languages in the database, it would be necessary to consult an expert for each language.

New words

The purpose of our work is to find a method to quantify the influence that one language has had on another, however there is no a generalized process that shows how much influence the things have. We built a first method to quantify the influence that one language has on another, this consisted of analyzing the times when a migrant word with a certain source language appears for the first time in the receiving language, finding relationships between those times and the words involved on the move. Migrant words with the same source language and that appear for the first time in the receiving language were referred to as **new migrant words**.

The previous process was carried out in three different ways, the first by taking a source language A and counting the new migrant words in the different receiving languages for each decade, this shown in which languages A has been influential; the second way was by taking language A as the receiving language, and counting how many new migrant words of different source languages are present in each decade, this shown the languages that have been influential on A ; finally, by taking a language A and a language B , and counting how many new migrant words from A are in B , and how many from B are in A , this shown the influence between A and B .

A complementary part of interpreting the influence between languages is to recognize the causes that originate the migrations of words. These were inferred after analyzing the meaning of the new migrant words in a given year of migration. According to [?] a semantic field is a set of associated words that share part of their meaning. Then, from classifying the migrant words within a semantic field, the cause or event that originated the migrations might be found, since the migrations occur in the same year (or in the years around) where the event occurred.

The fact of looking for relationships between the migrant words and their meaning, we decided to omit the functional words from the results, these words are those that help to structure a message according to the grammar of the language, such as articles, prepositions and conjunctions. In this way, all new migrant words are content words, that is, those that carry the information and meaning of the message.

Analysis of new migrant words

In order to quantify how many new words have migrated between the English, French, German, Italian and Spanish languages, Fig 1 shows, for each decade of the 20th century, the behavior of each language, being the source language or the receiving language of new words, representing how one language has influenced others and also how it has been influenced. Throughout the 20th century, the English language has contributed almost three times as many words as it has received, where its greatest apogee occurred in the 1940s, during which World War II occurred; and in the first decade of 2000, where access to technology was feasible for larger sectors of the population, after the development of globalization.

It is also notable that the French and German languages were mainly languages influenced by others, since they received more words than they contributed, this trend being more appreciable since 1940 after the end of World War II. In addition, since 1980, German has obtained 30 new words every decade compared to the previous decade, consequently it has been the largest receiver in the last 30 years.

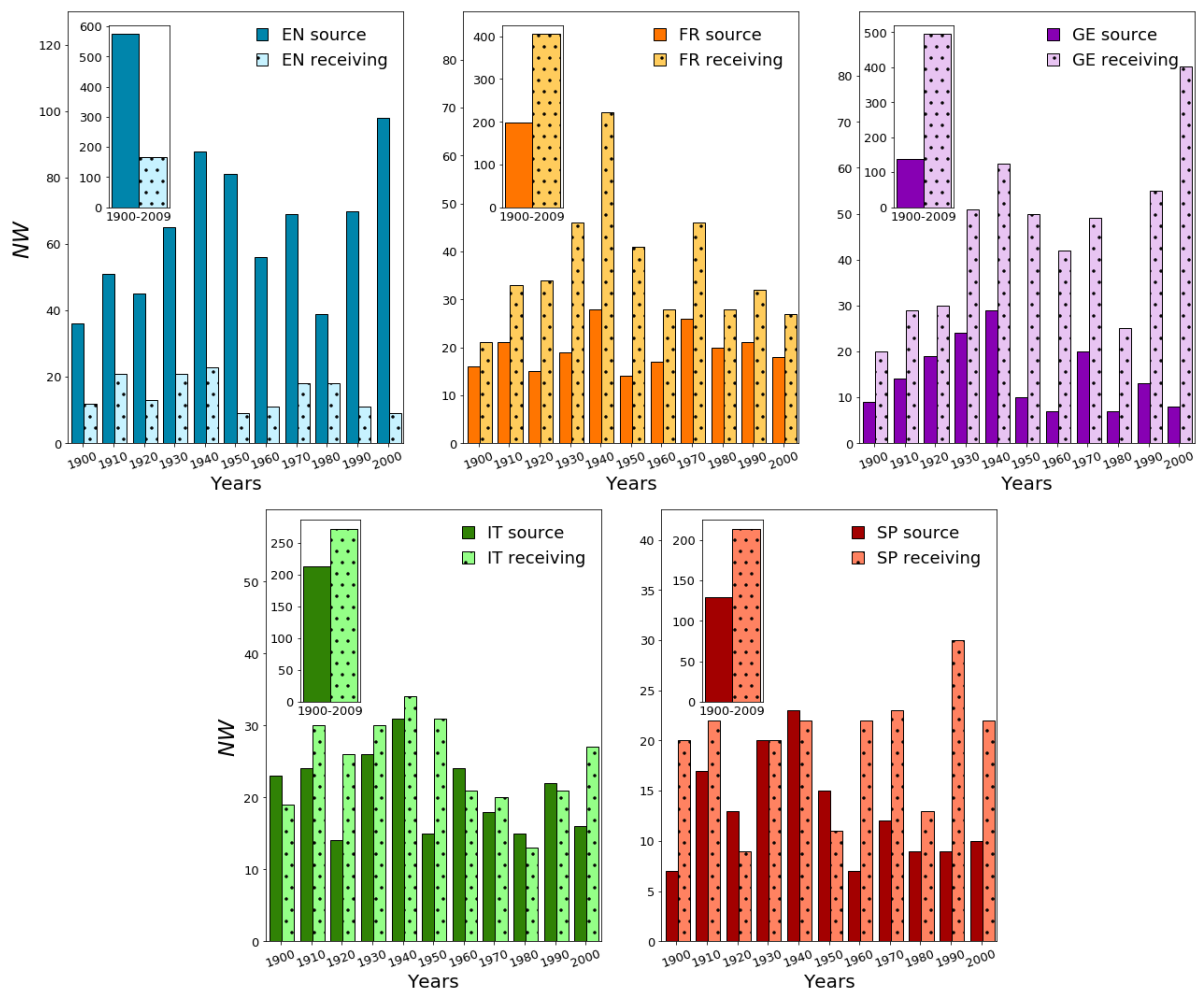


Fig 1. Languages as sources and receivers of new words. During the 20th century, only English language has been the language that has migrated more words than it has received. French, German, Italian and Spanish have received more words during the second half of the century, after the end of World War II.

The Italian language decade by decade, contributes almost the same number of words as those it receives, the only exceptions occurred in the decades of the 1920s and 1950s after the end of the First and Second World War, decades where it was mainly a receiving language. Finally, the Spanish language in 1930 and 1940, also contributed the same number of words as those it received; in the other decades it was primarily a receiving language.

The previous results only manage to show in which decades the languages contributed or received more words, it being evident that half a century after the end of World War II, the English language began to be primarily the source language that contributed the most words to others.

Another way to interpret the previous results is shown in Fig 2. This shows how the new migrant words that each source language contributed are distributed in the different receivers for each decade. To exhibit that migrant words are related by their meaning, some new migrant words that were related to certain semantic fields are display in table 1; specifying the language they come from, to which languages they

migrated and the decades they migrated.

138

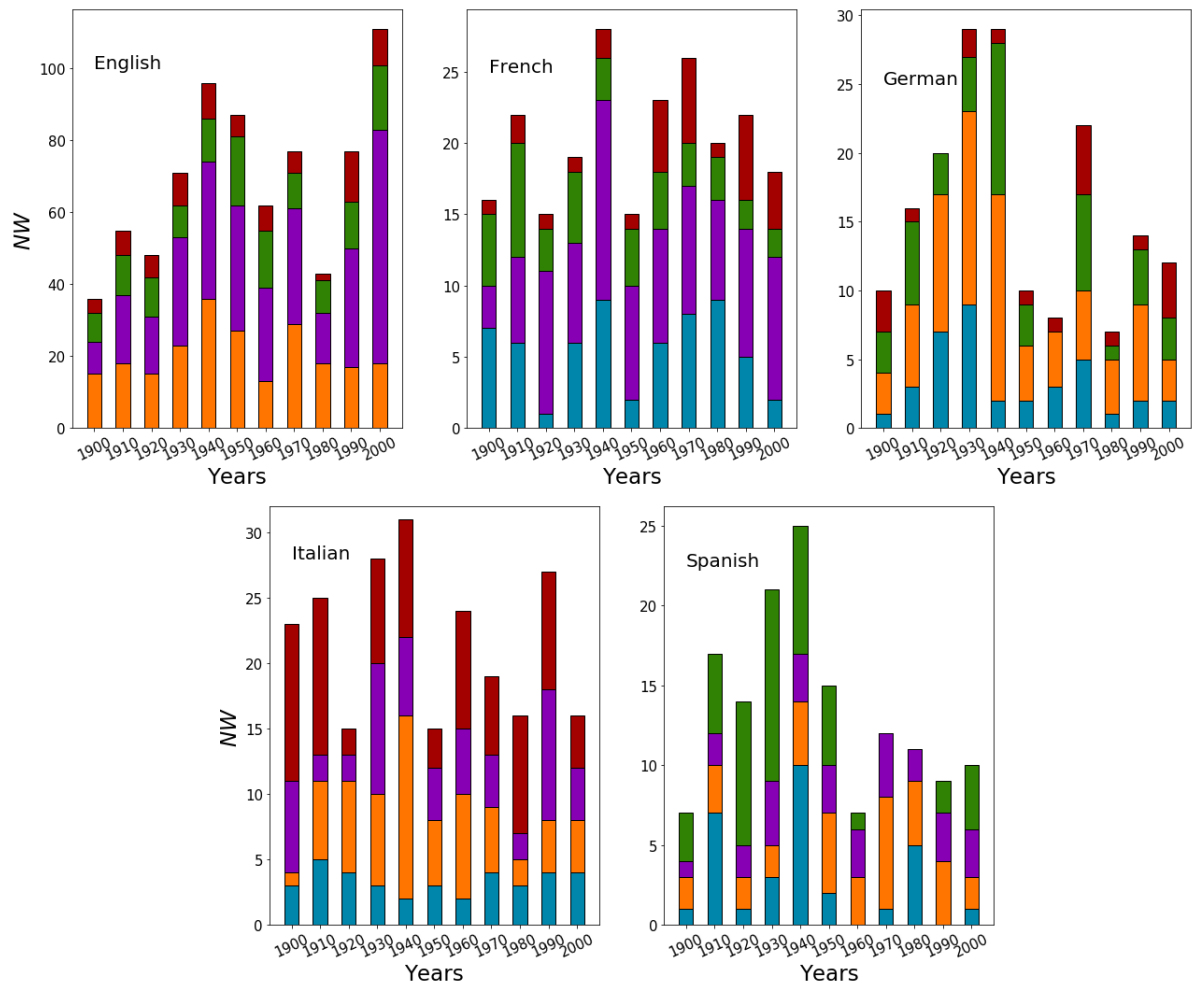


Fig 2. New migrant words that migrate from one source language to the different receiving languages (English ■, French ■, German ■, Italian ■, and Spanish ■).

Table 1. Table

Heading1				Heading2			
cell1row1	cell2 row 1	cell3 row 1	cell4 row 1	cell5 row 1	cell6 row 1	cell7 row 1	cell8 row 1
cell1row2	cell2 row 2	cell3 row 2	cell4 row 2	cell5 row 2	cell6 row 2	cell7 row 2	cell8 row 2
cell1row3	cell2 row 3	cell3 row 3	cell4 row 3	cell5 row 3	cell6 row 3	cell7 row 3	cell8 row 3

Comentarios

The English language as the source language, contributed more words to the French, German and Italian languages, between the 1930s and 1950s. In this period the words found are related to the semantic field of World War II; where countries such as the United States of America, the United Kingdom, France, Italy and Germany were involved. In the last twenty years (1980-2000), the French, German, Italian and Spanish languages have received from English words with content on economics, technology and globalization; we associate this increase with the fact that the English language has been positioned as the common language for communication and for spreading information.

139

140

141

142

143

144

145

146

In the case of the French language, its influence as a source language increased between the 1920s and 1940s, where the main recipient languages were English, French and Italian. The words that migrated to these languages were related to both the World War I and World War II, highlighting exonyms from French to cities, countries and political ideologies that were relevant in such warlike conflicts.

For German as the source language, the largest number of words migrated to other languages between the 1920s and 1940s. The role of Germany in the warlike conflicts of the 20th century, led to the migrations of the semantic field of war in these decades. Another important semantic field that characterizes words with German words, is the important surnames of German-speaking characters who excelled in certain academic areas, such as psychology, philosophy and music; this field shows an influence of the cultural type not only of the German language, but also of the German-speaking academy.

In Italian, the semantic field that characterizes its migrant words is again the World War II, showing an increase of the Italian language between 1920 and 1940. In addition to the words of the semantic field of World War II, in Spanish also migrated words related to political ideologies that were relevant in the aftermath of World War II.

The migrant words from Spanish that reached the other receiving languages, those whose content refers to the field of medicine stand out, these words being the ones that originated the greatest contribution of Spanish in the first half of the 20th century. Finally, names of Latin American countries which suffered from economic crises, were also found.

Among the multiple combinations of source language and recipient language, the most common semantic fields are of a historical nature, highlighting the one referring to the World War II, where words referring to it migrated in all languages.

Accumulated words

The previous results of classifying migrant words according to a semantic field, shows which areas have been influential for there to be a movement of words among languages, nevertheless there is still no answer neither on which language has been more influential over another, nor how the influence occurred.

We define as accumulated migrant words those words with source language A that had already appeared in a receiving language B , and for a given year they do so again. The difference between the new migrant words and the accumulated migrant words, is that they will only be new in the first year of appearance, later on they will become accumulated.

With the accumulated migrant words, we construct a method to obtain a measurable quantity with which to interpret the influence. The procedure to achieve such quantity begins with taking from a year t the list of the five thousand most used words of the recipient language B . Once this is done, we add the frequency $f(k)$ of the five thousand most used words of the receiving language B in year t , where k is the rank of each word

$$\sum_{k=1}^{5000} f(k). \quad (1)$$

Later in the list of words, we distinguish and add the frequency $f(j)$ of those with rank j that are accumulated migrants and that come from the source language A

$$\sum_j f(j). \quad (2)$$

It is called as the **Use** $U_{A \rightarrow B}(t)$ of A in B to the result of normalizing the previous quantity after dividing it by the frequency of the five thousand most used words, given by 1.

$$U_{A \rightarrow B}(t) = \frac{\sum_j f(j)}{\sum_{k=1}^{5000} f(k)}. \quad (3)$$

This new quantity is the one that will quantify the influence of A on B . It will be said that language A has influenced language B more if within a time interval Δt , the use of A in B increases; this will mean that the accumulated words have increased their frequency.

Finally, to quantify the change in use ΔU within a time interval Δt , the term ratio of words per year (wpy) will be used, obtained by dividing ΔU by Δt .

The use among languages

We obtained with the previous method, the accumulated migrant words among all the possible combinations from a source language to a receiving language, in the 269 years (1740-2009) of the dataset. Then in each combination we applied the equation 3 to realize the use in the years between 1900 and 2009. To simplify the reading we establish a set of abbreviations to distinguish the languages, EN for English, FR for French, GE for German, IT for the Italian and SP for Spanish. In addition, to name some combination of source and receiving languages, two abbreviations will be used, the first being the source language of the accumulated migrant words, and the second the receiving language of them, for example EN-GE, means the migrant words from English that are present in the German.

The migrant accumulated words between pairs of language were ordered by the year in which they appeared and in descending order in their frequency, therefore the first place in the ranking is occupied by the most frequent word in that year, the second for the second most frequent, and so on.

The Fig 3 shows the use of a source language, in all receiving languages, as the use in each case may have a different scale, the data were normalized by dividing them by their average value, this with the intention of observing in which time interval the use was more affected (increase or decrease more than in the others), since this will indicate that accumulated words were more or less influential.

English

The use of English in French, Italian and Spanish began to increase after 1930, while in German it was until 1990. We associate the cause of these increases with the emergence of the United States as a world power, after finishing the WWII and impose its economic model as well as the development of science and technology. The accumulated migrant words that are present in all receiving languages are *capital, dollar, investment, relations, market, company, development, financial, institutions, internet, windows and software*. Those can again be associated with semantic fields such as economics, technology, and globalization.

French

The increase in the influence of French in the other languages occurred in English between 1920 and 1970, in German between 1900 and 2009, and in Spanish between 1970 and 1995. In these years the words that increased its frequency are from the semantic fields of religion such as *dieu, évêque, dime, religion, saint* and *église* (iglesia);

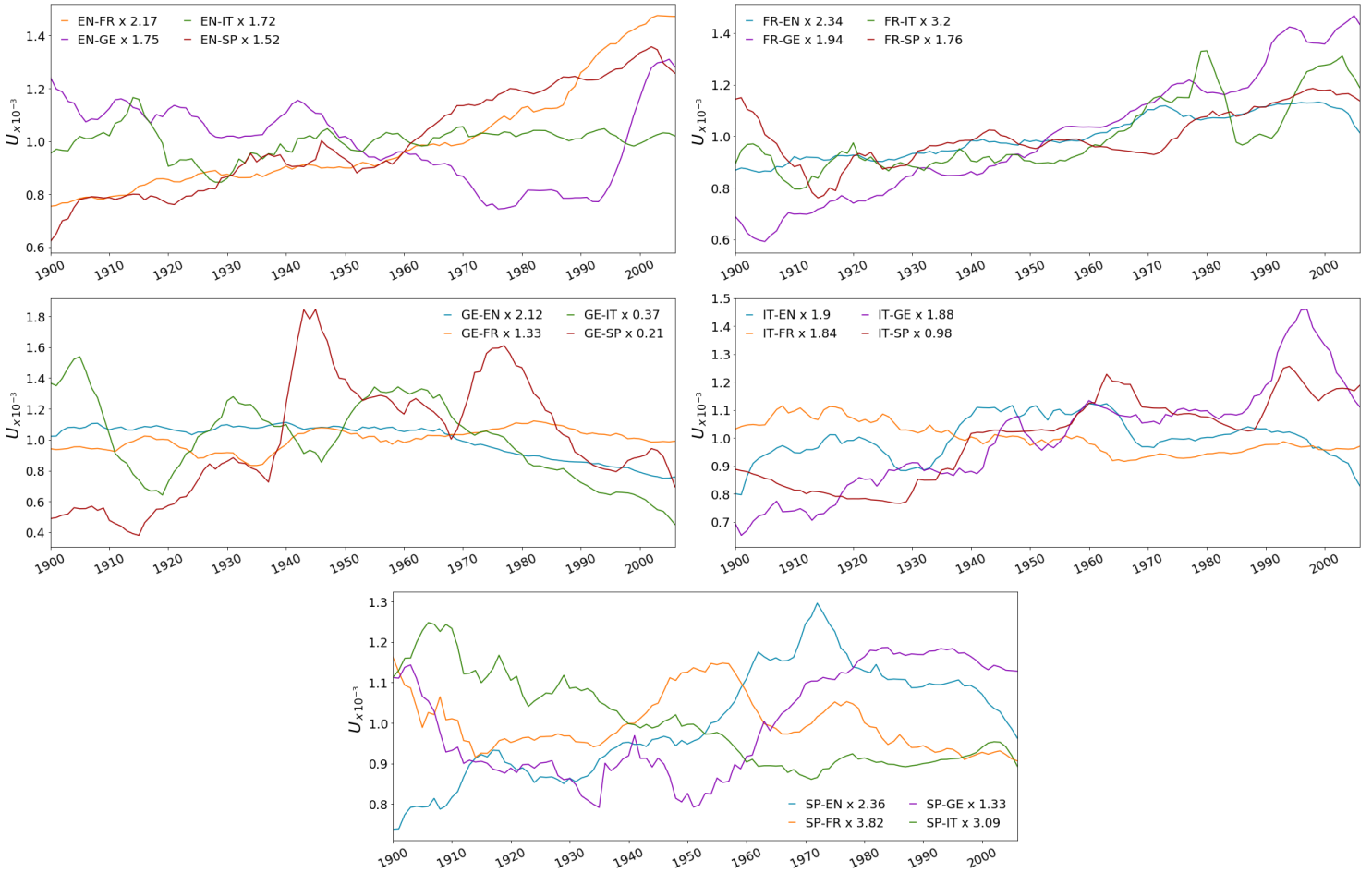


Fig 3. The Use among languages. All language pairs.....

while *reine*, *forteresse*, *napoleon*, *guerre*, *imperiale*, *bastille*, *royal* and *bourgeois* are from French Revolution.

In Italian between 1950 and 1970, in addition to the above words, *raisins*, *vin*, *vignoble* y *recolte* were found, the common meaning of which is the wine industry, a common industry in France and Italy.

German

Spanish between 1935 and 1940 was the language where German had the greatest increase among all receiving languages; followed by the French between 1930 and 1945; in both the words that are present are the surnames of German-speaking characters who excelled in academic areas, such as *marx*, *Freud*, *heidegger*, *nietzsche*, *hegel*, *engels* and *mozart*.

In English and Italian, the biggest change was between 1960 and 2009, where the use of German decreased. In this period the words that lost influence are referred to the World War II, among them *berlin*, *marx*, *hitler*, *lenin*, *testen* and *reich*.

Italian

The influence of Italian came mainly from two semantic fields, the WWII with *mussolini*, *fascismo*, *battaglia*, *regime*, *sociale* and *liberale*; and the religion with *santo*, *suora* (monja) and *cattedrale*. These semantic fields are responsbale for the increase in English between 19630 and 1940, in German between 1950 and 1995, and in Spanish between 1930 and 1960.

In French, Italian has the lowest influence, where the aforementioned words began to be less frequent between 1940 and 1960.

Spanish

The influence of Spanish in English between 1920 and 1970, was due tu historical and cultural facts, since countries in Latin American as *mexico*, *panama*, *chile*, *cuba*, *peru*, *colombia*, *argentina* and its capital *buenos aires*; and states in the United States where a large part of the Spanish-speaking population is concentrated as *california* and *florida*, were increased their frequency.

In German after WWII, and in French between 1930 and 1955, the mainly words involved in that increse are, *terapia*, *anemia*, *lepra*, *tumor*, *syphilis*, *virus* and *renal*.

Rank diversity

Since the accumulated migrant words are organized by year, and at the same time in each year the words are ordered in ascending order in rank, then over time, the same rank can be occupied by different words. One way to quantify how many different elements can occupy the same rank within the same corpus is through rank diversity $d(k)$.

Rank diversity has been used in datasets of the most used words in six Indo-European languages, and in sports and game classifications. Although in each case the criteria for establishing a ranking are different, in both there is a common result: for the same set of data that have had different rankings over time, it is true that the lowest ranks are always occupied by fewer elements, thereby as the range increases, the number of elements that occupy it also does them.

After applying the rank diversity in each source language and receiving language pair, the diversity values resemble a logarithmic curve, as can be seen in Fig 4. We proposed the curve $d(k) = \alpha \ln(k) + \beta$. to fit the rank diversity where the parameters α and β were obtained with a linear regression. With the corresponding linear regression of each language pair, the reliability of the fit with the coefficient of determination R^2 was verified. If a_k represents the value of the adjustment equation when evaluating it in the range k , d_k is the diversity obtained for the same range and \bar{d} is the average of all values of the diversity, then R^2 is expressed as:

$$R^2 = 1 - \sum_{k=1} \frac{(d_k - a_k)^2}{(d_k - \bar{d})^2}. \quad (4)$$

Fig 4 also shows the graph of equation $d(k) = 0.16 \ln(k) - 0.17$, obtained after averaging the coefficients α and β for each language pair. It is observed that the behavior of diversity increases as the rank also increases, regardless of whether the corpus has few or many ranks (14 in German-Spanish, 290 in Spanish-Italian). With this, it can be concluded that, the migrant accumulated words in the middle and high ranks are the ones that tend to change their position the most within a ranking over time.

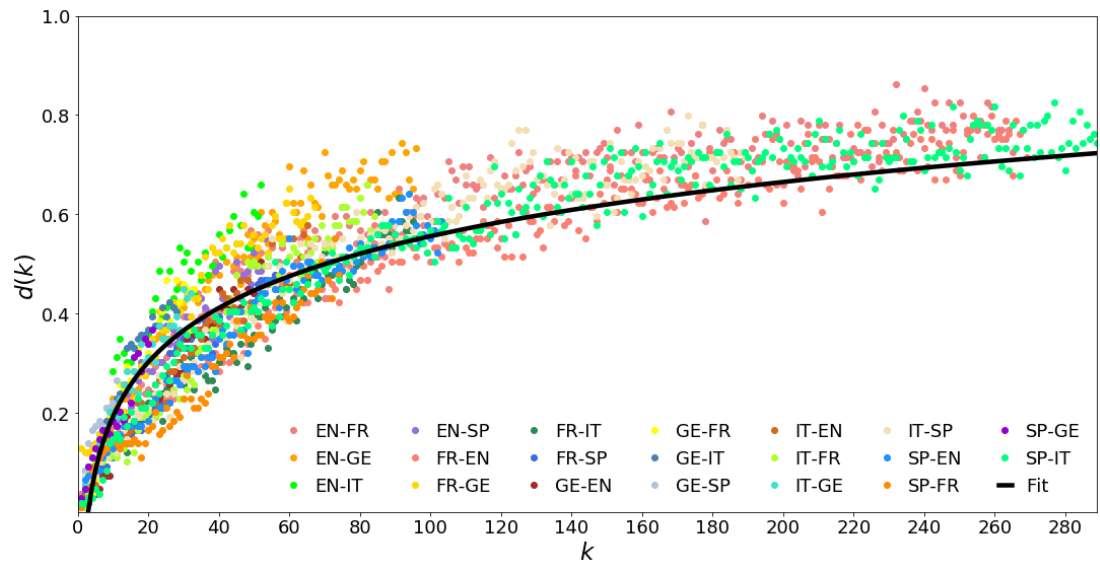


Fig 4. Rank diversity of accumulated migrant words among languages. All language pairs show a logarithmic behavior regardless of the number of ranks where the rank diversity was applied. The reliability of the fit to a logarithmic curve of the 20 combinations is on average $R^2 = 0.88 \pm 0.04$.

Supporting information

289

Acknowledgments

290

References

1. Conant GC, Wolfe KH. Turning a hobby into a job: how duplicated genes find new functions. *Nat Rev Genet.* 2008 Dec;9(12):938–950.
2. Ohno S. *Evolution by gene duplication.* London: George Alien & Unwin Ltd. Berlin, Heidelberg and New York: Springer-Verlag.; 1970.
3. Magwire MM, Bayer F, Webster CL, Cao C, Jiggins FM. Successive increases in the resistance of *Drosophila* to viral infection through a transposon insertion followed by a Duplication. *PLoS Genet.* 2011 Oct;7(10):e1002337.