

Supplementary Information Appendix for *Links that speak: The global language network and its association with global fame*

Shahar Ronen, Bruno Gonçalves, Kevin Z. Hu, Alessandro Vespignani,
Steven Pinker, César A. Hidalgo

Supplementary online material (SOM) and additional
visualizations are available on <http://language.media.mit.edu>

Table of Contents

S1 Data	2
S1.1 Twitter	2
S1.2 Wikipedia	4
S1.3 Book translations	7
S2 Language notation and demographics	8
S2.1 Notation.....	8
S2.2 Population	9
S2.3 Language GDP	10
S3 Additional calculations	11
S4 Language centrality: Eigenvector centrality vs. betweenness centrality	12
S5 Famous people per language	14
S5.1 Associating a famous person with languages	14
S5.2 Wikipedia	18
S5.3 Human Accomplishment	22
S5.4 Comparison of the famous people datasets	24
References for the SI Appendix	26

S1 Data

S1.1 Twitter

Twitter is a microblogging and online social networking service where users communicate using text messages of up to 140 characters long called *tweets*. As of December 2012, Twitter had over 500 million registered users from all over the world, tweeting in many different languages. Of these, 200 million users were active every month (1).

Tweets are attributed to their authors and can be used to identify polyglots and the language communities they connect, making Twitter a good source for representing the GLN of tens of millions of people. Registered Twitter accounts make up for 7% of world population, but its demographics may not reflect real-life demographics (2). For example, Twitter users in the United States are younger and hold more liberal opinions than the general public (3).

We collected 1,009,054,492 tweets between December 6, 2011 and February 13, 2012, through the Twitter *garden hose*, which gives access to 10% of all tweets. We detected the language of each tweet using the Chromium Compact Language Detector (CLD) (4), which was chosen for its wide language support and its relatively accurate detection of short messages (5, 6). However, any automated language detection is prone to errors (7), all the more so when performed on short, informal texts such as tweets. To reduce the effect of such errors, we applied the following methods.

Firstly, to improve detection, we removed *hashtags* (marks of keywords or topics, which start with a #), URLs, and *@-mentions* (references to usernames, which start with a @). Hashtags, URLs and *@-mentions* are often written in English or in another Latin script, regardless of the actual language of the tweet, and may mislead the detector.

Secondly, we used only tweets that CLD detected with a high degree of confidence. CLD suggests up to three possible languages for the text detected, and gives each option a score that indicates its certainty of the identification, 1 being the lowest and 100 being the highest. If the top option has a much higher score than the other options, CLD marks the identification as *reliable*. We only used tweets that CLD was able to detect with a certainty

over 90% and indicated a reliable detection. The 90% threshold was chosen as the optimal tradeoff between detection accuracy and number of tweets detected, based on a sample of 1 million tweets (see Figure S1A).

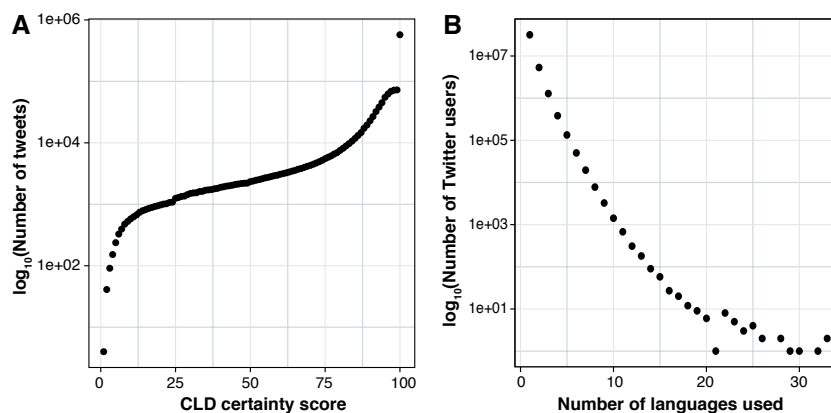


Figure S1 **A** number of tweets as function of certainty **B** Distribution of Twitter users by number of languages in which they tweet.

Thirdly, as mutually intelligible languages are difficult to distinguish, we merged similar languages. To do so, we converted the two-letter ISO 639-1 language codes (8) produced by CLD to three-letter ISO 639-3 codes (9), and merged them using the ISO 639-3 macrolanguages standard. See **Section S2.1** for further details and limitations.

Finally, to reduce the effect of individual detection errors, we considered for each user only languages in which he or she tweeted at least twice, and considered only users who made at least five tweets overall. We found that a large number of users tweeted in a relatively large number of languages, and we attribute some of this to inaccurate language detection. To prevent this from skewing the representation of the Twitter GLN, we discarded users who tweeted in more than five languages (Figure S1B). Five was chosen as the cutoff based on the impression of linguist Richard Hudson that five languages were the most spoken in a community; he coined the term *hyper-polyglots* for people who speak six languages or more (10). Some of these users might be bots, which are common on Twitter. Note however that multilingual Twitter bots are not considered a common phenomenon, and even if they were, a bot reading news in one language and re-tweeting them in another is certainly an indication of interaction between the two languages.

After applying the criteria listed above, we had a dataset of 548,285,896 tweets in 73 languages by 17,694,811 users, which is available on the SOM site. We used this dataset to

generate the Wikipedia GLN shown in Figure 1 of the main section. **Table S1** shows statistics for the languages with the most tweets in our Twitter dataset.

#	Language	Code	Tweets	Users	Tweets per user	% of total users
1	English	eng	255,351,176	10,859,465	23.5	61.37%
2	Japanese	jpn	91,669,691	2,602,426	35.2	14.71%
3	Malay	msa	49,546,710	1,651,705	30	9.33%
4	Portuguese	por	46,520,572	1,617,409	28.8	9.14%
5	Spanish	spa	44,195,979	2,043,468	21.6	11.55%
6	Korean	kor	11,674,755	289,982	40.3	1.64%
7	Dutch	nld	10,526,980	435,128	24.2	2.46%
8	Arabic	ara	9,993,172	366,643	27.3	2.07%
9	Thai	tha	7,449,790	154,171	48.3	0.87%
10	Turkish	tur	4,660,694	233,158	20	1.32%
11	Russian	rus	4,577,942	243,159	18.8	1.37%
12	French	fra	3,434,065	147,843	23.2	0.84%
13	Filipino	fil	1,905,619	257,611	7.4	1.46%
14	German	deu	1,705,256	73,897	23.1	0.42%
15	Italian	ita	1,586,225	89,242	17.8	0.50%
16	Swedish	swe	596,130	36,604	16.3	0.21%
17	Modern Greek	ell	526,527	30,609	17.2	0.17%
18	Chinese	zho	453,837	24,113	18.8	0.14%
19	Catalan	cat	236,424	32,376	7.3	0.18%
20	Norwegian	nor	170,430	16,500	10.3	0.09%

Table S1 Statistics for the twenty languages with the most tweets in our Twitter dataset. The full table is available on the SOM.

S1.2 Wikipedia

Wikipedia is a multilingual, web-based, collaboratively edited encyclopedia. As of March 2013, Wikipedia had 40 million registered user accounts across all language editions, of which over 300,000 actively contributed on a monthly basis (11). Wikipedia’s single sign-on mechanism lets editors use the same username on all language editions to which they contribute. This allows us to associate a contribution with a specific person and identify the languages spoken by that person.

We compiled our Wikipedia dataset as follows. Firstly, we collected information on editors and their contributions in different languages from the edit logs of all Wikipedia editions until the end of 2011. We collected only edits to proper articles (as opposed to user pages or talk pages), and only edits made by human editors. Edits by bots used by Wikipedia for basic maintenance tasks (e.g., fixing broken links, spellchecking, adding references to other pages) were ignored, as many of them make changes in an unrealistic

number of languages, potentially skewing the GLN. This initial dataset contained 643,435,467 edits in 266 languages by 7,344,390 editors.

Secondly, we merged the languages as we did for the Twitter dataset, discarding ten Wikipedia editions in the process. Two of them are more or less duplicates of other editions, namely *simple* (Simple English) of English and *be-x-old* (Classic Belarusian) of Official Belarusian. The remaining eight could not be mapped to standard ISO639-3 languages: *bh*, *cbk_zam*, *hz*, *map_bms*, *nah*, *nds_nl*, *tokipona*, *roa_tara*. These eight editions are small and contain together 220,575 edits by 318 contributors.

Finally, to reduce the effect of one-time edits, which may be cosmetic or technical and may not indicate knowledge of a language, we set the same thresholds as for our Twitter dataset. For each user we considered only languages in which he or she made at least two edits, and considered only users who made at least five edits overall. We also discarded editors who contributed to more than five languages, following the rationale explained in the Twitter section. We did so because a large number of users contributed to an unrealistic number of languages: hundreds of users contributed to over 50 language editions each, and dozens edited in over 250 languages each (see Figure S2). For example, one of the users we identified was a self-reported native speaker of Finnish (contributed 6,787 edits to this edition by the end of 2011), and an intermediate speaker of English (834 edits) and Swedish (20 edits). However, this user contributed to ten additional language editions, in particular Somali (149 edits) and Japanese (58 edits). Most of these contributions are maintenance work that does not require knowledge of the language, such as the addition of a redirection or the reversion of changes.

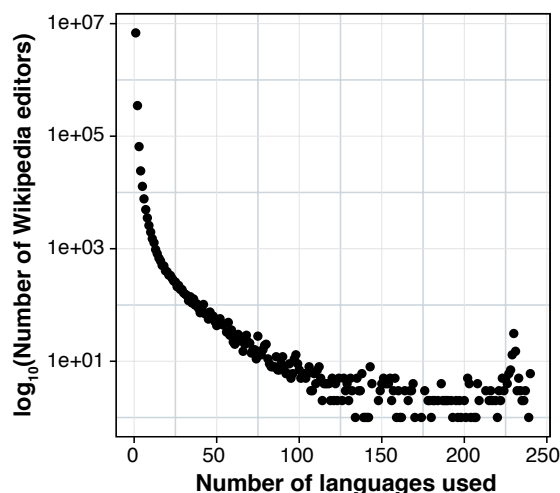


Figure S2 Distribution of Wikipedia editors by number of languages in which they contribute.

Table S2 below shows statistics for the languages with the most edits in our dataset. The final dataset consists of 382,884,184 edits in 238 languages by 2,562,860 contributors, and is available on the SOM site. We used this dataset to generate the Wikipedia GLN shown in Figure 1 of the main section.

#	Language	Code	Edits	Editors	Edits per user	% of total editors
1	English	eng	198,361,048	1,589,250	124.81	62.011%
2	German	deu	33,977,378	224,215	151.54	8.749%
3	French	fra	23,070,757	142,795	161.57	5.572%
4	Japanese	jpn	16,149,315	102,857	157.01	4.013%
5	Spanish	spa	13,645,596	145,487	93.79	5.677%
6	Russian	rus	12,445,887	81,925	151.92	3.197%
7	Italian	ita	11,923,658	72,981	163.38	2.848%
8	Chinese	zho	7,302,770	50,341	145.07	1.964%
9	Polish	pol	6,589,015	47,015	140.15	1.834%
10	Dutch	nld	6,393,791	46,951	136.18	1.832%
11	Hebrew	heb	5,467,149	18,998	287.77	0.741%
12	Portuguese	por	5,168,734	60,487	85.45	2.360%
13	Swedish	swe	3,521,224	30,498	115.46	1.190%
14	Finnish	fin	2,926,115	20,811	140.60	0.812%
15	Hungarian	hun	2,713,725	18,033	150.49	0.704%
16	Korean	kor	2,634,092	16,464	159.99	0.642%
17	Arabic	ara	2,178,719	18,258	119.33	0.712%
18	Turkish	tur	2,062,037	23,926	86.18	0.934%
19	Serbo-Croatian	hbs	2,030,039	10,901	186.23	0.425%
20	Ukrainian	ukr	1,839,988	10,028	183.49	0.391%

Table S2 Statistics for the twenty languages with the most edits in our Wikipedia dataset. The full table is available on the SOM site.

S1.3 Book translations

The Index Translationum is an international bibliography of book translations maintained by UNESCO (12). The online database contains information on books translated and published in print in about 150 countries since 1979. Some countries are missing data for certain years, such as the United Kingdom in the years 1995-2000 and 2009-2011 (13).

We retrieved a dump of the data on July 22, 2012, which contained 2,244,527 translations in 1,160 languages. After removing a few corrupt entries, we converted the language codes listed in the Index Translationum to standard three-letter ISO639-3 codes. The following entries were discarded from the dataset: 41 miscellaneous dialects of languages that were already listed (together accounting for under 100 translations total), 46 languages that could not be mapped to standard ISO639-3 codes (together accounting for about a thousand translations total), and 5 administrative codes (*mis*, *mul*, *und*, *zxx*, and *not supplied*; see ISO639-3 documentation (9)). The remaining languages were merged into macrolanguages (see **Section S2.1**).

Table S3 shows statistics for the languages with the most translations in our dataset. The final dataset contains 2,231,920 translations in 1,019 languages. We used this dataset to generate the book translations GLN shown in Figure 1 of the main section.

#	Language	Code	Translations from	Translations to	Total translations
1	English	eng	1,225,237	146,294	1,371,531
2	German	deu	201,718	292,124	493,842
3	French	fra	216,624	238,463	455,087
4	Spanish	spa	52,955	228,910	281,865
5	Russian	rus	101,395	82,772	184,167
6	Japanese	jpn	26,921	130,893	157,814
7	Dutch	nld	18,978	111,371	130,349
8	Italian	ita	66,453	59,830	126,283
9	Swedish	swe	39,192	71,688	110,880
10	Polish	pol	14,104	76,720	90,824
11	Portuguese	por	11,390	74,721	86,111
12	Danish	dan	21,239	64,799	86,038
13	Czech	ces	17,202	64,442	81,644
14	Chinese	zho	13,337	62,650	75,987
15	Hungarian	hun	11,256	54,989	66,245
16	Norwegian	nor	14,530	45,923	60,453
17	Serbo-Croatian	hbs	12,743	45,036	57,779
18	Finnish	fin	8,296	46,271	54,567
19	Modern Greek (1453-)	ell	4,862	27,422	32,284
20	Bulgarian	bul	3,667	25,742	29,409

Table S3 Statistics for the twenty languages with the most translations (to and from) in our Index Translationum dataset. The full table is available on the SOM site.

S2 Language notation and demographics

S2.1 Notation

Each of our three datasets uses a different system for identifying language names. For the sake of consistency, we converted the language identifiers to ISO 639-3 identifiers. ISO 639-3 is a code that aims to define three-letter identifiers for all known human languages (9). For example, English is represented as *eng*, Spanish as *spa*, Modern Greek as *ell* and Ancient Greek as *grc*.

Some languages are *mutually intelligible* or nearly mutually intelligible with others, such as Serbian and Croatian, Indonesian and Malaysian, and the various regional dialects of Arabic. Because of the similarity of mutually intelligible languages we do not consider their speakers as polyglots. Instead, we merged mutually intelligible languages to macrolanguages following the ISO 639-3 Macrolanguage Mappings (9). For example, we merged 29 varieties of Arabic into one Arabic macrolanguage (*ara*), and Malaysian, Indonesian, and 34 other Bhasa languages into a Malay macrolanguage (*msa*).

Another reason for consolidating languages is that the language detector we used to identify the language of tweets cannot distinguish between the written forms of many mutually intelligible languages, such as Indonesian and Malaysian and Serbian and Croatian. For this reason, we added a couple of merges that are not in the ISO 639-3 macrolanguage mappings: we consolidated Serbian, Croatian, and Bosnian into Serbo-Croatian (*hbs*) even though the latter had been deprecated as a macrolanguage, and merged Tagalog (*tgl*) with Filipino (*fil*) into one Filipino language that uses the identifier *fil*. Our full conversion table is available on the SOM site.

Languages belong to language families (14). We mapped languages to language families using the hierarchy in Ethnologue (15) complemented by information from articles from the English Wikipedia about the respective languages. We used the standard language family names and identifiers as defined by ISO 639-5 (16).

S2.2 Population

We use language speaker estimates from the June 14, 2012 version of Wikipedia Statistics page (17). These estimates include all speakers of a language, native and non-native alike. We converted language names to ISO 639-3 identifiers and merged them into macrolanguages as explained in **Section S2.1**.

In general, the number of speakers of a macrolanguage is the sum of speakers of its constituent languages. However, for the macrolanguages listed in Table S4 we determined that the estimated number of speakers for one of the individual languages that constitute them includes speakers of the other languages, and used that number as the speaker estimate for the entire macrolanguage. Refer to Table S5 for number of speakers for languages in our datasets.

Macrolanguage	ISO 639-3 identifier	Speaker estimate we use in our dataset	Individual languages according to Wikipedia (Wikipedia language code)	Wikipedia Statistics speaker estimate
Akan	aka	19 million	Akan (ak) Twi (tw)	19 million 15 million
Arabic	ara	530 million	Arabic (ar) Egyptian Arabic (arz)	530 million 76 million
Malay	msa	300 million	Malay (ms) Indonesian (id)	300 million 250 million
Serbo-Croatian	hbs	23 million	Serbo-Croatian (sh) Serbian (sr) Croatian (hr) Bosnian (bs)	23 million 23 million 6 million 3 million
Norwegian	nor	5 million	Norwegian (no) Nynorsk (nn)	5 million 5 million
Komi	kom	293,000	Komi (kv) Komi-Perniak (koi)	293,000 94,000

Table S4 Macrolanguages for which the estimated number of speakers is not the sum of the estimates for the individual languages that constitute them.

	Language	Code	Speakers (millions)	GDP per capita (\$)
1	Afrikaans	afr	13	10,373
2	Albanian	sqi	16	9,182
3	Arabic	ara	530	8,720
4	Armenian	hye	6	5,598
5	Azerbaijani	aze	27	11,902
6	Bashkir	bak	2	
7	Basque	eus	1	30,626
8	Belarusian	bel	6	15,028
9	Bengali	ben	230	2,457
10	Bulgarian	bul	12	13,488
11	Catalan	cat	9	30,626
12	Chinese	zho	1575	9,207
13	Czech	ces	12	27,062
14	Danish	dan	6	37,152
15	Dutch	nld	27	40,518
16	English	eng	1500	32,953
17	Esperanto	epo	1	
18	Estonian	est	1.07	20,380
19	Filipino	fil	90	4,073
20	Finnish	fin	6	36,236
21	French	fra	200	15,103
22	French (Old)	fro		
23	Galician	glg	4	30,626
24	Georgian	kat	4	5,491
25	German	deu	185	38,268
26	German (Middle High)	gmh		
27	Greek (Ancient)	grc		
28	Greek (Modern)	ell	15	26,693
29	Haitian	hat	12	1,235
30	Hebrew	heb	10	30,975
31	Hindi	hin	550	3,696
32	Hungarian	hun	15	18,672
33	Icelandic	isl	0.32	38,061
34	Italian	ita	70	30,623
35	Japanese	jpn	132	34,740
36	Kara-Kalpak	kaa	0.41	
37	Kazakh	kaz	12	13,001
38	Kirghiz	kir	5	2,372
39	Korean	kor	78	21,723
40	Latin	lat	0.01	

	Language	Code	Speakers (millions)	GDP per capita (\$)
41	Latvian	lav	2.15	15,662
42	Lithuanian	lit	4	18,856
43	Macedonian	mkd	3	10,367
44	Malay	msa	300	6,023
45	Malayalam	mal	37	3,694
46	Maltese	mlt	0.37	25,428
47	Maori	mri	0.157	27,668
48	Marathi	mar	90	3,694
49	Moldavian	mol	3.5	
50	Mongolian	mon	5	4,744
51	Norwegian	nor	5	53,471
52	Occitan	oci	2	
53	Persian	fas	107	9,826
54	Polish	pol	43	20,326
55	Portuguese	por	290	11,853
56	Romanian	ron	28	11,354
57	Russian	rus	278	15,487
58	Sanskrit	san	0.05	
59	Serbo-Croatian	hbs	23	12,908
60	Sinhala	sin	19	5,674
61	Slovak	slk	7	23,432
62	Slovenian	slv	2	28,642
63	Spanish	spa	500	16,777
64	Swahili	swa	50	1,415
65	Swedish	swe	10	40,265
66	Tajik	tgk	4	2,238
67	Tamil	tam	66	3,923
68	Tatar	tat	8	
69	Thai	tha	73	9,396
70	Tibetan	bod	7	
71	Turkish	tur	70	14,623
72	Turkmen	tuk	9	5,816
73	Uighur	uig	10	
74	Ukrainian	ukr	45	7,242
75	Urdu	urd	60	3,511
76	Uzbek	uzb	24	3,182
77	Vietnamese	vie	80	3,447
78	Welsh	cym	0.75	
79	Yiddish	yid	3	

Table S5 Population and GDP per capita for the languages used in the GLNs. Blank cells indicate dead languages or insufficient data.

S2.3 Language GDP

The GDP (*gross domestic product*) per capita for a language l measures the average contribution of a single speaker of language l to the world GDP, and is calculated by summing the contributions of speakers of l to the GDP of every country, and dividing the sum by the number of speakers of l . A similar method was used by Davis (18). Given a country c , let G_c be the GDP per capita (based on purchasing-power-parity) of that country (2011 values; retrieved from the IMF (18) with a few additions from the CIA World Factbook (19)). Also, given a language l , let N_{lc} be the number of native speakers of l in country c ,

obtained from Ethnologue (15) and The World Factbook (19). We calculated N_{lc} using the language demographics listed in Table S6. Thus, G_l , the GDP per capita for l is

$$G_l = \frac{\sum_c (G_c N_{lc})}{\sum_c N_{lc}}$$

The GDP per capita values in Table S5 are approximate, because the economic activity of a country is not distributed evenly by language. Moreover, a person may contribute in a language different than his or her native language: for example, many use English to communicate at their workplace although English is not their native language. Tables of GDP per capita and population by country and language are available on the SOM site.

S3 Additional calculations

In this section we briefly document two calculations used in the main text of the paper. First, we note that for all figures we use the number of multilingual speakers, or expressions, from a language. We estimate the number of multilingual speakers or expression from a language (N_i) as:

$$N_i = \sum_j M_{ij}$$

Also, we note that we estimate the eigenvector centrality of a language by using:

$$\lambda v_i = \sum_j M_{ij} v_j$$

and finding the eigenvector v , associated with the largest eigenvalue. Since the eigenvector associated with the largest eigenvalue could be positive or negative, we take the absolute value of the elements of this eigenvector as our measure of a language's eigenvector centrality.

S4 Language centrality: Eigenvector centrality vs. betweenness centrality

In this section we compare two measures of centrality, *eigenvector centrality* (the metric used in the main text) and *betweenness centrality*. The betweenness centrality of a node is the number of shortest paths from all nodes to all others that pass through that node³¹. This centrality value focuses on quantity rather than quality: all shortest paths that go through a node contribute equally to its betweenness score, regardless of the characteristics of the source and target nodes (e.g., the number of their neighbors or their identity). The eigenvector centrality of a node is the sum of its summed connections to others, weighted by their centralities (20). Eigenvector centrality thus takes into account the quality of a node's connections, by rewarding a node for being connected to "important" nodes. Each node is assigned a relative score based on its connections, and a connection to a high-scoring node contributes more to the eigenvector centrality score of the node being scored than a connection to a low-scoring node.

Figure S3 shows the correlation of eigenvector centrality and betweenness centrality for all languages and datasets. The correlation between the two centrality measures is $R^2=0.25$ for Twitter, $R^2=0.62$ for Wikipedia, and $R^2=0.39$ for book translations. A table with eigenvector and betweenness centralities of each language in the Twitter, Wikipedia and book translation GLNs is available on the SOM site.

The deviations between these two centrality measures are quite informative. For instance, according to betweenness centrality the most central language in the book translations GLN is Russian. Figure 1 in the main text shows why: Russian is the portal to a large number of languages that would otherwise be disconnected from the rest of the network (such as Tatar, Armenian and Kirghiz). All paths to these languages pass through Russian, contributing to Russian's high betweenness score. The same is not true for English, the language with the second-highest betweenness. While English is also highly connected, it is connected to many languages that are connected to others, and is thus located in a part of the network where there are alternative paths that reduce the betweenness of English. At the same time, the fact that English is connected to languages that are connected to others increases its eigenvector centrality.

We chose eigenvector centrality over betweenness, as the former is more suitable for identifying global languages according to our definition: a global language is a language that are connected to other hub languages (such as English in the example from the book translations network above), not a language that serve as the only gateway to many peripheral languages (such as Russian in the above example).

We also had a practical reason for preferring eigenvector centrality to betweenness centrality: the latter is a measure that is unable to differentiate among more peripheral languages, since most languages get a betweenness score of zero (see Figure S3). Eigenvector centrality, on the other hand, can help us differentiate between the positions of languages in the GLN at all levels of centrality, not only among the most central languages.

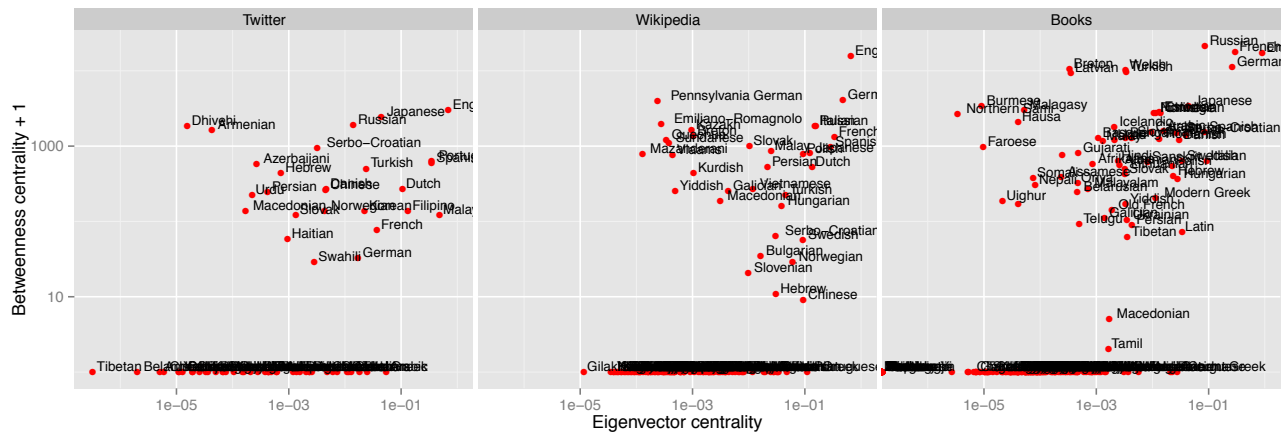


Figure S3 Comparison between eigenvector centrality and betweenness centrality, calculated as the total number of paths going through a node, for **A** The Twitter GLN **B** The Wikipedia GLN **C** The book translations GLN.

S5 Famous people per language

We measure the cultural impact of a language by the number of its speakers who made a long-lasting cultural impression on the world. We focus on these *famous* people, rather than on ideas or other forms of cultural expression, because people names are easier to identify and match across languages.

We use two separate methods to decide whether a person is famous. The first is having Wikipedia articles in at least 26 language editions, and the second is being included in the *Human Accomplishment* list (21), a list of nearly 4,000 influential people in the arts and sciences, from 800 BCE to 1950. As neither dataset contains information about the language used by the famous people it lists, we start this section by describing how we associated famous people with languages. Then, we dedicate a subsection to each dataset, in which we describe how the dataset was retrieved and prepared for use.

S5.1 Associating a famous person with languages

Ideally each language would be given a point for each famous person who spoke this language as his or her native language, or who used this language as the main language for his or her main contributions. Unfortunately, this information is not available in a structured format and finding it manually for each person does not scale well for thousands of people. Therefore, we determined a person's language affiliation using the current language demographics for his or her country of birth. Each famous person in our datasets equals one point, which is distributed across the languages spoken in his or her native country according to their population (15, 19). For example, Italian inventor Guglielmo Marconi counts as one point for Italian. Former Canadian Prime Minister Pierre Trudeau contributes 0.59 to English and 0.22 to French. We stress again that our scoring is based on national identity and not on cultural or linguistic identity. Trudeau was a native speaker of French while Leonard Cohen is a native speaker of English, but since both of them are Canadian, each one adds 0.59 points for English and 0.22 points for French, regardless of their native language. Refer to Table S6 for the language demographics of each country.

We determine a person's country of birth using present-day international borders. For example, we code Italy as the country of birth for author Ippolito Nievo, although Italy was unified only shortly before his death in 1861 and at the time of his birth his native Padua was part of the Austrian Empire. In some cases, this method produces unintuitive results. The Ancient Greek historian Herodotus was born in Halicarnassus (present-day Bodrum, Turkey) and would earn points for Turkish, while Mustafa Kemal Atatürk, founder of the Republic of Turkey, was born in Thessaloniki, present-day Greece, and would earn points for Greek. Because our language distribution statistics are from the last few years, we include only people born in 1800 and later, to reduce the effect of geopolitical and cultural changes on our mapping of countries to languages. To match the year limitation of the Human Accomplishment dataset, we also set 1950 as the latest year of birth for the Wikipedia dataset.

Despite some inaccuracies, using present-day countries provides a consistent mapping of people who lived over a period of several millennia to their contemporary countries. Moreover, using present-day countries allows us to use the present-day language distribution statistics for each country to identify the main languages spoken in a country and determine the language affiliation of each person.

1	Afghanistan	Persian 50%, Pushto 35%, Uzbek 6%, Turkmen 5%	26	Brunei	Malay 100%	51	Ecuador	Spanish 100%	76	Guinea-Bissau	Upper Guinea Crioulo 44%, Portuguese 14%
2	Albania	Albanian 95%, Greek (Modern) 3%	27	Bulgaria	Bulgarian 76.8%, Turkish 8.2%, Romany 3.8%	52	Egypt	Arabic 100%	77	Guyana	English 50%
3	Algeria	Arabic 80%, French 20%	28	Burkina Faso	French 100%	53	El Salvador	Spanish 100%	78	Haiti	Haitian 75%, French 25%
4	Andorra	Catalan 40%, Spanish 35%, Portuguese 15%, French 5.5%	29	Burma	Burmese 100%	54	Equatorial Guinea	Spanish 67.6%, French 20%	79	Honduras	Spanish 100%
5	Angola	Portuguese 70%	30	Burundi	French 50%, Rundi 50%	55	Eritrea	Tigrinya 55%, Tigre 16%, Estonian 67.3%, Russian 29.7%, Oromo 33.8%, Amharic 29.3%, Somali 6.2%, Tigre 5.9%, Sidamo 4%	80	Hong Kong	Chinese 95%, English 3.5%
6	Argentina	Spanish 98%	31	Cambodia	Central Khmer 95%	56	Estonia	Estonian 67.3%, Russian 29.7%, Oromo 33.8%, Amharic 29.3%, Somali 6.2%, Tigre 5.9%, Sidamo 4%	81	Hungary	Hungarian 93.6%
7	Armenia	Armenian 97.7%, Russian 0.9%	32	Cameroon	French 50%, English 50%	57	Ethiopia	Estonian 67.3%, Russian 29.7%, Oromo 33.8%, Amharic 29.3%, Somali 6.2%, Tigre 5.9%, Sidamo 4%	82	Iceland	Icelandic 100%
8	Aruba	Papiamentu 66.3%, Spanish 12.6%, English 7.7%, Dutch 5.8%	33	Canada	English 58.8%, French 21.6%	58	Faroe Islands	Faroese 100%	83	India	Hindi 41%, Bengali 8.1%, Telugu 7.2%, Marathi 7%, Tamil 5.9%, Urdu 5%, Gujarati 4.5%, Kannada 3.7%, Oriya 3.2%, Malayalam 3.2%, Panjabi 2.8%
9	Australia	English 78.5%, Chinese 2.5%, Italian 1.6%, Greek (Modern) 1.3%, Arabic 1.2%, Vietnamese 1%	34	Cape Verde	Portuguese 100%	59	Fiji	Fiji Hindi 45.3%, Fijian 39.3%	84	Indonesia	Malay 100%
10	Austria	German 88.6%, Serbo-Croatian 3.8%, Turkish 2.3%	35	Central African Republic	Sango 80%, French 20%	60	Finland	Finnish 91.2%, Swedish 5.5%	85	Iran	Persian 53%, Azerbaijani 18%, Kurdish 10%, Luri 6%, Arabic 2%
11	Azerbaijan	Azerbaijani 90.3%, Lezghian 2.2%, Russian 1.8%, Armenian 1.5%	36	Chad	Arabic 50%, French 50%	61	France	French 100%	86	Iraq	Arabic 80%, Kurdish 15%
12	Bahamas, The	English 100%	37	Chile	Spanish 100%	62	French Guiana	French 100%	87	Ireland	English 95%, Irish 2%
13	Bahrain	Arabic 100%	38	China	Chinese 100%	63	Gabon	French 75%, Fang 25%	88	Ile of Man	English 100%
14	Bangladesh	Bengali 98%	39	Colombia	Spanish 100%	64	Gambia, The	English 100%	89	Israel	Hebrew 80%, Arabic 15%
15	Barbados	English 100%	40	Congo, Democratic Republic of the	French 33%, Swahili 20%, Lingala 20%	65	Georgia	Georgian 71%, Russian 9%, Armenian 7%, Azerbaijani 6%	90	Italy	Italian 100%
16	Belarus	Russian 70.2%, Belarusian 23.4%	41	Congo, Republic of the	French 30%, Ibali Teké 17%, Lingala 13%	66	Germany	German 100%	91	Jamaica	English 100%
17	Belgium	Dutch 60%, French 40%	42	Costa Rica	Spanish 100%	67	Ghana	Akan 24.7%, English 21.3%, Ewe 12.7%, Abon 4.6%	92	Japan	Japanese 100%
18	Belize	English 41%, Spanish 32%	43	Cote d'Ivoire	French 50%, Baoulé 14%	68	Gibraltar	English 100%	93	Jersey	English 94.5%, Portuguese 4.6%
19	Benin	French 40%, Fon 39%, Yoruba 12%	44	Croatia	Serbo-Croatian 100%	69	Greece	Greek (Modern) 99%	94	Jordan	Arabic 100%
20	Bermuda	English 100%	45	Cuba	Spanish 100%	70	Greenland	Danish 100%	95	Kazakhstan	Kazakh 63%, Russian 24%
21	Bhutan	Tshangla 28%, Dzongkha 24%, Nepali 22%	46	Cyprus	Greek (Modern) 77%, Turkish 18%	71	Grenada	English 87%, French 2%	96	Kenya	Swahili 80%, English 20%
22	Bolivia	Spanish 60.7%, Quechua 21.2%, Aymara 14.6%	47	Czech Republic	Czech 95.4%, Slovak 1.6%	72	Guadeloupe	French 99%	97	Kiribati	Gilbertese 62.6%
23	Bosnia and Herzegovina	Serbo-Croatian 100%	48	Denmark	Danish 100%	73	Guam	English 38.3%, Chamorro 22.2%, Filipino 22.2%	98	Korea, North	Korean 100%
24	Botswana	Tswana 78.2%, Kalanga 7.9%, English 2.1%	49	Djibouti	Somali 38%, Arabic 20%, French 20%, Afar 13%	74	Guatemala	Spanish 60%	99	Korea, South	Korean 100%
25	Brazil	Portuguese 100%	50	Dominican Republic	Spanish 100%	75	Guinea	French 100%	100	Kosovo	Albanian 100%

Table S6 Language demographics by country. Values for each country add to 100% or less (continued next page)

101	Kuwait	Arabic 100%	126	Morocco	Arabic 90%	151	Russia	Russian 100%	176	Taiwan	Chinese 100%
102	Kyrgyzstan	Kirghiz 64.7%, Uzbek 13.6%, Russian 12.5%	127	Mozambique	Makhuwa 25.3%, Portuguese 10.7%, Tsonga 10.3%, Sena 7.5%, Lomwe 7%, Chuwabu 5.1%	152	Rwanda	Kinyarwanda 98%	177	Tajikistan	Tajik 100%
103	Laos	Lao 100%	128	Namibia	Afrikaans 60%, German 32%, English 7%	153	Saint Kitts and Nevis	English 100%	178	Tanzania	Swahili 100%
104	Latvia	Latvian 58.2%, Russian 37.5%	129	Nauru	Nauru 100%	154	Saint Lucia	English 100%	179	Thailand	Thai 100%
105	Lebanon	Arabic 80%, French 20%	130	Nepal	Nepali 47.8%, Maithili 12.1%, Bhojpuri 7.4%	155	Samoa	Samoan 90%, English 10%	180	Timor-Leste	Tetum 36.6%, English 31.4%, Portuguese 23.5%
106	Lesotho	Southern Sotho 100%	131	Netherlands	Dutch 100%	156	Saudi Arabia	Arabic 100%	181	Togo	French 30%
107	Liberia	English 20%	132	New Caledonia	French 97%	157	Senegal	Wolof 70%, French 10%	182	Tonga	Tonga (Tonga Islands) 70%, English 30%
108	Libya	Arabic 95%	133	New Zealand	English 91.2%, Maori 3.9%, Samoan 2.1%, Chinese 2.1%, French 1.3%, Hindi 1.1%	158	Serbia	Serbo-Croatian 90.1%, Hungarian 3.8%, Romany 1.1%	183	Trinidad and Tobago	English 90%
109	Lithuania	Lithuanian 82%, Russian 8%, Polish 5.6%	134	Nicaragua	Spanish 97.5%	159	Seychelles	Seselwa Creole French 91%, English 4.9%	184	Tunisia	Arabic 100%
110	Luxembourg	Luxembourgish 77%, French 6%, German 4%, English 1%	135	Niger	Hausa 49.6%, Zarma 25.5%, Tamashek 8.4%, Fulah 8.3%, French 5%	160	Sierra Leone	Krio 90%	185	Turkey	Turkish 85.4%, Kurdish 12%, Arabic 1.2%
111	Macedonia	Macedonian 66.5%, Albanian 25.1%, Turkish 3.5%, Romany 1.9%, Serbo-Croatian 1.2%	136	Nigeria	English 30%	161	Singapore	Chinese 58.8%, English 23%, Malay 14.1%, Tamil 3.2%	186	Turkmen-istan	Turkmen 72%, Russian 12%, Uzbek 9%
112	Madagascar	French 70%, Malagasy 30%	137	Norway	Norwegian 100%	162	Slovakia	Slovak 83.9%, Hungarian 10.7%, Romany 1.8%, Ukrainian 1%	187	Uganda	Ganda 14%, English 8%
113	Malawi	Nyanja 70%, Yao 10.1%, Tumbuka 9.5%	138	Oman	Arabic 100%	163	Slovenia	Slovenian 91.1%, Serbo-Croatian 4.5%	188	Ukraine	Ukrainian 67%, Russian 24%
114	Malaysia	Malay 100%	139	Pakistan	Panjabi 48%, Sindhi 12%, Lahnda 10%, Urdu 8%, Pushto 8%	164	Solomon Islands	English 2%	189	United Arab Emirates	Arabic 100%
115	Maldives	Dhivehi 100%	140	Palestinian Authority	Arabic 100%	165	Somalia	Somali 80%, Arabic 20%	190	United Kingdom	English 100%
116	Mali	Bambara 46.3%, French 10%, Fulah 9.4%, Soninke 6.4%	141	Panama	Spanish 100%	166	South Africa	Zulu 23.82%, Xhosa 17.64%, Afrikaans 13.35%, Pedi 9.39%, Tswana 8.2%, English 8.2%, Southern Sotho 7.93%	191	United States	English 82.1%, Spanish 10.7%
117	Malta	Maltese 90.2%, English 6%	142	Papua New Guinea	English 2%, Tok Pisin 1.8%	167	South Sudan	Arabic 50%	192	Uruguay	Spanish 100%
118	Martinique	French 100%	143	Paraguay	Guarani 50%, Spanish 50%	168	Spain	Spanish 74%, Catalan 17%, Galician 7%, Basque 2%	193	Uzbekistan	Uzbek 74.3%, Russian 14.2%, Tajik 4.4%
119	Mauritania	Arabic 100%	144	Peru	Spanish 84.1%, Quechua 13%, Aymara 1.7%	169	Sri Lanka	Sinhala 74%, Tamil 18%	194	Vanuatu	Bislama 23.1%, English 1.9%, French 1.4%
120	Mauritius	Bhojpuri 12.1%, French 3.4%, English 1%	145	Philippines	Filipino 100%	170	Sudan	Arabic 100%	195	Venezuela	Spanish 100%
121	Mexico	Spanish 98.5%	146	Poland	Polish 97.8%	171	Suriname	Dutch 60%	196	Vietnam	Vietnamese 100%
122	Moldova	Romanian 76.5%, Russian 11.2%, Ukrainian 4.4%, Gagauz 4%, Bulgarian 1.6%	147	Portugal	Portuguese 100%	172	Swaziland	Swati 98%	197	Virgin Islands	English 74.7%, Spanish 16.8%, French 6.6%
123	Monaco	French 100%	148	Puerto Rico	Spanish 90%, English 10%	173	Sweden	Swedish 100%	198	Yemen	Arabic 100%
124	Mongolia	Mongolian 90%	149	Qatar	Arabic 100%	174	Switzerland	German 63.7%, French 20.4%, Italian 6.5%, Serbo-Croatian 1.5%, Albanian 1.3%, Portuguese 1.2%, Spanish 1.1%, English 1%	199	Zambia	Bemba 30.1%, English 16%, Nyanja 10.7%, Tonga (Zambia) 10.6%, Lozi 5.7%
125	Montenegro	Serbo-Croatian 91.1%, Albanian 5.3%	150	Romania	Romanian 91%, Hungarian 6.7%, Romany 1.1%	175	Syria	Arabic 100%	200	Zimbabwe	Shona 70%, North Ndebele 20%, English 2.5%

S5.2 Wikipedia

Wikipedia is available in more than 270 language editions. As Wikipedia is collaboratively authored, each edition reflects the knowledge of the language community that contributed to it (22, 23). For example, an article about Plato in the Filipino Wikipedia indicates that Plato is known enough among speakers of Filipino to motivate some of them to write an article about him. While a Wikipedia article in just one language can be the result of short-lived fame within a limited community, a person with articles written about him or her in many languages has likely made a substantial cultural contribution that impacted people from a diverse linguistic and cultural background.

We compiled our *Wikipedia* dataset of famous people as follows. We started by retrieving a table of 2,345,208 people from *Freebase* (www.freebase.com), a collaboratively curated repository of structured data of millions of entities, such places and people. We used a data dump from November 4, 2012; the latest version of the table is available from Freebase (24). For each person, the table contains his or her name, date of birth, place of birth, occupation, and additional information. In addition, for each person with an article in the English Wikipedia, Freebase stores the *Wikipedia unique identifier* (known as *pageid* or *curid*) of the respective article, which we retrieved through the Freebase API (25). The *pageid* and the Wikipedia API (26) were used to find the number of language editions in which a person had an article. Then, the *pageid*, Wikipedia article name, and number of languages of each article were added to the table retrieved from Freebase.

We matched 991,684 people with the English Wikipedia, from which we selected 216,280 people with a defined date of birth, place of birth and gender. We then restricted this list to include only the 11,340 people who had articles in at least 26 Wikipedia language editions and a defined date of birth, place of birth and gender. We then validated the places of birth for all people on the list and converted them to a standardized format (e.g., entries such as “NYC”, “New York” or “New York City” were all converted to “New York, NY, US”). After examining biographical articles in all Wikipedia language editions, we found that there is no biography that appears in at least 26 languages or more that does not have an English version. Thus, by compiling biographies from the English Wikipedia we capture the famous people in any other Wikipedia language. The 26-language threshold generated a group that

is exclusive enough while still containing enough data points, and was within a reasonable size that allowed a comprehensive curation and normalization effort. For comparison, a 20-language threshold would give us 13,334 articles, and a 30-language threshold would give us 6,336 articles.

Next, we converted dates to a standard four-digit year format. While doing so, we fixed all BCE years, which the Freebase dump listed one year off. For example, Jesus's year of birth was listed as 3 BCE instead of 4 BCE. We then used the Google Geocoding API (27) to resolve the listed places of birth to latitude-longitude coordinates, and used the GeoNames database (www.geonames.com) to resolve the coordinates into the present-day name of the country in which each person was born. After dropping records with an ambiguous place of birth we remained with 10,773 people—to which we refer henceforth as the *Wikipedia 26* dataset. Finally, we converted countries to languages as described in Section 4.1 above. To increase the accuracy of the conversion, we selected from the *Wikipedia 26* dataset only the 4,886 people who were born after 1800 and before 1950.

The following tables show the number of famous people in the *Wikipedia 26* dataset for each country (Table S7) and language (Table S8).

	Country	People (all years)	People (1800-1950)		Country	People (all years)	People (1800-1950)		Country	People (all years)	People (1800-1950)
1	Afghanistan	21	10	67	Greece	140	34	133	Nigeria	23	6
2	Albania	15	7	68	Greenland	1		134	Norway	59	33
3	Algeria	17	11	69	Guadeloupe	4	1	135	Oman	2	1
4	Andorra	1		70	Guam	1	1	136	Pakistan	28	13
5	Angola	5	4	71	Guatemala	5	2	137	Palestinian State	14	2
6	Antigua and Barbuda	1	1	72	Guinea	5	3	138	Panama	4	3
7	Argentina	102	33	73	Guinea-Bissau	3	3	139	Paraguay	13	3
8	Armenia	12	4	74	Guyana	1		140	Peru	21	12
9	Aruba	1	1	75	Haiti	7	2	141	Philippines	19	16
10	Australia	95	28	76	Honduras	4	1	142	Poland	167	114
11	Austria	139	91	77	Hong Kong	5		143	Portugal	88	16
12	Azerbaijan	15	6	78	Hungary	81	58	144	Puerto Rico	6	
13	Bahrain	1	1	79	Iceland	15	8	145	Qatar	1	
14	Bangladesh	8	7	80	India	136	69	146	Romania	50	26
15	Barbados	1		81	Indonesia	8	7	147	Russia	369	240
16	Belarus	22	10	82	Iran	61	20	148	Rwanda	1	1
17	Belgium	103	40	83	Iraq	29	8	149	Saint Kitts and Nevis	1	
18	Benin	3	1	84	Ireland	73	29	150	Saint Lucia	2	2
19	Bermuda	1		85	Isle of Man	4	3	151	Samoa	1	1
20	Bhutan	4	1	86	Israel	73	20	152	Sao Tome and Principe	1	1
21	Bolivia	3	1	87	Italy	793	194	153	Saudi Arabia	35	9
22	Bosnia and Herzegovina	26	8	88	Jamaica	10	3	154	Senegal	10	2
23	Botswana	4	3	89	Japan	137	75	155	Serbia	60	12
24	Brazil	137	53	90	Jersey	1		156	Seychelles	1	1
25	Brunei	1	1	91	Jordan	7	3	157	Sierra Leone	1	
26	Bulgaria	29	8	92	Kazakhstan	10	6	158	Singapore	7	4
27	Burkina Faso	2	1	93	Kenya	10	8	159	Slovakia	24	6
28	Burma	7	7	94	Korea, North	6	4	160	Slovenia	15	3
29	Burundi	1		95	Korea, South	37	17	161	Solomon Islands	1	
30	Cambodia	5	2	96	Kosovo	7		162	Somalia	8	3
31	Cameroon	11	2	97	Kuwait	3	2	163	South Africa	43	22
32	Canada	106	46	98	Kyrgyzstan	5	4	164	South Sudan	1	1
33	Cape Verde	4	1	99	Laos	1	1	165	Spain	298	77
34	Central African Republic	1	1	100	Latvia	18	11	166	Sri Lanka	6	5
35	Chad	2		101	Lebanon	13	6	167	Sudan	4	4
36	Chile	27	13	102	Lesotho	1		168	Suriname	5	2
37	China	94	37	103	Liberia	5	2	169	Swaziland	1	
38	Colombia	17	3	104	Libya	11	2	170	Sweden	135	61
39	Congo, Democratic Republic of the	7	3	105	Lithuania	28	19	171	Switzerland	102	56
40	Congo, Republic of	2	1	106	Luxembourg	8	4	172	Syria	19	2
41	Costa Rica	3	1	107	Macedonia	15	3	173	Taiwan	10	4
42	Cote d'Ivoire	15	3	108	Madagascar	2	1	174	Tajikistan	1	
43	Croatia	56	10	109	Malawi	4	4	175	Tanzania	3	3
44	Cuba	13	9	110	Malaysia	6	4	176	Thailand	7	5
45	Cyprus	9	5	111	Maldives	3	1	177	Timor-Leste	3	3
46	Czech Republic	105	53	112	Mali	8	4	178	Togo	5	2
47	Denmark	99	39	113	Malta	3	2	179	Tonga	2	1
48	Djibouti	1		114	Martinique	3	2	180	Trinidad and Tobago	5	2
49	Dominican Republic	2	1	115	Mauritania	1		181	Tunisia	18	7
50	Ecuador	4	1	116	Mauritius	1	1	182	Turkey	184	35
51	Egypt	68	24	117	Mexico	56	23	183	Turkmenistan	3	1
52	El Salvador	3	1	118	Micronesia, Federated States	1	1	184	Uganda	5	3
53	Equatorial Guinea	1	1	119	Moldova	5	2	185	Ukraine	100	58
54	Eritrea	1	1	120	Monaco	4	1	186	United Arab Emirates	5	4
55	Estonia	15	9	121	Mongolia	8	1	187	United Kingdom	1,140	508
56	Ethiopia	10	6	122	Montenegro	10	4	188	United States	2,291	1,221
57	Faroe Islands	1	1	123	Morocco	14	7	189	Uruguay	23	7
58	Finland	63	34	124	Mozambique	6	3	190	Uzbekistan	9	1
59	France	857	397	125	Namibia	2	2	191	Vanuatu	1	1
60	French Guiana	1		126	Nauru	1		192	Venezuela	12	3
61	Gabon	3	3	127	Nepal	4	3	193	Vietnam	10	9
62	Gambia, The	1		128	Netherlands	162	56	194	Virgin Islands	2	1
63	Georgia	21	12	129	New Caledonia	2		195	Yemen	6	2
64	Germany	740	407	130	New Zealand	17	9	196	Zambia	3	3
65	Ghana	17	4	131	Nicaragua	5	5	197	Zimbabwe	7	4
66	Gibraltar	1		132	Niger	1	1			10,773	4,886

Table S7 Number of people with articles in at least 26 Wikipedia language editions, by *country*.

	Language	Code	People (all years)	People (1800-1950)		Language	Code	People (all years)	People (1800-1950)
1	Afrikaans	afr	6.94	4.14	33	Latvian	lav	10.48	6.4
2	Albanian	sqi	26.87	8.34	34	Lithuanian	lit	22.96	15.58
3	Arabic	ara	273.07	94.46	35	Macedonian	mkd	9.97	2
4	Armenian	hye	13.42	4.84	36	Malay	msa	15.99	12.56
5	Azerbaijani	aze	25.79	9.74	37	Malayalam	mal	4.35	2.21
6	Basque	eus	5.96	1.54	38	Maltese	mlt	2.71	1.8
7	Belarusian	bel	5.15	2.34	39	Maori	mri	0.66	0.35
8	Bengali	ben	18.86	12.45	40	Marathi	mar	9.52	4.83
9	Bulgarian	bul	22.35	6.18	41	Modern Greek	ell	147.22	38.08
10	Catalan	cat	51.06	13.09	42	Mongolian	mon	7.2	0.9
11	Chinese	zho	115.6	44.24	43	Norwegian	nor	59	33
12	Czech	ces	100.17	50.56	44	Persian	fas	42.83	15.6
13	Danish	dan	100	39	45	Polish	pol	164.89	112.56
14	Dutch	nld	226.86	81.26	46	Portuguese	por	235.69	74.92
15	English	eng	3300.8	1617.77	47	Romanian	ron	49.33	25.19
16	Estonian	est	10.1	6.06	48	Russian	rus	429.38	272.91
17	Filipino	fil	19.22	16.22	49	Serbo-Croatian	hbs	152.84	36.92
18	Finnish	fin	57.46	31.01	50	Sinhala	sin	4.44	3.7
19	French	fra	997.7	455.51	51	Slovak	slk	21.82	5.88
20	Galician	glg	20.86	5.39	52	Slovenian	slv	13.66	2.73
21	Georgian	kat	14.91	8.52	53	Spanish	spa	774.64	305.48
22	German	deu	929.09	524.1	54	Swahili	swa	12.4	10
23	Haitian	hat	5.25	1.5	55	Swedish	swe	138.47	62.87
24	Hebrew	heb	58.4	16	56	Tajik	tgk	1.4	0.04
25	Hindi	hin	55.95	28.39	57	Tamil	tam	9.33	5.1
26	Hungarian	hun	84.01	57.13	58	Thai	tha	7	5
27	Icelandic	isl	15	8	59	Turkish	tur	164.86	33.64
28	Italian	ita	801.15	198.09	60	Turkmen	tuk	3.21	1.22
29	Japanese	jpn	137	75	61	Ukrainian	ukr	67.46	39.01
30	Kazakh	kaz	6.3	3.78	62	Urdu	urd	9.04	4.49
31	Kirghiz	kir	3.23	2.59	63	Uzbek	uzb	8.9	1.98
32	Korean	kor	43	21	64	Vietnamese	vie	10.95	9.28

Table S8 Number of people with articles in at least 26 Wikipedia language editions, by *language*.

S5.3 Human Accomplishment

The book *Human Accomplishment: The Pursuit of Excellence in the Arts and Sciences, 800 B.C. to 1950* (21) ranks the contribution of 3,869 people to different fields of arts and science. Each listed person is ranked on a scale of 1 to 100 for his or her contribution to one or more of the following fields: art, literature, music, philosophy, astronomy, biology, chemistry, earth sciences, mathematics, medicine, physics and technology. People who contributed to more than one field were ranked separately for each field. For example, Isaac Newton received the highest score of 100 for his contribution in physics, and a score of 88.93 for his contribution in mathematics. For each person, the *Human Accomplishment* tables contain his or her name, ranking in all relevant fields, year of birth, year of death, year flourished, country of birth and country of work.

To find the number of famous people for each language, we converted countries of birth to languages as explained in **Section S5.2**. In most cases, we used the countries of birth as listed on *Human Accomplishment*. However, the dataset occasionally provided a geographical or cultural region, rather than a country, as a place of birth: *Balkans*, *Latin America*, *Sub-Saharan Africa*, *Arab World*, *Ancient Greece* and *Rome*. We replaced the first three with the specific places of birth for the respective people, as listed on *Wikipedia 26*, and converted them to languages based on their present-day countries. We did not resolve *Arab World*, *Ancient Greece* or *Rome* to specific locations, but instead converted them directly to *Arabic*, *Ancient Greek*, or *Latin*, respectively. As with the *Wikipedia 26* dataset, we increased the accuracy of the country-to-language mapping by selecting only the 1,655 people born between 1800 and 1950. Doing so also removed native speakers of Latin and Ancient Greek.

The following tables show the number of famous people in the *Human Accomplishment* dataset for each country (Table S9) and language (Table S10).

Country	People (all years)	People (1800-1950)	Country	People (all years)	People (1800-1950)
1 <i>Ancient Greece</i>	134	N/A	25 Japan	169	57
2 <i>Arab World</i>	86	14	26 Kenya	1	1
3 Argentina	2	2	27 Mexico	5	4
4 Australia	4	4	28 Montenegro	1	1
5 Austria	75	48	29 Netherlands	84	31
6 Belgium	82	27	30 New Zealand	3	3
7 Brazil	3	3	31 Nicaragua	1	1
8 Bulgaria	1	1	32 Norway	23	22
9 Canada	11	11	33 Peru	1	1
10 Chile	3	3	34 Poland	25	21
11 China	237	22	35 Portugal	11	4
12 Croatia	5	3	36 Romania	5	4
13 Cuba	3	3	37 Rome	55	N/A
14 Czech Republic	48	28	38 Russia	134	118
15 Denmark	37	20	39 Serbia	2	2
16 Finland	6	5	40 Slovakia	4	4
17 France	542	236	41 Slovenia	2	2
18 Germany	536	267	42 South Africa	1	1
19 Greece	9	6	43 Spain	76	26
20 Guatemala	1	1	44 Sweden	44	21
21 Hungary	21	18	45 Switzerland	64	32
22 Iceland	2	1	46 United Kingdom	531	230
23 India	93	16	47 United States	297	272
24 Italy	389	58	Total	3869	1655

Table S9 Number of people listed on human accomplishment, by *country*.

Language	Code	People (all years)	People (1800-1950)	Language	Code	People (all years)	People (1800-1950)
1 Afrikaans	afr	0.13	0.13	23 Japanese	jpn	169	57
2 Albanian	sqi	0.88	0.47	24 Latin	lat	55	
3 Arabic	ara	86.05	14.05	25 Malayalam	mal	2.98	0.51
4 Basque	eus	1.52	0.52	26 Maori	mri	0.12	0.12
5 Bengali	ben	7.53	1.3	27 Marathi	mar	6.51	1.12
6 Bulgarian	bul	0.77	0.77	28 Norwegian	nor	23	22
7 Catalan	cat	12.92	4.42	29 Polish	pol	24.45	20.54
8 Chinese	zho	237.16	22.16	30 Portuguese	por	14.77	7.38
9 Czech	ces	45.79	26.71	31 Romanian	ron	4.55	3.64
10 Danish	dan	37	20	32 Russian	rus	134	118
11 Dutch	nld	133.2	47.2	33 Serbo-Croatian	hbs	11.61	8.11
12 English	eng	788.1	466.26	34 Slovak	slk	4.12	3.8
13 Finnish	fin	5.47	4.56	35 Slovenian	slv	1.82	1.82
14 French	fra	590.27	255.74	36 Spanish	spa	104.02	63.01
15 Galician	glg	5.32	1.82	37 Swahili	swa	0.8	0.8
16 German	deu	643.22	329.91	38 Swedish	swe	44.33	21.27
17 Greek (Ancient)	grc	134		39 Tamil	tam	5.49	0.94
18 Greek (Modern)	ell	8.96	5.99	40 Turkish	tur	1.81	1.19
19 Hindi	hin	38.16	6.59	41 Ukrainian	ukr	0.04	0.04
20 Hungarian	hun	20.5	17.62	42 Urdu	urd	4.65	0.8
21 Icelandic	isl	2	1	43 Vietnamese	vie	0.04	0.04
22 Italian	ita	393.22	60.14				

Table S10 Number of people listed on human accomplishment, by *language*.

S5.4 Comparison of the famous people datasets

The two datasets we use—*Wikipedia 26* and *Human Accomplishment*—were compiled in different ways. Wikipedia is written by a large number of volunteers with different backgrounds from all over the world, while Human Accomplishment is the work of a single author, the American political scientist Charles Murray. Naturally, both sources exhibit certain biases despite the efforts taken by their authors.

To understand these biases, we compared the cultural significance attributed by each dataset to the listed individuals. We define the cultural significance of a person as the number of languages in which his/her Wikipedia biography is available (for entries on Wikipedia 26), or the score that Murray gave this individual (*Human Accomplishment* entries are given a score from 1 to 100 based on their contribution in their respective field or fields). Figure S4 shows the correlation between these two measurements. One notable observation is that the cultural contribution the Charles Murray attributes to people born in Asia (measured by their score on his list) is higher than their cultural contribution according to *Wikipedia 26* (measured by the number of languages in which a Wikipedia biography is available). Murray is also less likely than Wikipedia to acknowledge the contribution of left-wing liberals.

The moderate correlation ($R^2=0.25$) shows that using these two lists of famous individuals provides a more balanced perspective than the exclusive use of Wikipedia. While the two datasets are substantially different, there is a consistent correlation between the number of famous people in a language according to either dataset and the centrality of that language, attesting to the robustness of our method.

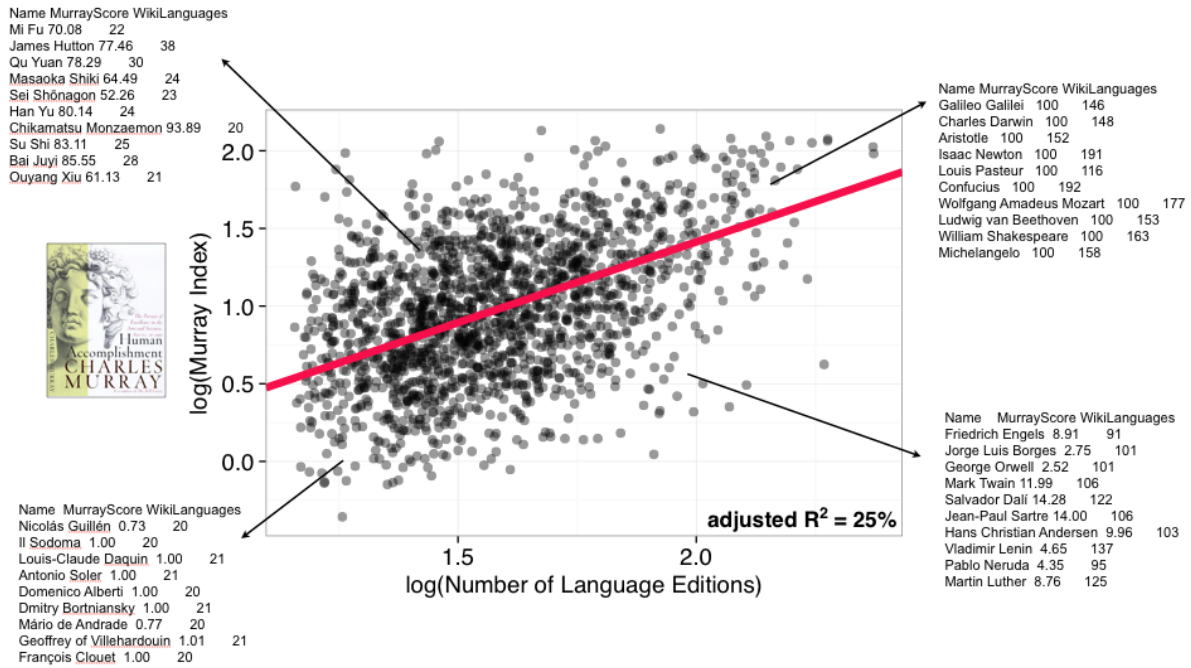


Figure S4 Correlation of the *Wikipedia 26* and *Human Accomplishment* datasets

References for the SI Appendix

1. Rodriguez S (2012) Another Milestone for Twitter: 200 Million Monthly Active Users. *Los Angel Times*. Available at: <http://www.latimes.com/business/technology/la-fi-tn-twitter-200-million-monthly-active-users-20121219,0,3316419.story> [Accessed March 12, 2013].
2. Boyd D, Crawford K (2011) Six Provocations for Big Data. *SSRN ELibrary*. Available at: http://papers.ssrn.com/sol3/papers.cfm?abstract_id=1926431 [Accessed November 1, 2012].
3. Pew Internet & American Life Project (2013) *Twitter Reaction to Events Often at Odds with Overall Public Opinion* (Pew Internet & American Life Project) Available at: <http://www.pewresearch.org/2013/03/04/twitter-reaction-to-events-often-at-odds-with-overall-public-opinion/> [Accessed March 6, 2013].
4. McCandless M (2011) *Chromium Compact Language Detector* Available at: <http://code.google.com/p/chromium-compact-language-detector/>.
5. Mocanu D, et al. (2013) The Twitter of Babel: Mapping World Languages through Microblogging Platforms. *PLoS ONE* 8(4):e61981.
6. Graham M, Hale SA, Gaffney D (2013) Where in the world are you? Geolocation and language identification in Twitter. *Prof Geogr*.
7. Herring SC, et al. (2007) in *System Sciences, 2007. HICSS 2007. 40th Annual Hawaii International Conference on* Available at: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4076532 [Accessed December 13, 2012].
8. International Information Centre for Terminology (2002) ISO 639-1 Registration Authority. Available at: http://www.infoterm.info/standardization/iso_639_1_2002.php.
9. SIL International (2007) ISO 639-3 Registration Authority. Available at: <http://www.sil.org/iso639-3> [Accessed June 14, 2012].
10. Erard M (2012) *Babel No More: The Search for the World's Most Extraordinary Language Learners* (Free Press, New York).
11. Meta-Wiki List of Wikipedias. Available at: http://meta.wikimedia.org/wiki/List_of_Wikipedias [Accessed March 10, 2013].
12. UNESCO Index Translationum: World Bibliography of Translation. Available at: <http://www.unesco.org/xtrans/bsform.aspx> [Accessed July 22, 2012].
13. UNESCO Contributions from Countries. *Index Transl*. Available at: <http://www.unesco.org/xtrans/bscontrib.aspx> [Accessed September 1, 2012].

14. Ruhlen M (1991) *A Guide to the World's Languages: Classification* (Stanford University Press).
15. Lewis MP (2009) *Ethnologue: Languages of the World* (SIL international, Dallas, TX). 16th Ed. Available at: <http://www.ethnologue.com/16> [Accessed November 13, 2012].
16. Library of Congress (2008) ISO 639-5 Registration Authority. Available at: <http://www.loc.gov/standards/iso639-5>.
17. Zachte E (2012) Wikipedia Statistics. *Wikimedia Stat*. Available at: <http://stats.wikimedia.org/EN/Sitemap.htm> [Accessed June 14, 2012].
18. International Monetary Fund (2012) *World Economic Outlook Database, April 2012* Available at: <http://www.imf.org/external/pubs/ft/weo/2012/01/weodata/index.aspx> [Accessed July 17, 2012].
19. Central Intelligence Agency (2011) *The World Factbook* (Central Intelligence Agency, Washington, DC).
20. Bonacich P (1987) Power and Centrality: A Family of Measures. *Am J Sociol* 92(5):1170–1182.
21. Murray CA (2003) *Human Accomplishment: The Pursuit of Excellence in the Arts and Sciences, 800 B.C. to 1950* (HarperCollins, New York).
22. Graham M (2011) in *Critical Point of View: A Wikipedia Reader*, eds Lovink GW, Tkacz N, pp 269–282.
23. Hecht B, Gergle D (2009) in *Proceedings of the fourth international conference on Communities and technologies, C&T '09*. (ACM, New York, NY, USA), pp 11–20. Available at: <http://doi.acm.org/10.1145/1556460.1556463> [Accessed November 5, 2012].
24. Freebase (2012) person.tsv. Available at: <http://download.freebase.com/datadumps/latest/browse/people/person.tsv> [Accessed November 9, 2012].
25. Freebase Wiki Freebase API. Available at: http://wiki.freebase.com/wiki/Freebase_API [Accessed March 10, 2013].
26. Wikimedia MediaWiki API. Available at: <https://www.mediawiki.org/wiki/API>.
27. Google The Google Geocoding API v3. Available at: <https://developers.google.com/maps/documentation/geocoding/>.