

# Summary of findings

## Main analyses of Wikipedia data.

See the main text for a summary of the main analyses using the Wikipedia data.

## Analyses of the estimates from the Common Crawl data

Mixed effects model: The correlation between semantic alignment and cultural similarity was significant ( $\beta=0.182$ ,  $\chi^2(1)=4.94$ ,  $p=0.026$ ). See the figure 1 below:

MRM results:

	Estimate	p-value
Intercept	0.011	0.84
Cultural distance	0.086	0.172
Language family	0.001	0.987
Geographic distance	0.043	0.318
Comparison count	0.871	0.0001

Table 1: MRM analysis predicting semantic alignment (Common Crawl), with family control.  $R^2=0.761$

	Estimate	p-value
Intercept	0.034	0.0311
Cultural distance	0.072	0.315
ASJP	-0.168	0.0009
Geographic distance	0.075	0.108
Comparison count	0.841	0.0001

Table 2: MRM analysis predicting semantic alignment (Common Crawl), with ASJP control.  $R^2=0.744$

Mantel tests:

	Var1	Var2	r	llim	ulim	p
2	Cultural	Linguistic	-0.084	-0.201	0.000885	0.277
3	Cultural	Historical	-0.315	-0.447	-0.204	0.0211
4	Cultural	Geographic	-0.461	-0.586	-0.312	0.00295
5	Linguistic	Historical	-0.25	-0.357	-0.114	0.015
6	Linguistic	Geographic	0.136	-0.0134	0.235	0.121
7	Historical	Geographic	0.405	0.31	0.519	0.00101
71	Linguistic	Cultural **	-0.0831	-0.152	-0.00152	0.27

Table 3: Mantel tests (Common Crawl). \*\* = partial Mantel test, controlling for historical and geographical distance.

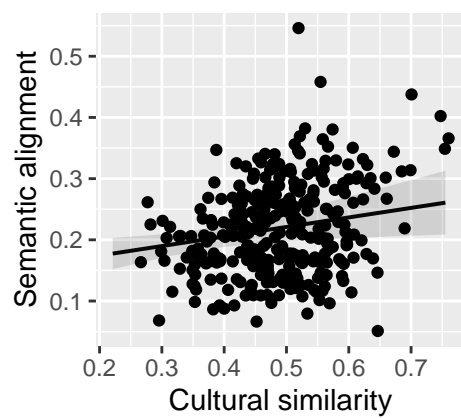


Figure 1: Semantic alignment and cultural similarity for data using the Common Crawl alignments

## Analyses of the estimates from the Subtitles data

Mixed effects model: The correlation between semantic alignment and cultural similarity was not significant ( $\beta = 0.0606$ ,  $\chi^2(1) = 0.74$ ,  $p = 0.39$ ). See the figure 2 below:

MRM results:

	Estimate	p-value
Intercept	0.073	0.517
Cultural distance	0.125	0.413
Language family	-0.041	0.838
Geographic distance	-0.013	0.887
Comparison count	0.806	0.0002

Table 4: MRM analysis predicting semantic alignment (Subtitles), with family control.  $R^2 = 0.744$

	Estimate	p-value
Intercept	0.083	0.297
Cultural distance	-0.023	0.884
ASJP	-0.282	0.0083
Geographic distance	0.019	0.854
Comparison count	0.831	0.0003

Table 5: MRM analysis predicting semantic alignment (Subtitles), with ASJP control.  $R^2 = 0.803$

Mantel tests:

	Var1	Var2	r	llim	ulim	p
2	Cultural	Linguistic	0.351	0.241	0.538	0.0884
3	Cultural	Historical	-0.17	-0.286	-0.0153	0.155
4	Cultural	Geographic	-0.34	-0.57	-0.172	0.0276
5	Linguistic	Historical	-0.352	-0.503	-0.0822	0.0301
6	Linguistic	Geographic	-0.273	-0.462	-0.0818	0.0806
7	Historical	Geographic	0.346	0.181	0.522	0.01
71	Linguistic	Cultural **	0.281	0.129	0.466	0.135

Table 6: Mantel tests (Subtitles). \*\* = partial Mantel test, controlling for historical and geographical distance.

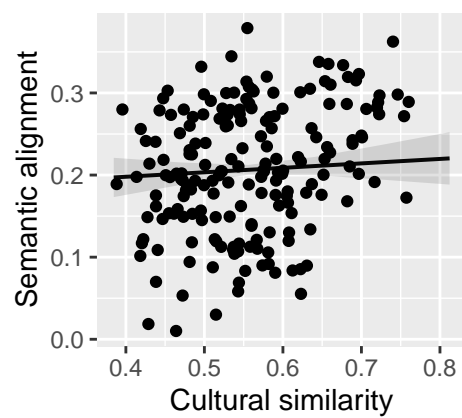


Figure 2: Semantic alignment and cultural similarity for data using the Subtitles alignments

## Analysis of numerals

The analyses of numerals found:

- 1 and 2 have lower alignment due to often being grammaticalised as indefinite or dual marker (Givon, 1981).
- Numbers 3-12 generally have high alignment (mean local alignment = 0.87), and higher numbers decline in alignment up to 1000.
- There are also language-specific differences due to how numerals are constructed (e.g. base, combination rules, see Calude & Verkerk, 2016), or for irregular forms (e.g. 50, 60, 70, 80 and 90 in Danish).
- Some number words have alternative associations due to homophones (e.g. the Hungarian 7 is used directly to mean ‘week’, and ‘neuf’ in French means ‘9’ or ‘new’).
- The historical distance between languages did not predict much of the variation.

See the main text for a discussion.