

1 Speakers of diverse languages structure their utterances for efficient communication

2 Josef Klafka<sup>1</sup> & Daniel Yurovsky<sup>1,2</sup>

3 <sup>1</sup> Carnegie Mellon University

4 <sup>2</sup> University of Chicago

5 Author Note

6 Correspondence concerning this article should be addressed to Josef Klafka, 5000  
7 Forbes Ave. E-mail: jklafka@andrew.cmu.edu

## Abstract

Optimal coding theories of language predict that speakers should keep the amount of information in their utterances relatively uniform under the constraints imposed by their language. But how much of a role do these constraints provide, and does it vary across languages? We find a consistent non-uniform shape which characterizes both spoken and written sentences of English but is tempered by predictive context. We then show that other languages are also characterized by consistent but non-English shaped curves related to their typological features, but that sufficient context produces more uniform shapes across languages. Thus, producers of language appear to structure their utterances in similar near-uniform ways despite varying linguistic constraints.

*Keywords:* information theory; communication; efficiency; syntax; typology; language development; computational modeling

Speakers of diverse languages structure their utterances for efficient communication

## Introduction

One of the defining features of human language is its power to transmit information. We use language for a variety of purposes like greeting friends, making records, and signaling group identity. These purposes all share a common goal: Transmitting information that changes the mental state of our listener (Austin, 1975). For this reason, we can describe language as a cryptographic code, one that allows speakers to turn their intended meaning into a message that can be transmitted to a listener, and subsequently converted by the listener back into an approximation of the intended meaning (Shannon, 1948).

How should we expect this code to be structured? If language has evolved as a code for information transmission, its structure should reflect this process of optimization (Anderson & Milson, 1989). The optimal code would have to work with two competing pressures: (1) For listeners to easily and successfully decode messages sent by the speaker, and (2) For speakers to easily code their messages and transmit them to a listener with minimal effort and error. A fundamental constraint on both of these processes is the linear order of spoken language—sounds are produced one at a time and each is unavailable perceptually once it is no longer being produced.

Humans accommodate this linear order constraint through incremental processing: People process speech continuously as it arrives, predicting upcoming words and building expectations about the meaning of an utterance in real time rather than at its conclusion (Kutas & Federmeier, 2011; Pickering & Garrod, 2013; Tanenhaus, Spivey-Knowlton, Eberhard, & Sedivy, 1995). This solution creates new guidance for speakers. Since prediction errors can lead to severe processing costs and difficulty integrating new information on the part of listeners, speakers should seek to minimize prediction errors. However, the cost of producing more predictable utterances is using more words. Thus, the most efficient strategy is for speakers seeking to minimize their production costs is to produce utterances that are

just at the prediction capacity of listeners without exceeding this capacity (Aylett & Turk, 2004; Genzel & Charniak, 2002). In other words, speakers should maintain a constant transmission of information, with the optimal rate of information transfer as close to the listener's fastest decoding rate as possible. The hypothesis that speakers follow this optimal strategy is known as the *Uniform Information Density* hypothesis.

Using information theory, a mathematical framework for formalizing predictability, researchers have tested and confirmed this optimal coding prediction across several levels and contexts in language production. For example, Genzel and Charniak (2002) provided a clever indirect test of Uniform Information Density across sentences in a paragraph. They showed that the predictability of successive sentences, when analyzed in isolation, decreases, as would be expected if readers use prior sentences to predict the content of future sentences. Thus, based on the increasing amount of context, they found that total predictability remains constant. At the level of individual words, Mahowald, Fedorenko, Piantadosi, and Gibson (2013) showed that speakers use shorter alternatives of more predictable words, maximizing the amount of information in each word while minimizing the time spent on those words.

Other research has suggested that efficient encoding impacts how speakers structure units between words and sentences. The inclusion of complementizers in relative clauses (Jaeger & Levy, 2007) and the use of contractions (Frank & Jaeger, 2008) are two situations in sentence formation in which speakers can omit or reduce words to communicate more efficiently and maximize use of the communication channel without exceeding the listener's capacity.

How languages evolve is shaped by efficient communication as well. Piantadosi, Tily, and Gibson (2011) showed that more easily predictable words in a language may tend to become shorter over time, maximizing the amount of information transmitted over the communication channel at every second by speakers in each language. Semantic categories of words across languages can also evolve to be structured efficiently. Categories such as kinship

terms (Kemp & Regier, 2012) maintain a trade-off between informativeness and complexity. Structure in language evolves from a trade-off between efficient and learnable encoding on the one hand and an expressive and descriptive lexicon on the other (Kirby, Tamariz, Cornish, & Smith, 2015). Languages may come to efficiently describe the particular environment in which they are spoken over the course of evolution: features of the world that are relevant to speakers become part of a language, while irrelevant features are disregarded (Perfors & Navarro, 2014).

However, despite this literature using the predictive coding model of language, one level has not yet been studied in depth: how speakers structure each individual utterances. This level may show the strongest effects of variation between languages. While speakers can make bottom-up choices such as controlling which of several near-synonyms they produce, they cannot control the grammatical properties of their language. Properties of a language, like canonical word order, impose top-down constraints on how speakers can structure what they say. While speakers may produce utterances as uniform in information density as their languages will allow, these top-down constraints may create significant and unique variation across languages.

How significant are a language's top-down constraints on determining how its speakers structure their speech? Yu, Cong, Liang, and Liu (2016) analyzed how the information in words of English sentences of a fixed length varies with their order in the sentence (e.g. first word, second word, etc). They found a surprising non-linear shape, and argued that this shape may arise from top-down grammatical constraints in the English language. We build on these ideas, asking (1) Whether this shape depends on listener's predictive models, (2) Whether this shape varies across linguistic contexts, and (3) Whether this shape is broadly characteristic of a diverse set of languages or varies predictably from language to language. We find that languages are characterized by highly-reliable but cross-linguistically variable information structures that co-vary with top-down linguistic features. Listeners' predictive

coding flattens these shapes across languages, in accord with predictions of the Uniform Information Density hypothesis.

## Methods

We measure information structure within languages, using a universal information metric proposed for the study of information transmission more generally by Shannon (1948) and applied to words specifically by Levy (2008): lexical surprisal. We can compute surprisal with the predictability of the word based on previously heard or seen words in its context, as in the formula below. The surprisal of a word is inversely proportional to the predictability of a word, such that less common and less predictable words carry more information. For example, “flower” has less information than “azalea” because “flower” is much more common than “azalea”. Though the two words have the same length in number of letters, it is more difficult to process “azalea” when reading it here than when reading “flower”. Frequency is intimately tied information content in words, with much of the differences between words frequencies being explained by information content cross-linguistically (Piantadosi et al., 2011). The surprisal of a word is also correlated with the processing cost of a word, shown by evidence from e.g. eye-tracking (Smith & Levy, 2013) and ERP (Frank, Otten, Galli, & Vigliocco, 2015) studies.

However, when reading or listening, people don’t just consider each word as an isolated linguistic signal. Listeners use the words they have already heard to predict and decode the word they are currently hearing. Following this incremental processing paradigm, we can also condition the surprisal of a word in its context. Ideally, we would like to measure the predictability of each word in an utterance using all of the information available to that word. For example, in an utterance of twenty words, we would like to use the previous 19 words of context to predict the 20th word. However, we would need to train on a corpus of many trillion word tokens to predict with this amount of context. Regardless of computational constraints, we want to directly compare how predictable each word is

124 regardless of its position in an utterance.

125 We therefore use a simplifying *Markov assumption*: we condition our next predictions  
126 on a fixed-size context window instead of all preceding words. Although these models may  
127 seem to use an inconsequential amount of context when predicting the next word, bigram  
128 and trigram models introduce a great deal of improvement over unigram models across tasks  
129 (Chen & Goodman, 1999). Models which incorporate more than two words of context have  
130 issues with overfitting to the corpus and only predicting observed sequences, often  
131 generalizing poorly.

132 When we use a word or two of context in our surprisal calculations, then the set of  
133 reasonable final items in our ngrams is greatly restricted. “Flower” may contain less  
134 information than “azalea” when we consider the words independently of their context, but  
135 with context this can be reversed. Flower appears in a variety of contexts, and so the  
136 information content of a word like “flower” in a particular context may be higher than  
137 “azalea”. If you only have azaleas in your garden, then hearing someone say “in that garden,  
138 look at the flowers” may be higher surprisal for you: you expect them to say “azalea”. This  
139 prediction does not require many words for context. For example, in the sentence “I take my  
140 coffee with cream and sugar”, when hearing “cream and”, a listener might automatically  
141 predict “sugar”, but there are few possible continuations with even the two words “cream  
142 and”. Hearing “I” restricts the next word to a verb, or possibly an adverb, and since the  
143 listener has heard the speaker refer to themselves in the first person singular, their set of  
144 possible completions is significantly restricted.

145 We train two types of ngram language models independently on a corpus. One of our  
146 models is frequency-based: we do not incorporate context into our surprisal calculations. To  
147 incorporate context into our models, we train bigram and trigram language models, which  
148 incorporate one and two words of context for each processed word, respectively. The  
149 frequency-based surprisal metric gives us an idea of when in their utterances speakers say

frequent i.e. independently information-rich words. The context-based surprisal metric shows us how speakers tend to distribute the information in utterances relative to real-time processing in communication. We expect a priori that our frequency-based surprisal curve will be flat. No one part of the sentence will on average have words that are more frequent than another across utterance lengths. Similarly, we expect that there will be a small smoothing effect for our contextual surprisal metric such that the word in each position of an utterance is more predictable than its frequency-based counterpart.

### Estimating information

To estimate how information is distributed across utterances, we computed the lexical surprisal of each word under two different models. First, following Yu et al. (2016), we estimated a unigram model which considers each word independently:

$$\text{surprisal}(\text{word}) = -\log P(\text{word})$$

This unigram surprisal measure is a direct transformation of the word’s frequency and thus less frequent words are more surprising. Simply the less often a person has seen a word, the more information that word holds.

Second, we estimated a trigram model in which the surprisal of a given word ( $w_i$ ) encodes how unexpected it is to read it after reading the prior two words ( $w_{i-1}$  and  $w_{i-2}$ ):

$$\text{surprisal}(w_i) = -\log P(w_i | w_{i-1}, w_{i-2})$$

This metric encodes the idea that words that are low frequency in isolation (e.g. “meatballs”) may become much less surprising in certain contexts (e.g. “spaghetti and meatballs”) but more surprising in others (e.g. “coffee with meatballs”). The difficulty of



correctly estimating these probabilities from a corpus grows combinatorically with the number of prior words, and in practice trigram models perform well as an approximation (see e.g. Chen & Goodman, 1999; Smith & Levy, 2013).

**Model details.** We estimated the surprisal for each word type in a corpus using the KenLM toolkit (Heafield, Pouzyrevsky, Clark, & Koehn, 2013). Each utterance was padded with a special start-of-sentence token “ $\langle s \rangle$ ” and end of sentence token “ $\langle /s \rangle$ ”. Trigram estimates did not cross sentence boundaries, so for example the surprisal of the second word in an utterance was estimated as  $(\text{surprisal}(w_2) = -P(w_2|w_i, \langle s \rangle))$ . Naïve trigram models will underestimate the surprisal of words in low-frequency trigrams (e.g. if the word “meatballs” appears only once in the corpus following exactly the words “spaghetti and”, it is perfectly predictable from its prior two words).

To avoid this underestimation, we used modified Kneser-Ney smoothing as implemented in the KenLM toolkit (Heafield et al., 2013). Briefly, this smoothing technique discounts all ngram frequency counts, which reduces the impact of rare ngrams on probability calculations, and interpolates lower-order ngrams into the calculations. These lower-order ngrams are weighted according to the number of distinct contexts they occur as a continuation (e.g. “Francisco” may be a common word in a corpus, but likely only occurs after “San” as in “San Francisco”, so it receives a lower weighting). For a thorough explanation of modified Kneser-Ney smoothing, see Chen and Goodman (1999).

**Aggregating curves.** To develop a characteristic information curve for sentences in the corpus, we needed to aggregate sentences that varied dramatically in length (Fig ??A). We used Dynamic Time Warping Barycenter Averaging (DBA), an algorithm for finding the average of sequences that share an underlying pattern but vary in length (Petitjean, Ketterlin, & Gançarski, 2011). DBA inverts standard dynamic time warping, discovering a latent invariant template from a set of sequences.

We used DBA to discover the short sequence of surprisal values that characterized the

surprisal curves common to sentences of varying sentence lengths. We first averaged individual sentences of the same length together and then applied the DBA algorithm to this set of average sequences. DBA requires a parameter specifying the length of the template sequence.

**Optimizing the size hyperparameter for barycenter averaging.** Discussion of optimizing the expectation-maximization approach the barycenters use here.

**Application to the corpus.** Once we have fitted our language model, we can compute the surprisal of a continuation by simply taking the negative log-probability of that word’s ngram probability. To find the average information for a given position in a corpus, we take all utterances of a given length, and for each word position in utterances of that length, we compute the average of the surprisals for all of the non-unique words that occur in that position, conditioned or not conditioned on context. By computing these averages for each word position in an utterance, we compute a low-dimensional approximation to the average distribution of information in the corpus. With the surprisal metric, we base the information contained in each word on how often the word is encountered in its context in the corpus. As long as the corpus is representative of the language or population we study, then the distribution of information is approximated for that language or population as a whole.

## **Frequency-based and contextual information curves in written English: the British National Corpus**

We first turn to working with written English in the British National Corpus (BNC; Leech, 1992). The BNC is a collection of spoken and written records (90% written) from the turn of the century, intended to be a representative sample of British English. Using their word entropy metric without context, Yu et al. (2016) found a distinctive three-step distribution for information in written English sentences in the corpus. The first word tended to contain little information. While the middle words of sentences each had more information than the first word, they found a flat and non-increasing rate of information

transmission across the middle of sentences. The final word contained the most, though not most, of the information out of any in the sentence, with a noticeable spike in information. They found the same distribution across sentence lengths, from sentences with 15 words to sentences with 45 words.

We replicate the Yu et al. (2016) result using the surprisal metric in place of the entropy metric. We use the frequency-based or “contextless” surprisal metric, which determines the average distribution of information based on word frequencies in a corpus. A priori we expect that the frequency-based metric will produce a flat distribution of information across word positions in the BNC. We find the same frequency-based information trajectory as Yu et al. with little information in the first words of utterances and the most information in the final word, see Figure 1.

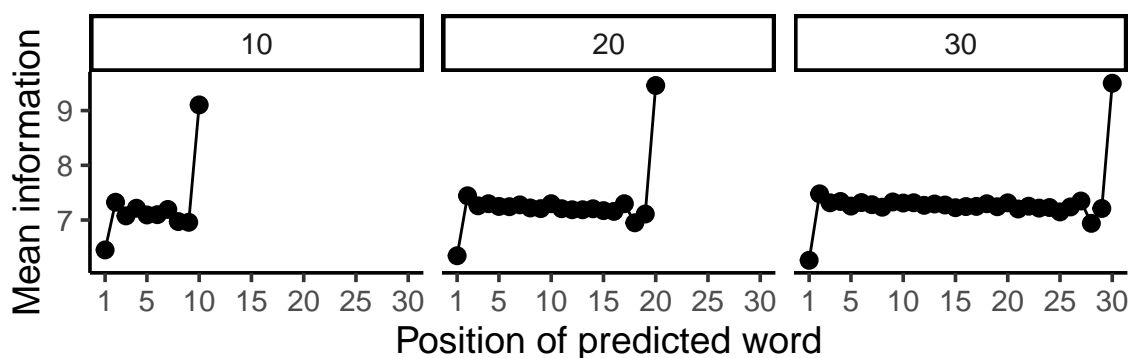


Figure 1. BNC frequency-based information curves

We have found a unique average distribution of information that appears to characterize the English language regardless of utterance length. This distribution indicates that in English, the words we speak or write at the beginnings of utterances have little information, while the words we speak or write at the ends of utterances have a lot of information. The words in the middle of utterances have a medial amount of information, without the increasing trend in information from word to word that we might expect from (Genzel & Charniak, 2002).

What about context? So far we've only discussed the frequency-based metric, considering words on their own without any explicit incorporation of prior context. As previously discussed, listeners decode information and process what they hear incrementally, using prior heard words to ease the comprehension process. We now include two words of context (trigrams) for each word in our measurements. We observe a flattening effect of context across both modalities and all speaker populations. After the first word or two, where the listener does not have access to prior context, then they decode information at a flat and more or less uniform rate. The contextual information curves for the BNC are in Figure 2. We also computed bigram curves with one word of context for each prediction: these bigram curves resemble the trigram curves.

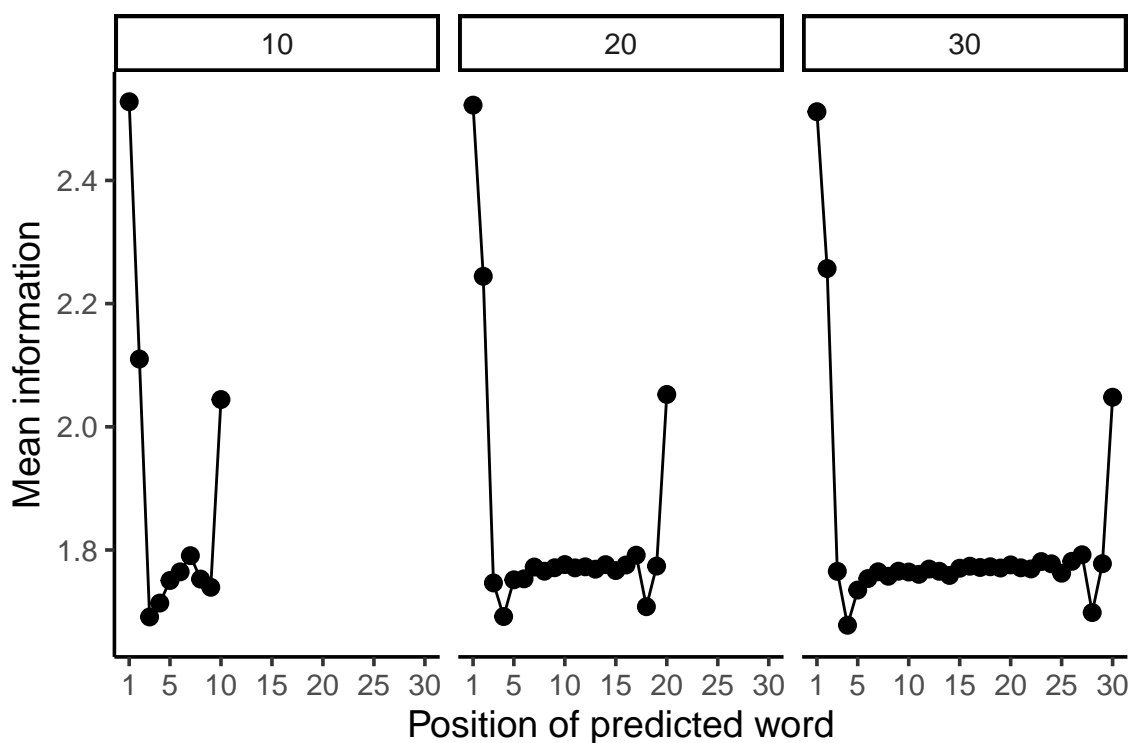


Figure 2. British National Corpus trigram context-based information curves

Speakers produce information at a more or less constant rate, avoiding peaks or troughs in their information distribution, except at the beginnings of utterances, where listeners may not have any context to predict what the speaker is going to say. Speakers of

English tend towards a characteristic and uneven distribution of word information based on frequencies within their utterances. Their interlocutors, however, once they have a word or two of context, decode information at a more or less constant and optimal rate.

### **Frequency-based and contextual information curves in spoken English: the British National Corpus**

We found that English speakers and writers used the same robust and distinctive distribution of information within each utterance, regardless of the number of words in their utterances. To determine if this distribution truly characterizes all speakers of the English language as a whole, we wanted to examine speech from English-speaking children who are producing their very first multi-word utterances. We hypothesize the three-step distribution of information we found for English will characterize child speech and child-directed speech. This approach also allows us to analyze parent speech to children, to examine speech more generally and understand if English speech gives rise to the same information distribution as English writing.

We use the North American English collection from CHILDES (MacWhinney, 2000), which consists of about 2.6 million utterances from 522 children and their parents across 49 corpora. We obtained this collection using the childesr frontend to the childes-db database (Sanchez et al., 2019). The utterances in the Providence corpus are on average significantly shorter than those in the BNC; over 95% of the utterances in the North American English collection are 10 words or fewer. Unlike the written BNC, which we split by sentence, we split our CHILDES corpus by conversational turns and pauses using the built-in utterance breaks for each corpus.

We observe the same distinctive distribution of information for parents and children in the child speech corpus as we did for adults in the BNC. The distribution of information we find at the level of individual words in English, therefore, characterizes the English language

as a whole and not only adult utterances, not only written utterances. See 3.

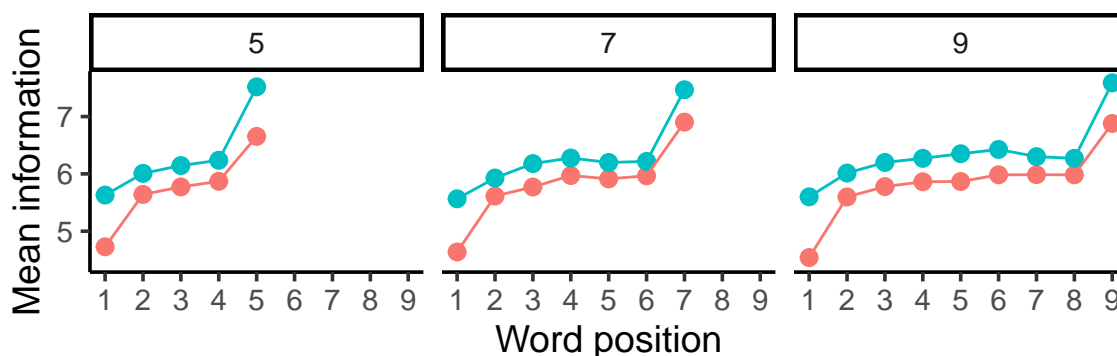


Figure 3. North American English frequency-based information curves. Lines around each point indicate 95% confidence intervals computed with non-parametric bootstrap

What about predictive processing in child-directed speech and child speech? When incorporating one or two words of predictive context, we observe the same trend as in the BNC. Beyond the first couple of words, once your interlocutor has enough context to predict with some accuracy what you will say next, then you decode information from their speech stream at a constant and optimal rate. This applies to parents and children speaking to one another, as well as adults speaking and writing to one another. See 4.

### Frequency-based and contextual information curves across languages: a qualitative analysis

So far, we have only looked at the distribution of information in words in English, both with and without context. We have examined child speech and child-directed speech at a variety of ages, as well as writing samples selected to be representative of British English as a whole. But this only captures the picture for English.

We now turn to a small number of typologically diverse languages, and conduct the same analysis, using monolingual adult-child speech corpora from CHILDES (MacWhinney, 2000) to compare the results from these languages directly to our results from English. We use corpora for Spanish, German, French, Mandarin Chinese and Japanese. Similar to our

English child speech collection, all of the language collections consist mainly of shorter utterances: most utterances in the corpora are under 10 words long. Mandarin and Japanese are not natively written using the Latin alphabet, and moreover words are not segmented in their native scripts. Instead of the native scripts, we use transliterations from the corpus for each of the Mandarin and Japanese utterances into pinyin for Mandarin and romanji for Japanese. In these transliterations, words are previously segmented.

We observe a distinct and characteristic frequency-based information trajectory for each language, robust across each utterance length. We see the same distribution of information for both parents and children. The parent often has more information on average at each word position in their utterances. This is an effect of the surprisal metric: parents speak more utterances than their children in most of the corpora, which inflates the number of tokens they use and increases the surprisal of hearing a rare word. We include the frequency-based information curve from the North American English CHILDES collection for comparison. See Figure 5

English, Spanish, French and German feature similar information curve shapes, with slight variations. The German information curve features lower information for longer towards the beginnings of utterances, possibly due to the grammatical restriction that the second word in German utterances must be a verb (V2). Spanish features a larger spike in the amount of information in the final word of utterances. For Japanese and Mandarin, we observe completely different frequency-based information curve trajectories. The Japanese frequency-based information curve trajectory begins high and finishes low, the mirror image of the German and Romance language information curves. The Mandarin curve begins low and finishes low, but features high information in the middle of utterances. We hypothesize this may be due to Japanese and Mandarin speakers typically ending their utterances with particles, which are common and thus contain little information on their own.

For the rigram information curves, we see the same contextual smoothing effect as in

English. While the frequency-based information curves may depend based on the language, the contextual information curves show the same trajectory cross-linguistically. Using more than two words of context is difficult for parent-child speech corpora because the utterances are so short on average (less than 10 words). Based on our results from the CHILDES collections, we hypothesize that the frequency-based information curves may vary based on the genealogy and typology of the languages in question. However, this does not extend to the information curves with two words of context in particular, where all languages we have seen so far are characterized by the same information distribution. See Figure 6.

### **Language structures and large-scale data analysis: methods**

To make a claim about how languages on a larger scale, we need to use larger corpora and a much larger number of languages. We pulled corpora for 159 diverse languages from Wikipedia, each of which had at least 10,000 articles on the knowledge base. We split each article into sentences; the variance in sentence lengths for Wikipedia was significantly larger than for the CHILDES corpora we used in the previous section. Most sentences in Wikipedia contained between 10 and 30 words, unlike the CHILDES corpora which mainly contained utterances with under 10 words. We excluded the small fraction of utterances with more than 50 words since they were small in number and, from manual inspection, uncharacteristic of typical written sentences.

How do we quantitatively analyze information curves for more than 40 difference sentence lengths for each language, adding up to several thousand information curves total? We used two different strategies, which yielded identical results upon analysis. Each strategy gave us a five-dimensional vector for each language in a Wikipedia “slope space”. For the first strategy, we split each sentence length by number of words into fifths, and computed surprisal values for the closest word position to each quintile. We then computed the slopes between the surprisal values at neighboring quintiles, yielding five slope values for each curve. For the second strategy, we split each sentence length by number of words into sections



based on those areas of the information curves that had seemed most important in our CHILDES analysis: between the first and second word; between the second and third word; between the third word and third-to-last word; between the third-to-last word and the second to last word; and between the second-to-last word and the last word. We then computed surprisal values at each of these positions, and computed slopes between the surprisal values at each section, giving us another five slope values for each language summarizing the information curves. We computed frequency-based and trigram contextual information curves for each language, using the aggregation strategies described above.

We include illustrations of these two strategies using the frequency-based curve from the British National Corpus in Figure 7

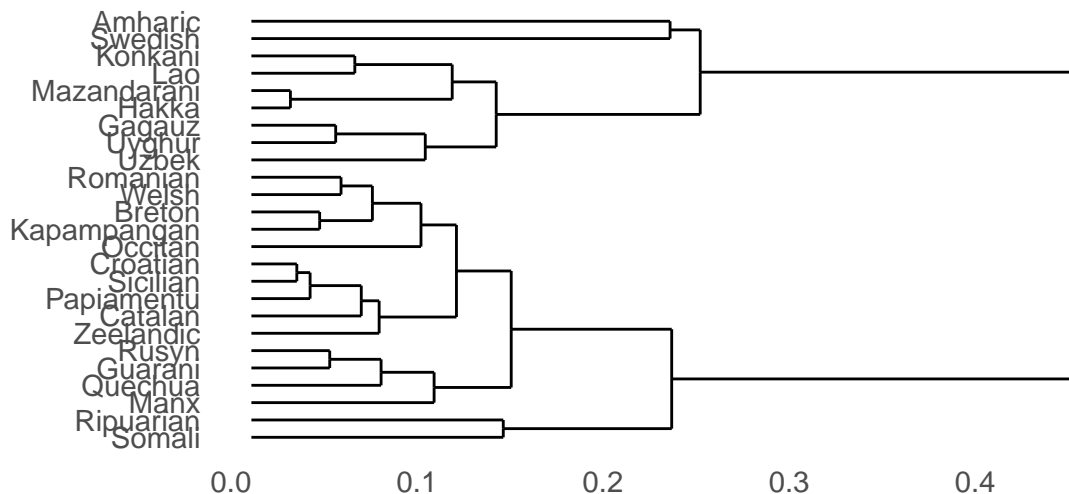
To more rigorously described the typological differences between languages, we used data from the World Atlas of Language Structures (WALS; Dryer & Haspelmath, 2013). The WALS database has data for 144 typological features in 2569 languages from across the world. These features describe aspects of morphology, syntax, phonology, etymology and semantics—in short the features describe the structures in each language. As WALS is a compiled database from dozens of papers from different authors, most of the features and languages are fairly sparse. Even limiting ourselves to the 159 language corpora we pulled from Wikipedia and 122 features from WALS, there are nearly 20000 individual possible data values, fewer than half of which were already computed for those languages in the WALS database.

To fill in the missing data for the features we selected using statistical imputation, we used Multiple Imputation Multiple Correspondence Analysis (MIMCA; Audigier, Husson, & Josse, 2017). MIMCA begins with mean imputation, converts the categorical WALS features into a numerical contingency table with dummy coding, then repeatedly performs principle components analysis and reconstructs the contingency table. Our final result from the MIMCA algorithm was a fully imputed table with 122 feature values for each language.

However, the WALS features describe specific structural differences between languages, while our surprisal metric is word-based. To target lexical differences between languages, we computed the average normalized Levenshtein distance (LDN; Holman et al., 2008) over the 40 item Swadesh list (Swadesh, 1955), retrieved from the ASJP database (Wichmann et al., 2016). The Swadesh list is designed to include near-universal words that target basic cognitive concepts, and are useful in determining the genealogical similarities and differences between languages. The results of classifying languages using the Swadesh list and LDN are correlated with those using WALS features, but the Swadesh list and LDN do not suffer from the same sparsity problem as WALS (Holman et al., 2008).

### Language structures and large-scale data analysis: results

We ran a hierarchical clustering algorithm on the frequency-based information curves using the `hclust` package from the R stats core library (Team & others, 2013). We used the complete linkage algorithm for hierarchical clustering, with distances between information curves between languages computed using cosine distance between their embeddings in the slope space. The complete linkage algorithm at every step pairs each language or cluster of languages with its closest neighboring language or cluster. A sample from the dendrogram is shown in Figure ?? . From a quick glance, the unigram information curves appear to reproduce some of the genealogical relationships between languages, although the dendrogram does not exactly replicate language genealogy for all 159 languages. This suggests using a first-pass quantitative method that the information curves do correspond in some measure to language families, but language families do not explain all of the variation and relationships between frequency-based information curves.



A sample of the contextual information curves (computed using two words of context) are plotted in Figure 8, and all contextual information curves for the languages we used follow the same pattern. The first few words in utterances for each language are surprising, but after even two words of predictive context for each word, the amount of information in each word flattens. Regardless of language, speakers produce information at a constant, optimal rate.

For our first quantitative analysis, we examined the effects of individual typological features on the shapes of the unigram information curves. We ran logistic regressions using the lme4 package in R (Bates, Mächler, Bolker, & Walker, 2014), checking whether the cosine distance between two languages' embeddings in the slope space played a role in determining if those two languages had the same value for a given WALS feature. Individual WALS features do not necessarily have ordinal values. Some, such as the “Number of Cases” feature, are easy to quantify and order. Others are more difficult. For example, how does one order “relative clauses appear after the nouns they modify”, “relative clauses appear before the nouns they modify” and “free order of relative clauses and nouns”? We chose the identify relation to avoid deciding on the basis of individual features. We found that 100 out of the 122 features from WALS we examined were statistically significant ( $p < .001$ ) in determining whether two languages had the same frequency-based information curve shape.

The results for some important features are in Figure 9.

We next compared how the cosine distance between two languages related to how many WALS features they had in common.  $r^2$  value is .005734, which suggests that in aggregate there is not a correlation between how many WALS features languages have in common and the similarity of their frequency-based information curves. Figure 10 displays the results. This result is surprising based on the significance of many WALS features in predicting the shapes of the frequency-based information curves, and we return to this result in the general discussion.

For lexical features, we see a stronger correlation between the similarity of two languages in terms of their average LDN and the cosine distance between their information curves. Figure 11. We see a higher  $r^2$  value here of .026, indicating that there is more correspondence between a language's lexical similarity to another language and their similarity in information curves. From these typological and lexical investigations, we conclude that the shape of a language's frequency-based information curve covaries with its typological and lexical similarity to other languages. However, most of the variation in frequency-based information curve shapes is not explained by typological properties in language.

### **Mutual information between WALS features and barycenters**

We use a final measure to quantify how much variation in WALS features explains variation in barycenters: mutual information between each coordinate of the barycenters and each WALS feature and feature group. We use the categorical-continuous measure derived and described in Ross (2014). How much does knowing the value of each WALS feature tell you about the value of each one of the barycenter coordinates? We reduce to a pairwise measure. The algorithm uses a variant of the k-nearest-neighbors regression algorithm. It essentially asks: how many points with the same WALS label are in a neighborhood around

each language's e.g. first coordinate?

Discussion of results.

Results and plot.

## Discussion

By considering the distribution of information at the level of utterances and sentences, we join together the information-theoretic work focusing on sub-word units and words, and that focusing on paragraphs. In doing so, we show that frequency and context-based metrics complement one another in studying efficiency and information in language. We directly link linguistic efficiency in a language to the genealogy and properties of that language. We provide evidence for a novel linguistic universal: low processing cost for listeners beyond the first words in utterances, driven by high average word predictability in conversation. With consideration to language acquisition, we observe that children tend to distribute information in their utterances according to the their language's frequency-based information curve as soon as they form multi-word utterances.

Throughout this work we have averaged the surprisal values at each position. Averaging removes variation, which in turn may obscure trends in the data. As discussed in the methods section, the surprisal metric has historically been used for calculating the information and processing cost for individual utterances, and our use of the metric here is actually a step forward rather than a step back. Future work can investigate variation in how speakers distribute information in individual utterances.

The WALS database we used to investigate typological variation in the information curves is overall sparse. We imputed well over 50% of the WALS features for most of our 159 languages, although all of the languages had at least 20 features evaluated in WALS. A large part of this is due to WALS being a collection of a number of different studies, instead of a

systematic effort to catalogue variation across the world's languages. Additionally, WAL  
features are meant to describe specific microvariations in languages, not to provide a  
comprehensive typological representation of each language compared to each other language.  
This may be why the Swadesh list provided a higher correlation for describing the differences  
in information curves: Swadesh (1955) intended the list to allow researchers to more  
comprehensively compared and contrast lexical differences between languages. For our  
Wikipedia analysis, we also reduce all of a language's variation down to a five-dimensional  
vector. These information curve representations show a surprising amount of variation  
despite the degree of compression.

## References

- Anderson, J. R., & Milson, R. (1989). Human memory: An adaptive perspective. *Psychological Review*, 96(4), 703.
- Audigier, V., Husson, F., & Josse, J. (2017). MIMCA: Multiple imputation for categorical variables with multiple correspondence analysis. *Statistics and Computing*, 27(2), 501–518.
- Austin, J. L. (1975). *How to do things with words*. Oxford university press.
- Aylett, M., & Turk, A. (2004). The smooth signal redundancy hypothesis: A functional explanation for relationships between redundancy, prosodic prominence, and duration in spontaneous speech. *Language and Speech*, 47(1), 31–56.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2014). Fitting linear mixed-effects models using lme4. *arXiv Preprint arXiv:1406.5823*.
- Chen, S. F., & Goodman, J. (1999). An empirical study of smoothing techniques for language modeling. *Computer Speech & Language*, 13(4), 359–394.
- Dryer, M. S., & Haspelmath, M. (Eds.). (2013). *WALS online*. Leipzig: Max Planck Institute for Evolutionary Anthropology. Retrieved from <https://wals.info/>
- Frank, A. F., & Jaeger, T. F. (2008). Speaking rationally: Uniform information density as an optimal strategy for language production. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 30).
- Frank, S. L., Otten, L. J., Galli, G., & Vigliocco, G. (2015). The erp response to the amount of information conveyed by words in sentences. *Brain and Language*, 140, 1–11.
- Genzel, D., & Charniak, E. (2002). Entropy rate constancy in text. In *Proceedings of the*

493       40th annual meeting of the association for computational linguistics (pp. 199–206).

494   Heafield, K., Pouzyrevsky, I., Clark, J. H., & Koehn, P. (2013). Scalable modified kneser-ney  
495       language model estimation. In *Proceedings of the 51st annual meeting of the*  
496       *association for computational linguistics (volume 2: Short papers)* (pp. 690–696).

497   Holman, E. W., Wichmann, S., Brown, C. H., Velupillai, V., Müller, A., Bakker, D., &  
498       others. (2008). Advances in automated language classification. *Quantitative*  
499       *Investigations in Theoretical Linguistics*, 40–43.

500   Jaeger, T. F., & Levy, R. P. (2007). Speakers optimize information density through syntactic  
501       reduction. In *Advances in neural information processing systems* (pp. 849–856).

502   Kemp, C., & Regier, T. (2012). Kinship categories across languages reflect general  
503       communicative principles. *Science*, 336(6084), 1049–1054.

504   Kirby, S., Tamariz, M., Cornish, H., & Smith, K. (2015). Compression and communication  
505       in the cultural evolution of linguistic structure. *Cognition*, 141, 87–102.

506   Kutas, M., & Federmeier, K. D. (2011). Thirty years and counting: Finding meaning in the  
507       n400 component of the event-related brain potential (erp). *Annual Review of*  
508       *Psychology*, 62, 621–647.

509   Leech, G. N. (1992). 100 million words of english: The british national corpus (bnc).

510   Levy, R. (2008). Expectation-based syntactic comprehension. *Cognition*, 106(3), 1126–1177.

511   MacWhinney, B. (2000). *The chldes project: The database* (Vol. 2). Psychology Press.

512   Mahowald, K., Fedorenko, E., Piantadosi, S. T., & Gibson, E. (2013). Info/information  
513       theory: Speakers choose shorter words in predictive contexts. *Cognition*, 126(2),  
514       313–318.



- Perfors, A., & Navarro, D. J. (2014). Language evolution can be shaped by the structure of the world. *Cognitive Science*, 38(4), 775–793.
- Petitjean, F., Ketterlin, A., & Gançarski, P. (2011). A global averaging method for dynamic time warping, with applications to clustering. *Pattern Recognition*, 44(3), 678–693.
- Piantadosi, S. T., Tily, H., & Gibson, E. (2011). Word lengths are optimized for efficient communication. *Proceedings of the National Academy of Sciences*, 108(9), 3526–3529.
- Pickering, M. J., & Garrod, S. (2013). An integrated theory of language production and comprehension. *Behavioral and Brain Sciences*, 36(4), 329–347.
- Ross, B. C. (2014). Mutual information between discrete and continuous data sets. *PloS One*, 9(2).
- Sanchez, A., Meylan, S. C., Braginsky, M., MacDonald, K. E., Yurovsky, D., & Frank, M. C. (2019). Chiles-db: A flexible and reproducible interface to the child language data exchange system. *Behavior Research Methods*, 51(4), 1928–1941.
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27(3), 379–423.
- Smith, N. J., & Levy, R. (2013). The effect of word predictability on reading time is logarithmic. *Cognition*, 128(3), 302–319.
- Swadesh, M. (1955). Towards greater accuracy in lexicostatistic dating. *International Journal of American Linguistics*, 21(2), 121–137.
- Tanenhaus, M. K., Spivey-Knowlton, M. J., Eberhard, K. M., & Sedivy, J. C. (1995). Integration of visual and linguistic information in spoken language comprehension. *Science*, 268(5217), 1632–1634.

- 537 Team, R. C., & others. (2013). R: A language and environment for statistical computing.
- 538 Wichmann, S., Müller, A., Wett, A., Velupillai, V., Bischoffberger, J., Brown, C. H., . . .
- 539 others. (2016). The asjp database. *Max Planck Institute for the Science of Human*
- 540 *History, Jena*.
- 541 Yu, S., Cong, J., Liang, J., & Liu, H. (2016). The distribution of information content in
- 542 english sentences. *arXiv Preprint arXiv:1609.07681*.

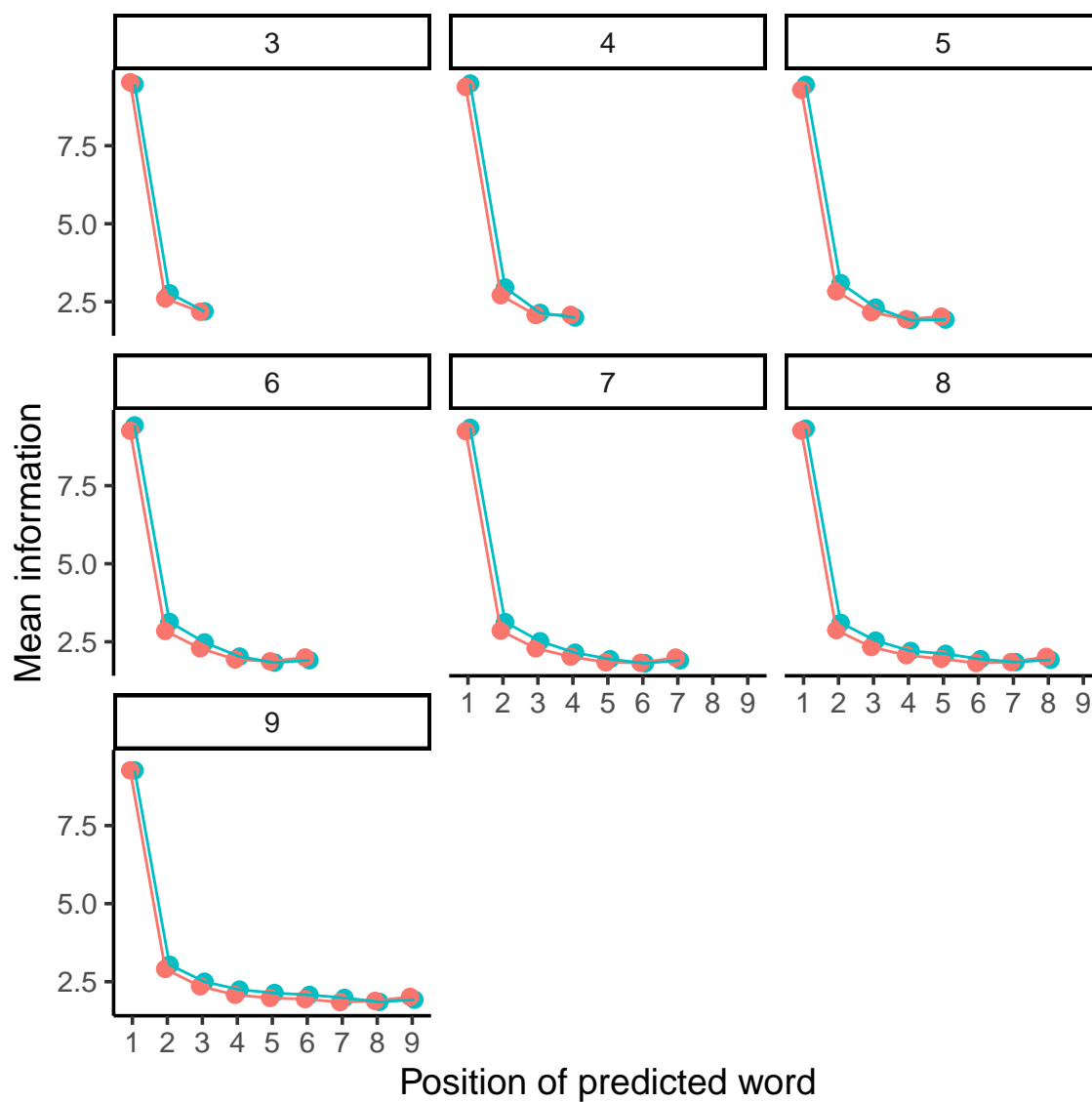


Figure 4. North American English context-based information curves. Lines around each point indicate 95% confidence intervals computed with non-parametric bootstrap

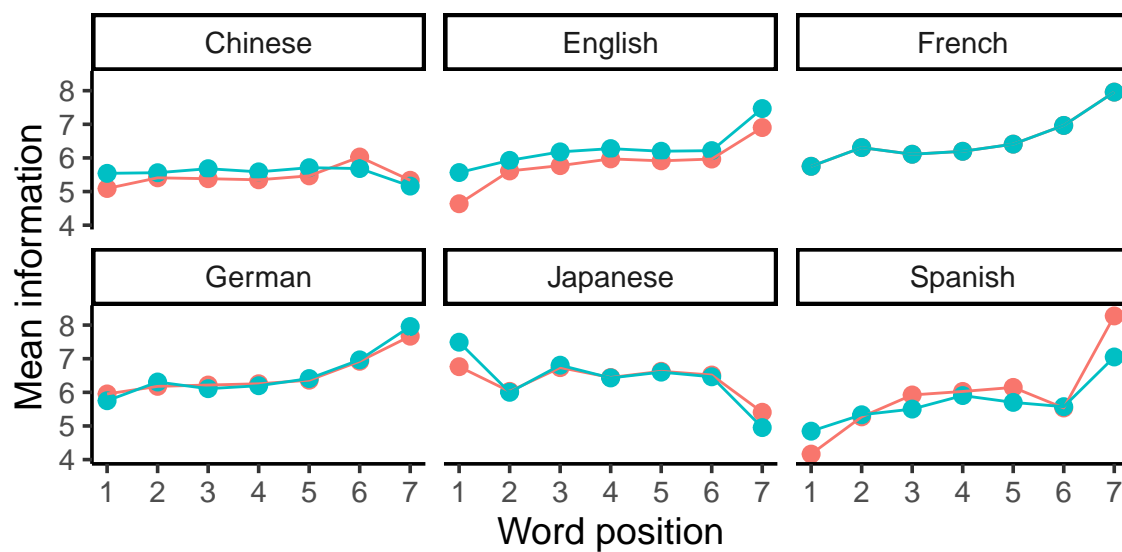


Figure 5. CHILDES frequency-based information curves

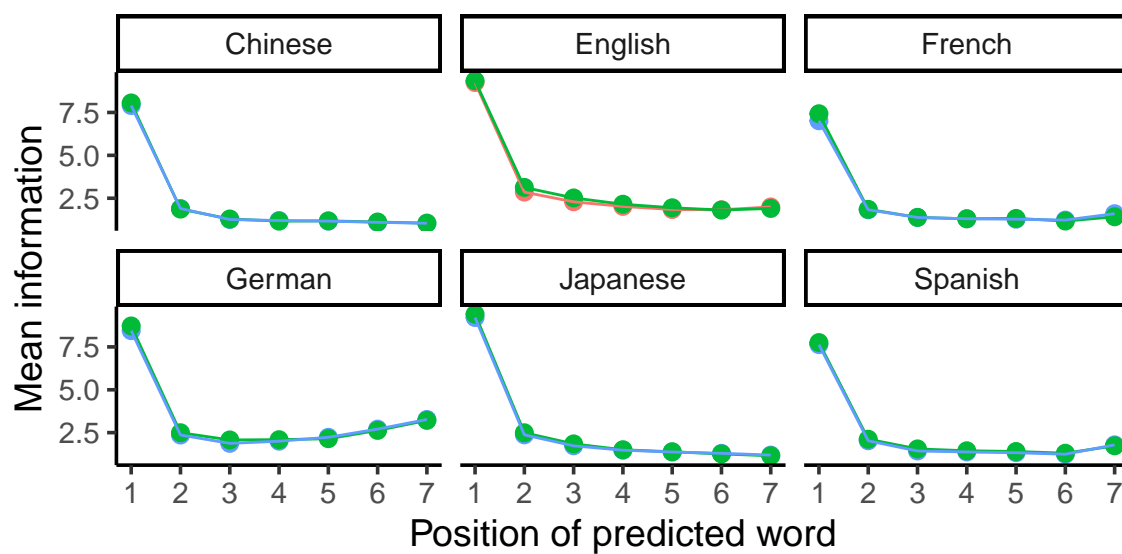


Figure 6. CHILDES trigram context-based information curves

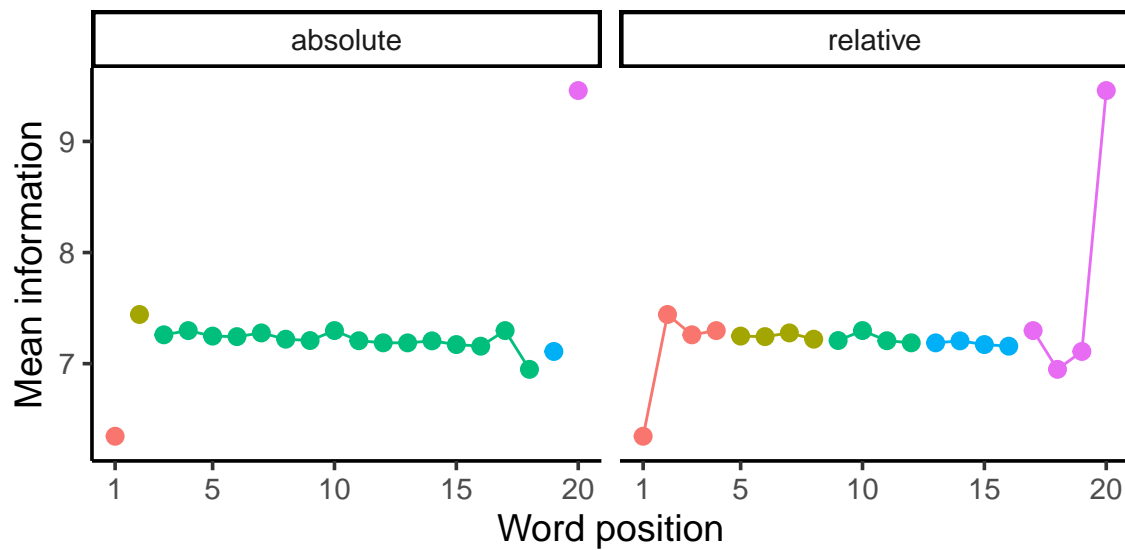


Figure 7. Illustration of slope treatments for Wikipedia information curves: relative on top and absolute on bottom

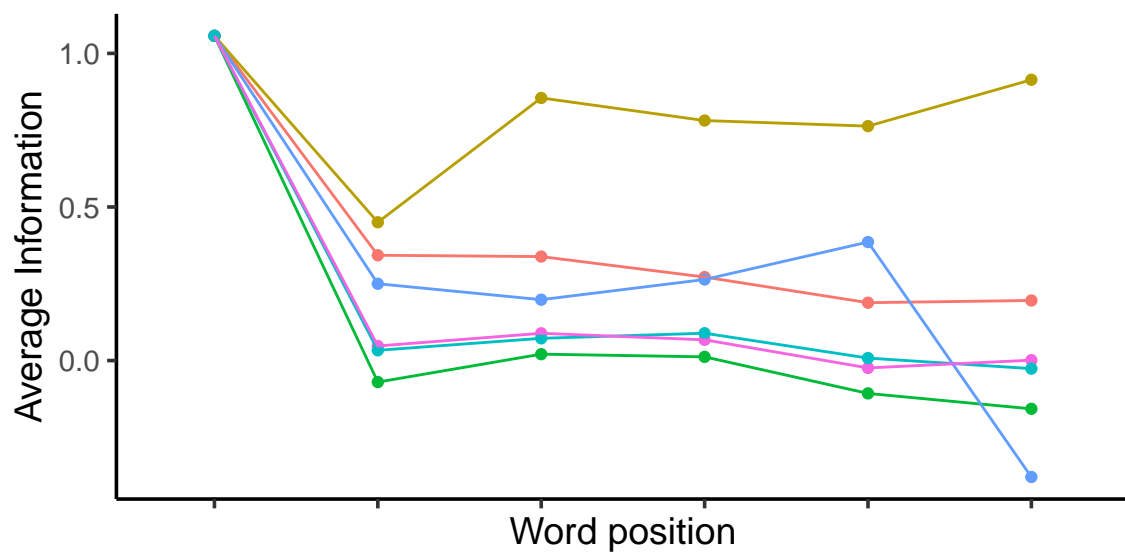


Figure 8. Some trigram information curves from the Wikipedia data

	Some Important Features	Examples
1	Order of Relative Clause and Noun	Noun–RC; RC–noun
2	Number of Cases	No cases; 6 cases
3	Order of Subject, Object and Verb	SOV; VSO
4	The Morphological Imperative	Only singular; sing. and plural
5	Definite Articles	Affix; distinct word
6	Position of Case Affixes	Prefixes; suffixes

Figure 9. Some linear model results from Wikipedia and WALS features

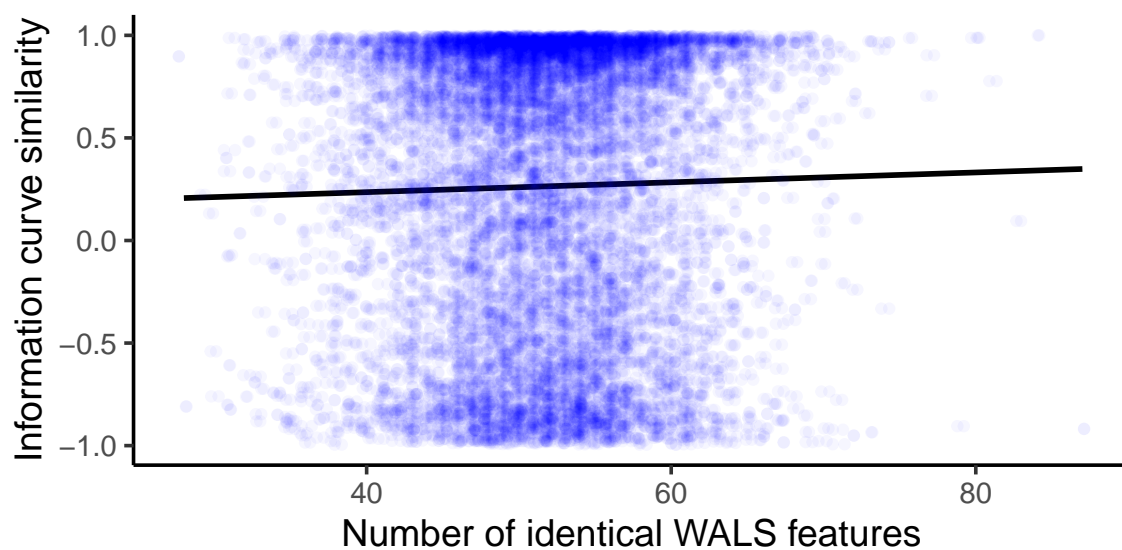


Figure 10. wals features vs cosine similarity

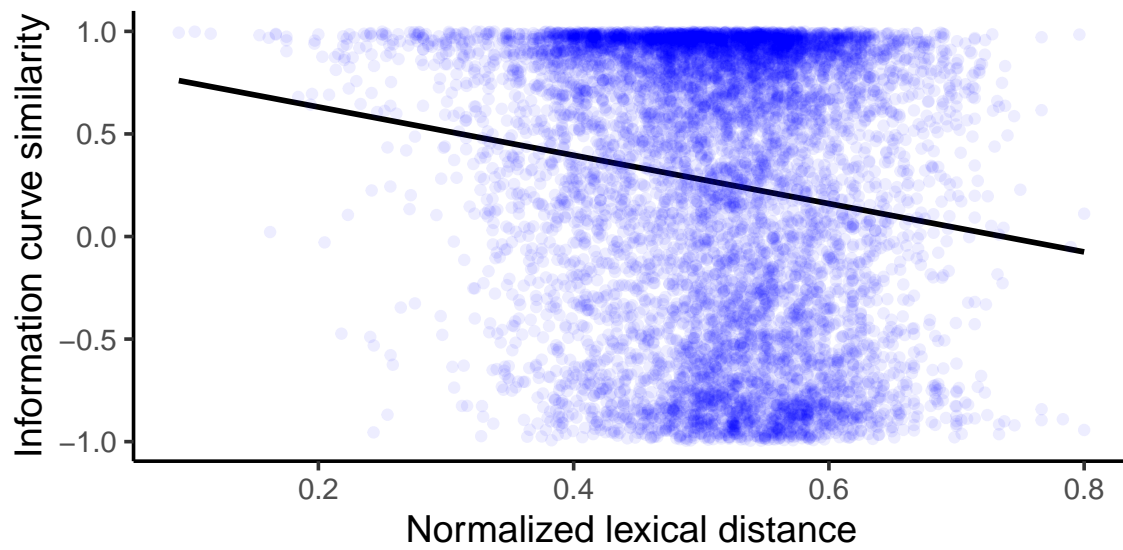


Figure 11. ldn features vs cosine similarity