

Is syntax optimized for periodic neurobiological sampling? Evidence from 21 languages

Chia-Wen Lo¹, Mark Anderson², Lorenzo Titone¹ & Lars Meyer^{1,3}

¹Research Group Language Cycles, Max Planck Institute for Human Cognitive and Brain Sciences, Leipzig, German

²Norwegian Computing Center, Oslo, Norway

³Clinic for Phoniatrics and Pedaudiology, University Clinic Münster, German

lo@cbs.mpg.de

Human memory is short-lived. To avoid information loss, humans need to chunk continuous speech into short units, each containing multiple words. Temporal limitations in line with this were observed in auditory short-term memory (Baddeley et al., 1975) and sentence comprehension in particular (Frazier and Fodor, 1978). A recent proposal has suggested that timing constraints may stem from the limited duration of underlying electrophysiological processing windows. In electroencephalography studies, it has been shown that the wavelength of slow-frequency rhythmic brain activity that may synchronize with stimuli—so-called *neural oscillations*—serves the formation of multi-word chunks (Meyer et al., 2016), but also limits chunk duration (Henke and Meyer, 2021).

If the wavelength of periodic brain activity indeed sets a neuronal timing constraint on multi-word chunking, this should be reflected in the languages of the world: Multi-word chunks should exhibit temporal regularity. In other words, neurobiology should set a limit on the evolution of the dependency structure in the languages of the world. This limit should be observable as periodicity and a tendency for isochrony across different languages. If this were the case, the distribution of inter-chunk intervals (ICIs) should be highly non-uniform, with a consistent median ICI across languages. That is, non-uniform ICIs would mean that chunks tend towards isochrony. We here employed the Universal Dependencies (UD; Nivre et al., 2020) to test this for 21 languages.

We first applied a chunker defined by dependency grammar, employing the Universal part-of-speech tags and dependency annotations to generate multi-word chunks for each sentence (Lo et al., 2023; Figure 1). The chunker generates base-level subtrees, which allow for language-specific annotations. Then the potential candidate chunks are extracted and thus result in some overlapping chunks. Final output chunks are selected by the higher value of mutual information between nodes if there are overlapping chunks. Essentially, output chunks are locally saturated dependency graph-lets; these would implicitly allow the listener for establishing all dependency relations within a single filling of the working memory buffer. We then synthesized the chunks with the Google WaveNet text-to-speech engine to get real-time audio and computed ICIs.

To test periodicity and consistency across languages, we assess the median of ICIs. We expect consistent medians across languages, as brains and processing windows are the same across humans and determined by genetic factors, irrespective of their language. Hence, all languages may have to conform to this duration constraint. We then assess whether ICI distributions are non-uniform, exhibiting low variance. We consider the Fano factor and the coefficient of variation, both of which are often used to measure temporal regularity of binary event series (e.g., neuronal spike trains). When the Fano factor of a binary time series (here: chunk onsets vs. non-onsets) is low, this means that chunk rate is less variable than expected by non-onsets; when the coefficient of variation is low, this means that chunk spacing is less variable than expected by random chunks. Post-hoc, to capture typological variance, a series of correlational analysis are also performed between ICI metrics and word-order predictability from Hahn et al. (2020). We specifically want to test whether ICIs and variance relate to word order flexibility.

Across languages, we found non-uniform distributions of ICIs across languages (Figure 2) and the average of median ICIs across languages is 1.1 seconds. This suggests that chunks within and across languages are periodic, at least to some extent. As for variance of distribution for each language, we found decreased ICI variances and Fano factors relative to surrogate data generated from 1,000 permutations of random chunks (Table 1). These results suggest that chunk spacing is less variable than expected for a fully random, non-isochronous chunk rate, and also less variable than expected when assuming randomness for each language. Correlational analysis (Figure 3) revealed positive correlations between word-order predictability and ICI median ($r(19) = .48, p = .03$) and Fano factor ($r(19) = .49, p = .03$), but not coefficient of variation ($r(19) = .21, p = .36$); these did not withstand FDR correction.

Our results provide evidence that multi-word chunks are periodic across languages. This periodicity may have evolved to fit the wavelength of endogenous oscillatory activity, thus optimizing chunking in the face of memory limitations. Our correlational results suggest that there may be trade-off between periodicity and speech content. Languages with more flexible word orders might have shorter chunks and more periodicity, suggesting that the ontogenetic neural constraint can be bent by cultural differences, within limits.

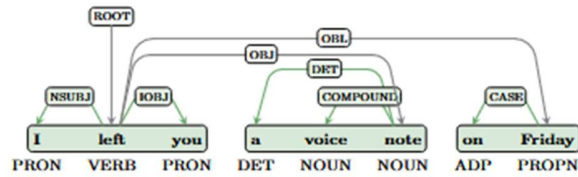


Figure 1. Formalization of chunks defined by dependency grammar.

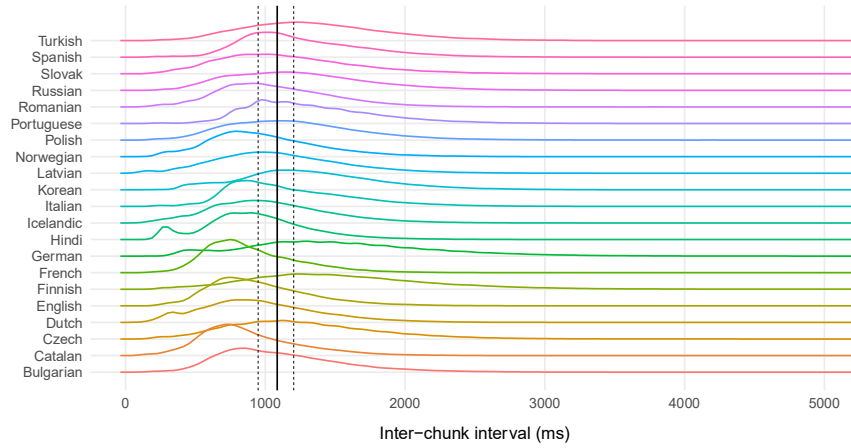


Figure 2. ICI histograms by language.

Language	Sentences (n)	ICI (ms)	Coefficient of variation		Fano factor	
			Observed	Cutoff	Observed	Cutoff
Bulgarian	10,120	1,040	0.41	0.47	4.88E+09	5.00E+09
Catalan	15,348	820	0.42	0.48	3.41E+09	3.41E+09
Czech	75,634	1,202	0.43	0.48	4.50E+08	4.75E+09
German	173,199	1,464	0.44	0.49	3.21E+09	3.28E+09
English	12,190	905	0.42	0.48	5.61E+09	5.79E+09
Spanish	16,482	1,180	0.36	0.43	2.14E+08	2.18E+09
Finnish	12,917	1,471	0.45	0.50	4.08E+09	4.80E+09
French	15,272	837	0.39	0.45	3.85E+09	3.91E+09
Hindi	16,572	910	0.41	0.46	3.03E+09	3.07E+09
Icelandic	42,112	1,000	0.41	0.46	4.74E+09	4.80E+09
Italian	13,123	1,085	0.38	0.45	3.68E+09	3.77E+09
Korean	25,550	1,284	0.38	0.45	3.50E+09	3.63E+09
Latvian	14,573	1,064	0.42	0.48	6.07E+09	6.15E+09
Dutch	11,889	930	0.45	0.50	7.43E+09	7.56E+09
Norwegian	17,527	949	0.41	0.47	5.30E+09	5.42E+09
Polish	20,648	1,160	0.40	0.48	4.35E+09	4.46E+08
Portuguese	10,705	1,254	0.36	0.43	2.36E+09	2.42E+09
Romanian	25,193	1,025	0.38	0.44	3.92E+09	3.97E+09
Russian	77,748	1,200	0.40	0.48	3.93E+09	4.04E+09
Slovak	8,257	1,098	0.43	0.49	6.34E+09	6.57E+09
Turkish	15,456	1,315	0.38	0.48	4.40E+09	4.53E+08

Table 1. Coefficient of variation and Fano factor by language.

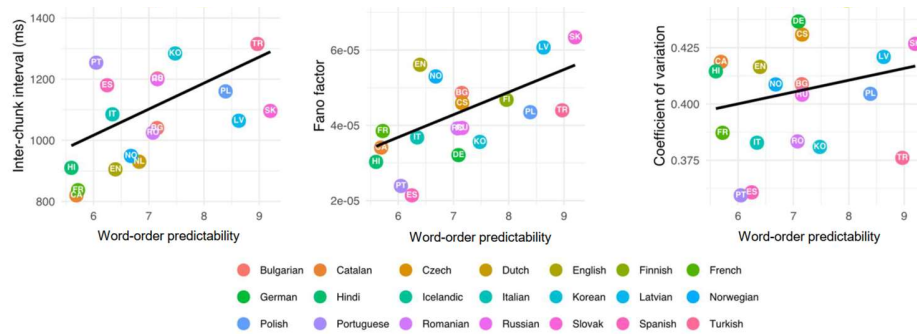


Figure 3. Correlations with word-order predictability.