

1 Oscillatory tracking of pseudo-rhythmic speech is constrained by linguistic
2 predictions

3 Sanne Ten Oever^{1,2,3,*} & Andrea E. Martin^{1,2}

4

5 ¹ Language and Computation in Neural Systems group, Max Planck Institute for Psycholinguistics,
6 Nijmegen, The Netherlands

7 ² Language and Computation in Neural Systems group, Donders Centre for Cognitive Neuroimaging,
8 Nijmegen, The Netherlands

9 ³ Department of Cognitive Neuroscience, Faculty of Psychology and Neuroscience, Maastricht University,
10 The Netherlands

11

12 *Corresponding author: sanne.tenoever@mpi.nl

13

14 **Abstract**

15 Neuronal oscillations putatively track speech in order to optimize sensory processing. However, it
16 is unclear how isochronous brain oscillations can track pseudo-rhythmic speech input. Here we
17 propose that oscillations can track pseudo-rhythmic speech when considering that speech time is
18 dependent on predictions flowing from internal language models. We show that the temporal
19 dynamics of speech are dependent on the predictability of words in a sentence. A computational
20 model including oscillations, feedback, and inhibition is able to track the natural pseudo-rhythmic
21 speech input. As the model processes, it generates temporal phase codes, which are a candidate
22 mechanism for carrying information forward in time. The model is optimally sensitive to the
23 natural temporal speech dynamics and can explain empirical data on temporal speech illusions. Our
24 results reveal that speech tracking does not only rely on the input acoustics but instead entails an
25 interaction between oscillations and constraints flowing from internal language models.

26

27 **Introduction**

28 Speech is a biological signal that is characterized by a plethora of temporal information. The
29 temporal relation between subsequent speech units allows for the online tracking of speech in order
30 to optimize processing at relevant moments in time [1-7]. Neural oscillations are a putative index
31 of such tracking [3, 8]. The existing evidence for neural tracking of the speech envelope is consistent
32 with such a functional interpretation [9, 10]. In these accounts, the most excitable optimal phase of
33 an oscillation is aligned with the most informative time-point within a rhythmic input stream [8,
34 11-14]. However, the range of onset time difference between speech units seems more variable than
35 fixed oscillations can account for [15-17]. As such, it remains an open question how is it possible
36 that oscillations can track a signal that is at best only pseudo-rhythmic [16].

37 Oscillatory accounts tend to focus on the prediction in the sense of predicting “when,”
38 rather than predicting “what”: oscillations function to align the optimal moment of processing given
39 that timing is predictable in a rhythmic input structure. If rhythmicity in the input stream is
40 violated, oscillations must be modulated to retain optimal alignment to incoming information. This
41 can be achieved through phase resets [15, 18], directly coupling of the acoustics to oscillations [19],
42 or the use of many oscillators at different frequencies [2]. However, the optimal or effective time
43 of processing stimulus input might not only depend on when you predict something to occur, but
44 also on what stimulus is actually being processed [20-23].

45 What and when are not independent, and certainly not from the brain’s-eye-view. If
46 continuous input arrives to a node in an oscillatory network, the exact phase at which this node
47 reaches threshold activation does not only depend on the strength of the input, but also on how
48 sensitive this node was to start with. Sensitivity of a node in a language network (or any neural
49 network) is naturally affected by predictions in the what domain generated by an internal language
50 model [24-27]. If a node represents a speech unit that is likely to be spoken next, it will be more
51 sensitive and therefore active earlier, that is, on a less excitable phase of the oscillation. In the
52 domain of working memory, this type of phase precession has been shown in rat hippocampus [28,
53 29] and more recently in human electroencephalography [30]. In speech, phase of activation and
54 perceived content are also associated [31-35] and phase has been implicated in tracking of higher-
55 level linguistic structure [18, 36, 37]. However, the direct link between phase and the predictability
56 flowing from a language model has yet to be established.

57 The time of speaking/speed of processing is not only a consequence of how predictable a
58 speech unit is within a stream, but also a cue for the interpretation of this unit. For example,
59 phoneme categorization depends on timing (e.g., voice onsets, difference between voiced and
60 unvoiced phonemes), and there are timing constraints on syllable durations (e.g. the theta syllable
61 [19, 38] that affect intelligibility [39]. Even the delay between mouth movements and speech audio
62 can influence syllabic categorizations [20]. Most oscillatory models use oscillations for parsing, but
63 not as a temporal code for content [40-43]. However, the time or phase of presentation does
64 influence content perception. This is evident from two temporal speech phenomena. In the first
65 phenomena, the interpretation of an ambiguous short /α/ or long vowel /a:/ depends on speech rate
66 (in Dutch; [44-46]). Specifically, when speech rates are fast the stimulus is interpreted as a long
67 vowel and vice versa for slow rates. However, modulating the entrainment rate effectively changes
68 the phase at which the target stimulus - which is presented at a constant speech rate - arrives (but
69 this could not be confirmed in [47]). A second speech phenomena shows the direct phase-
70 dependency of content [31, 34]. Ambiguous /da-/ga/ stimuli will be interpreted as a /da/ on one
71 phase and a /ga/ on another phase. This was confirmed in both a EEG as well as a behavioral study.
72 An oscillatory theory on speech tracking should account for how temporal properties in the input
73 stream can alter what is perceived.

74 In the speech production literature, there is strong evidence that the onset times (as well as
75 duration) of an uttered word is modulated by the frequency of that word in the language [48-52]
76 showing that internal language models modulate the access to or sensitivity of a word node [24, 53].
77 This word-frequency effect relates to the access to a single word. However, it is likely that during
78 ongoing speech internal language models use the full context to estimate upcoming words [54]. If
79 so, the predictability of a word in context should provide additional modulations on speech time.
80 Therefore, we predict that words with a high predictability in the producer's language model should
81 be uttered relatively early. In this way word-to-word onset times map to the predictability level of
82 that word within the internal model. Thus, not only the processing time depends on the
83 predictability of a word (faster processing for predictable words; see [55, 56] and [57] showing that
84 speech time in noise matters), but also the production time (earlier uttering of predicted words).

85 Language comprehension involves the mapping of speech units from a producer's internal
86 model to the speech units of the receiver's internal model. In other words, one will only understand
87 what someone else is writing or saying if one's language model is sufficiently similar to the speakers

88 (and if we speak in Dutch, fewer people
89 will understand us). If the producer's
90 and receiver's internal language model
91 have roughly matching top-down
92 constrains they should similarly
93 influence the speed of processing (either
94 in production or perception; Figure 1A-
95 C). Therefore, if predictable words arrive
96 earlier (due to high predictability in the
97 producer's internal model), the receiver
98 also expects the content of this word to
99 match one of the more predictable ones
100 from their own internal model (Figure
101 1C). Thus, the phase of arrival depends
102 on the internal model of the producer
103 and the expected phase of arrival
104 depends on the internal model of the
105 receiver (Figure 1D). If this is true,
106 pseudo-rhythmicity is fully natural to
107 the brain and it provides a means to use
108 time or arrival phase as a content
109 indicator. It also allows the receiver to be
110 sensitive to less predictable words when
111 they arrive relatively late. Current
112 oscillatory models of speech parsing do
113 not integrate the constraints flowing
114 from an internal linguistic model into
115 the temporal structure of the brain
116 response. It is therefore an open question
117 whether the oscillatory model the brain
118 employs is actually attuned to the
119 temporal variations in natural speech.

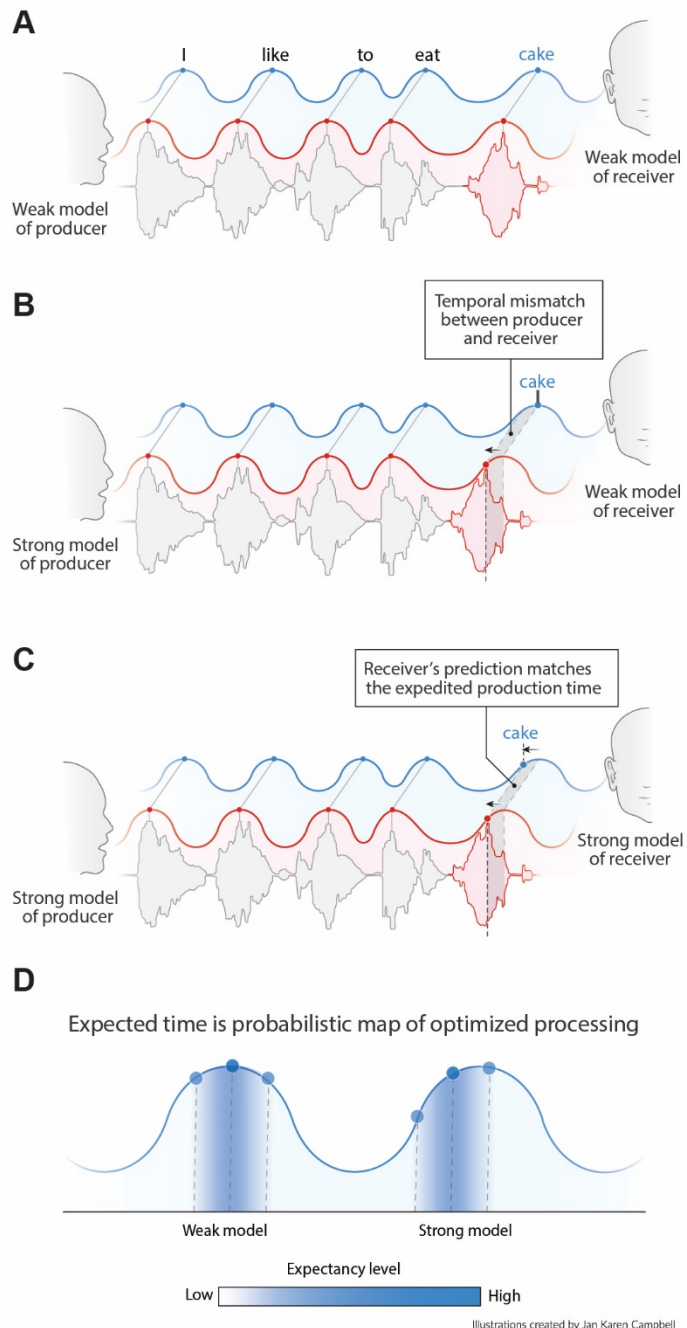


Figure 1. Proposed interaction between speech timing and internal linguistic models. A) Isochronous production and expectation when there is a weak internal model (even distribution of node activation). All speech units arrive around the most excitable phase B) When the internal model of the producer does not align with the model of the receiver temporal alignment and optimal communication fails. C) When both producer and receiver have a strong internal model, speech is non-isochronous and not aligned to the most excitable phase, but fully expected by the brain. D) Expected time is a constraint distribution which center can be shifted due to linguistic constraints.

120 Here, we propose that neural oscillations can track pseudo-rhythmic speech by taking into
121 account that speech timing is a function of linguistic constrains. As such we need to demonstrate
122 that speech statistics are influenced by linguistic constrains as well as showing how oscillations can
123 be sensitive to this property in speech. We approach this hypothesis as follows: First, we
124 demonstrate that in natural speech timing depends on linguistics predictions (*temporal speech*
125 *properties*). Then, we model how oscillations can be sensitive to these linguistic predictions
126 (*modeling speech tracking*). Finally, we validate that this model is optimally sensitive to the natural
127 temporal properties in speech and displays temporal speech illusions (*model validation*). Our results
128 reveal that tracking of speech needs to be viewed as an interaction between ongoing oscillations as
129 well as constraints flowing from an internal language model [21, 24]. In this way, oscillations do
130 not have to shift their phase after every speech unit and can remain at a relatively stable frequency
131 as long as the internal model of the speaker matches the internal model of the perceiver.

132

133 Results

134 Temporal speech properties: word frequency influences word duration

135 To extract the temporal properties in naturally spoken speech we used the Corpus Gesproken
136 Nederlands (CGN; (Version 2.0.3; 2014)). This corpus consists of elaborated annotations of over 900
137 hours of spoken Dutch and Flemish words. We focus here on the subset of the data of which onset
138 and offset timings were manually annotated at the word level in Dutch. Cleaning of the data
139 included removing all dashes and backslashes. Only words were included that were part of a Dutch
140 word2vec embedding (github.com/coosto/dutch-word-embeddings; needed for later modeling) and
141 required to have a frequency of at least 10 in the corpus. All other words were replaced with an
142 <unknown> label. This resulted in 574,726 annotated words with 3096 unique words. 2848 of the
143 words were recognized in the Dutch Wordforms database in CELEX (Version 3.1) in order to
144 extract the word frequency as well as the number of syllables per word. Mean word duration was
145 0.392 seconds, with an average standard deviation of 0.094 seconds (Supporting Figure 1A). By
146 splitting up the data in sequences of 10 sequential words we could extract the average word,
147 syllable, and character rate (Figure Supporting Figure 1B). The reported rates fall within the
148 generally reported ranges for syllables (5.2 Hz) and words (3.7 Hz; [5, 58]).

149 We predict that knowledge about the language statistics influences the duration of speech
150 units. As such we predict that more prevalent words will have on average a shorter duration (also
151 reported in [50]). In Figure 2A the duration of several mono- and bi-syllabic words are listed with
152 their word frequency. From these examples it seems that words with higher word frequency
153 generally have a shorter duration. To test this statistically we entered word frequency in an
154 ordinary least square regression with number of syllables as control. Both number of syllables
155 (coefficient = 0.1008, $t(2843) = 75.47$, $p < 0.001$) as well as word frequency (coefficient = -0.022,
156 $t(2843) = -13.94$, $p < 0.001$) significantly influence the duration of the word. Adding an interaction
157 term did not significantly improve the model ($F(1,2843) = 1.320$, $p = 0.251$; Figure 2B+C). The effect
158 is so strong that words with a low frequency can last three times as long as high frequency words
159 (even within mono-syllabic words). This indicates that word frequency could be an important part
160 of an internal model that influences word duration.

161 The previous analysis probed us to expand on the relation between word duration and
162 length of the words. Obviously, there is a strong correlation between word length and mean word
163 duration (number of characters 0.824, $p < 0.001$; number of syllables: $\rho = 0.808$, $p < 0.001$; for

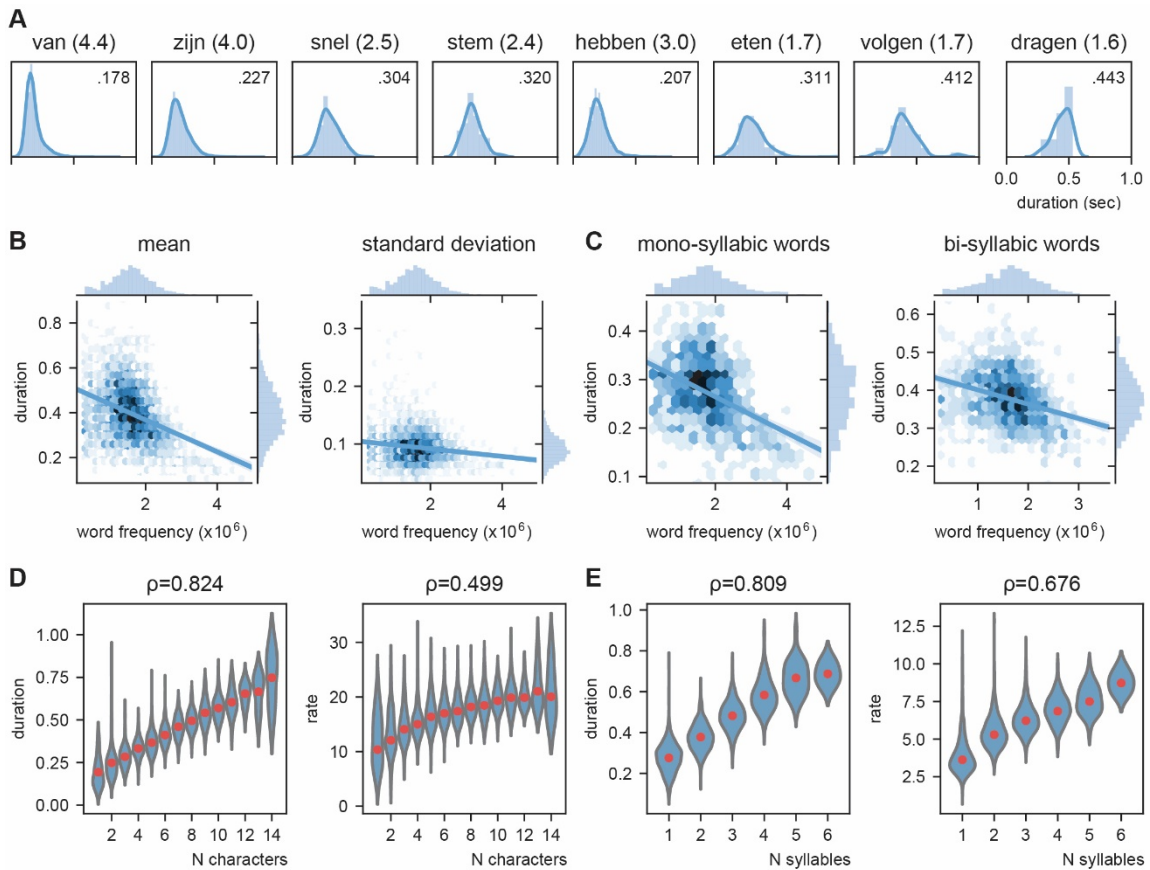


Figure 2. Word frequency modulates word duration. A) Example of mono- and bi-syllabic words of different word frequencies in brackets (van=from, zjn=be, snel=fast, stem=voice, hebben=have, eten=eating, volgen=next, toekomst=future). Text in the graph indicates the mean word duration. B) Relation between word frequency and duration. Darker colors mean more values. C) same as B) but separately for mono- and bi-syllabic words. D) Relation character amount and word duration. The longer the words, the longer the duration (left). The increase in word duration does not follow a fixed number per character as duration as measured by rate increases. E) same as D) but for number of syllables. Red dots indicate the mean.

164 number of syllables already shown above; Figure 2D+E). In contrast, this correlation is present, but
 165 much lower for the standard deviation of word duration (number of characters: $\rho = 0.269$, $p < 0.001$;
 166 number of syllables: $\rho = 0.292$, $p < 0.001$). Finding a strong correlation does not imply that for every
 167 time unit increase in the word length, the duration of the word also increases with the same time
 168 unit, i.e., bi-syllabic words do not necessarily have to last twice as long as mono-syllabic words.
 169 Therefore, we recalculated word duration to a rate unit considering the number of syllables/
 170 characters of the word. Thus a 250 ms mono- versus bi-syllabic word would have a rate of 4 versus
 171 8 Hz respectively. Then we correlated character/syllabic rate with word duration. If word duration
 172 increases monotonically with character/syllable length there should be no correlation. We found
 173 that the syllabic rate varies between 3 and 8 Hz as previously reported (Figure 2E right; [5, 58]).
 174 However, the more syllables there are in a word, the higher this rate ($\rho = 0.676$, $p < 0.001$). This
 175 increase was less strong for the character rate ($\rho = 0.499$, $p < 0.001$; Figure 2D right).

176 These results show that the syllabic/character rate depends on the number of characters
177 /syllables within a word and is not an independent temporal unit [38]. This effect is easy to explain
178 when assuming that the prediction strength of an internal model influences word duration:
179 transitional probabilities of syllables are simply more constrained within a word than across words
180 [59]. This will reduce the time it takes to utter/perceive any syllable which is later in a word.
181 Unfortunately, the CGN does not have separate syllable annotations to investigate this possibility
182 directly. However, we can investigate the effect of transitional probabilities and other statistical
183 regularities flowing from internal models across words (see next section and [17] for statistical
184 regularities in syllabic processing).

185

186 **Temporal speech properties: word-by-word predictability predicts word onset differences**

187 The brain's internal model likely provides predictions about what linguistic features and
188 representations, and possibly about which specific units, such as words, to expect next when
189 listening to ongoing speech [21, 24]. As such, it is also expected that word-by-word onset delays are
190 shorter for words that fit the internal model (i.e. those that are expected; [54]). To investigate this
191 possibility, we created a simplified version of an internal model predicting the next word using
192 recurrent neural nets (RNN). We trained an RNN to predict the next word from ongoing sentences
193 (Figure 3A). The model consisted of an embedding layer (pretrained; [github.com/coosto/dutch-](https://github.com/coosto/dutch-word-embeddings)
194 [word-embeddings](https://github.com/coosto/dutch-word-embeddings)), a recurrent layer with a tanh activation function, and a dense output layer with
195 a softmax activation. To prevent overfitting, we added a 0.2 dropout to the recurrent layers and the
196 output layer. An adam optimizer was used at a 0.001 learning rate and a batch size of 32. We
197 investigated four different recurrent layers (GRU and LSTM at either 128 or 300 units; see
198 Supporting Figure 4). The final model we use here includes a LSTM with 300 units. Input data
199 consist of 10 sequential words (label encoding) within the corpus (of a single speaker; shifting
200 the sentences by one word at a time), and an output consisted of a single word. A maximum of four
201 unknown labeled words (words not included in the word2vec estimations. Four was chosen as it
202 was < 50% of the words). was allowed in the input, but not in output. Validation consisted of a
203 randomly chosen 2% of the data.

204 The output of the RNN reflects a probability distribution in which the values of the RNN
205 sum up to one and each word has its own predicted value (Figure 3A). As such we can extract the

206 predicted value of the
 207 uttered word and relate the
 208 RNN prediction with the
 209 stimulus onset delay relative
 210 to the previous word. We
 211 entered word prediction in a
 212 regression using the
 213 stimulus onset difference
 214 between the current word in
 215 the sentence and the
 216 previous word (i.e. onset
 217 difference of words). We
 218 added the control variables
 219 bigram (using the NLTK
 220 toolbox based on the
 221 training data only),
 222 frequency of previous word,
 223 syllable rate (rate of the full
 224 sentence input), and mean
 225 duration of previous word
 226 (all variables that can
 227 account for part of the
 228 variance that affects the
 229 duration of the last word).
 230 We only used the test data
 231 (total of 7361 sentences,
 232 excluding all word not
 233 present in Celex. 4837

234 sentences). Many of the variables were skewed to the right, therefore we transformed the data
 235 accordingly (see Table 1; results were robust to changes in these transformation).

236 All predictors except word frequency of the previous word showed a significant effect
 237 (Table 1). The variance explained by word frequency was likely captured by the mean duration

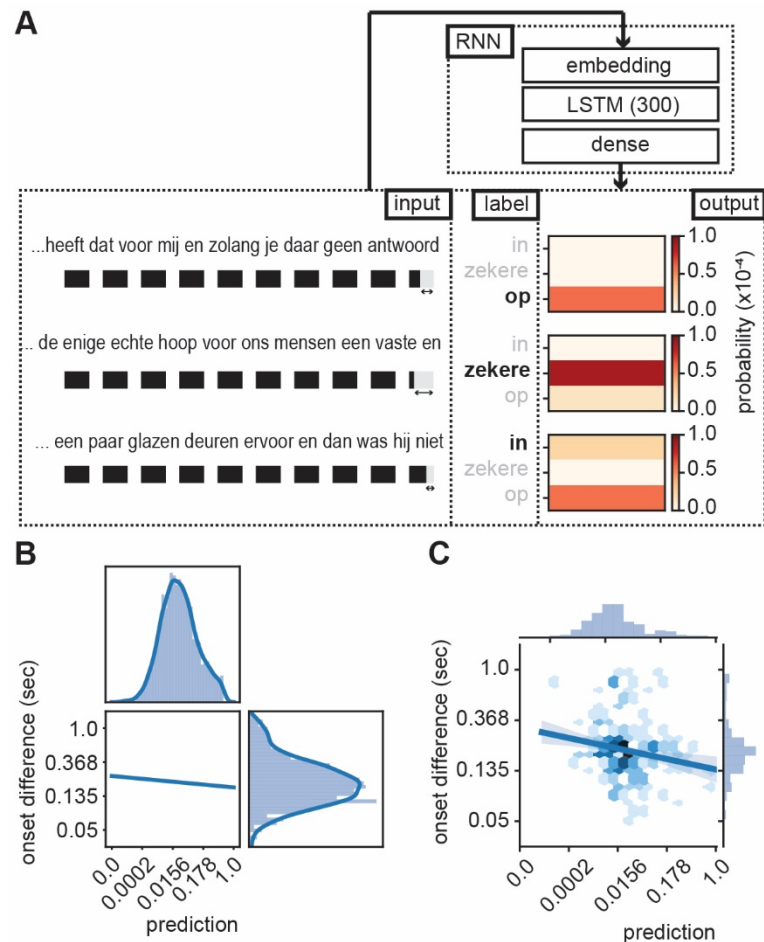


Figure 3. RNN output influence word onset differences. A) Sequences of ten words were entered in an RNN in order to predict the content of the next word. Three examples are provided of input data with the label (bold word) and probability output for three different words. The regression model showed a relation between the duration of last word in the sequence and the predictability of the next word such that words were systematically shorter when the next word was more predictable according to the RNN output (illustrated here with the shorted black boxes). B) Regression line estimated at mean value of word duration and bigram. C) Scatterplot of prediction and onset difference of data within ± 0.5 standard deviation of word duration and bigram. Note that for B and C the axes are linear on the transformed values. Translation of the sentences in A from top to bottom: “... that it has for me and while you have no answer [on]”, “... the only real hope for us humans is a firm and [sure]”, “... a couple of glass doors in front and then it would not have been [in]”.

238 variable of the
239 previous word
240 which is
241 correlated to
242 word frequency.
243 The RNN
244 predictor could
245 capture more

Table 1. Summary of regression model for logarithm of onset difference of words

Variable	Trans	B	β	SE	t	p	VIF
Intercept	x	0.9719		0.049	19.764	<0.001	
RNN prediction	$x^{(1/6)}$	-0.3370	-0.0862	0.047	-7.163	<0.001	1.5
Bigram	$\log(x)$	-0.0118	-0.0316	0.005	-2.424	0.015	1.8
Word frequency W-1	x	0.0049	0.0076	0.009	0.546	0.585	2.0
Mean duration W-1	$\log(x)$	1.1206	0.7003	0.022	50.326	<0.001	2.0
Syllable Rate	x	-0.1033	-0.2245	0.004	-23.014	<0.001	1.0

Model $R^2 = 0.542$. Trans = transformation, W-1 = previous word, B = unstandardized coefficient, β = standardized coefficient, SE = standard error, t = t value, p = p value, VIF = variance inflation factor

246 variance than the bigram model suggesting that word duration is modulated by the level of
247 predictability within a fuller context than just the conditional probability of the current word given
248 the previous word (Figure 3B+C). Importantly, it was necessary to use the trained RNN model as a
249 predictor; entering the RNN predictions after the first epoch did not results in a significant predictor
250 ($t(4837) = -1.191$, $p = 0.234$). Also adding the predictor word frequency of the current word did not
251 add significant information to the model ($F(1, 4830) = 0.2048$, $p = 0.651$). These results suggest that
252 words are systematically lengthened (or pauses are added. However, the same predictors are also
253 significant when excluding sentences containing pauses) when the next word is not strongly
254 predicted by the internal model.

255

256 **Modeling speech tracking: Speech Tracking in a Model Constrained Oscillatory Network (STiMCON)**

257 In order to investigate how much of these duration effects can be explained using an oscillator
258 model, we created the model Speech Tracking in a Model Constrained Oscillatory Network
259 (STiMCON). STiMCON in its current form will not be exhaustive; however, it can extract how
260 much an oscillating network can cope with asynchronies by using its own internal model
261 illustrating how the brain's language model and speech timing interact [60]. The current model is
262 capable of explaining how top-down predictions can influence the processing time as well as
263 provide an explanation for two known temporal illusions in speech.

264 STiMCON consists of a network of semantic nodes of which the activation A of each level
265 in the model l is governed by:

$$266 \quad A_{l,T} = C_{l-1 \rightarrow l} * A_{l-1,T} + C_{l+1 \rightarrow l} * A_{l+1,T} + inhib(Ta) + osc(T) \quad (1)$$

267 in which C represents the connectivity patterns between different hierarchical levels, T the time
 268 in a sentence, and Ta the vector of times of an individual node in an inhibition function (in
 269 milliseconds). The inhibition function is a gate function:

$$270 \quad \text{inhib}(Ta) = \begin{cases} -3 * \text{BaseInhib}, & Ta < 20 \\ 3 * \text{BaseInhib}, & 20 \leq Ta < 100 \\ \text{BaseInhib}, & Ta > 100 \end{cases} \quad (2)$$

271 in which BaseInhib is a constant for the base level of inhibition (negative value, set to -0.2). As such
 272 nodes are by default inhibited, as soon as they get activated above threshold (activation threshold
 273 set at 1) Ta sets to zero. Then, the node will have suprathreshold activation, which after 20
 274 milliseconds returns to increased inhibition until the base level of inhibition is returned. The
 275 oscillation is a constant oscillator:

$$276 \quad \text{osc}(T) = Am * e^{2\pi i \omega T + i\varphi} \quad (3)$$

277 in which Am is the amplitude of the oscillator, ω the frequency, and φ the phase offset. As such we
 278 assume a stable oscillator which is already aligned to the average speech rate (see [15, 19] for phase
 279 alignment models). The model used for the current simulation has one an input layer (l-1 level) and
 280 one single layer of semantic word nodes (l level) that receives feedback from a higher level layer
 281 (l+1 level). As such only the word (l) level is modeled according to equation 1-3 and the other levels
 282 form fixed input and feedback connection patterns.

283

284 **Modeling speech tracking: language models influence time of activation**

285 To illustrate how STIMCON can explain how processing time depends on the prediction of internal
 286 language models, we instantiated a language model that had only seen three sentences and five
 287 words presented at different probabilities (I eat cake at
 288 0.5 probability, I eat nice cake at 0.3 probability, I eat
 289 very nice cake at 0.2 probability; Table 2). This language
 290 model will serve as the feedback arriving from the l+1-
 291 level to the l-level. The l-level consists of five nodes that
 292 each represent one of the words and receives
 293 proportional feedback from l+1 according to Table 2
 294 with a delay of $0.9 * \omega$ seconds, which then decays at 0.01

Table 2. Example of a language model

	I	eat	very	nice	cake
I	0	1	0	0	0
eat	0	0	0.2	0.3	0.5
very	0	0	0	1	0
nice	0	0	0	0	1
cake	0	0	0	0	0

This model has seen three sentences at different probabilities. Rows represent the prediction for the next word, e.g. /I/ predicts /eat/ at a probability of 1, but after /eat/ there is a wider distribution.

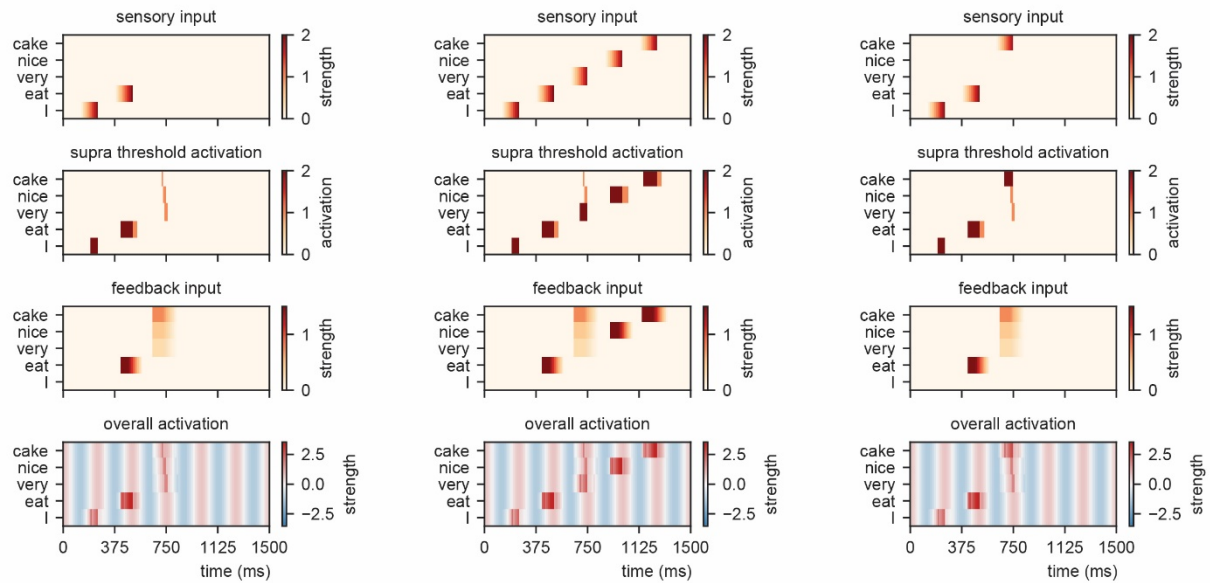


Figure 4. Model output for different sentences. For the supra-threshold activation dark red indicates activation which included input from $I+I$ as well as $I-I$, orange indicates activation due to $I+I$ input.

295 unit per millisecond and influences the l-level at a proportion of 1.5. This feedback is only initiated
 296 when supra-activation arrives due to l-1-level bottom-up input. Each word at the l-1-level input is
 297 modelled as a linearly function to the individual nodes lasting length of 125 milliseconds (half a
 298 cycle, ranging from 0-1 arbitrary units). As such, the input is not the acoustic input itself but rather
 299 reflects a linear increase representing the increasing confidence of a word representing the specific
 300 node. φ is set such that the peak of a 4 Hz oscillation aligns to the peak of sensory input of the first
 301 word. Sensory input is presented at a base stimulus onset asynchrony of 250 milliseconds (i.e. 4 Hz).

302 When we present this model with different sensory input at an isochronous rhythm of 4
 303 Hz it is evident that the timing at which different nodes reach activation depends on the level of
 304 feedback that is provided (Figure 4). For example, while the /I/-node needs a while to get activated
 305 after the initial sensory input, the /eat/-node is activated earlier as it is pre-activated due to
 306 feedback. After presenting /eat/ the feedback arrives at three different nodes and the activation
 307 timing depends on the stimulus that is presented (earlier activation for /cake/ compared to /very/).

308

309 **Modeling speech tracking: time of presentation influences processing efficiency**

310 To investigate how the time of presentation influences the processing efficiency we presented the
 311 model with /I eat XXX/ in which the last word was varied in content (either /I/, /very/, /nice/, or
 312 /cake/), intensity (linearly ranging from 0 to 1), and onset delay (ranging between -125 to +125

313 relative to isochronous presentation). We extracted the time at which the node matching the
 314 stimulus presentation reached activation threshold first (relative to stimulus onset, and relative to
 315 isochronous presentation).

316 Figure 5A shows the output. When there is no feedback (i.e. at the first word /I/
 317 presentation), a classical efficiency map can be found in which processing is most optimal (possible
 318 at lowest stimulus intensities) at isochronous (in phase with the stimulus rate) presentation and
 319 then drops to either side. For nodes that have feedback, input processing is possible at earlier times
 320 relative to isochronous presentation and parametrically varies with prediction strength (earlier for

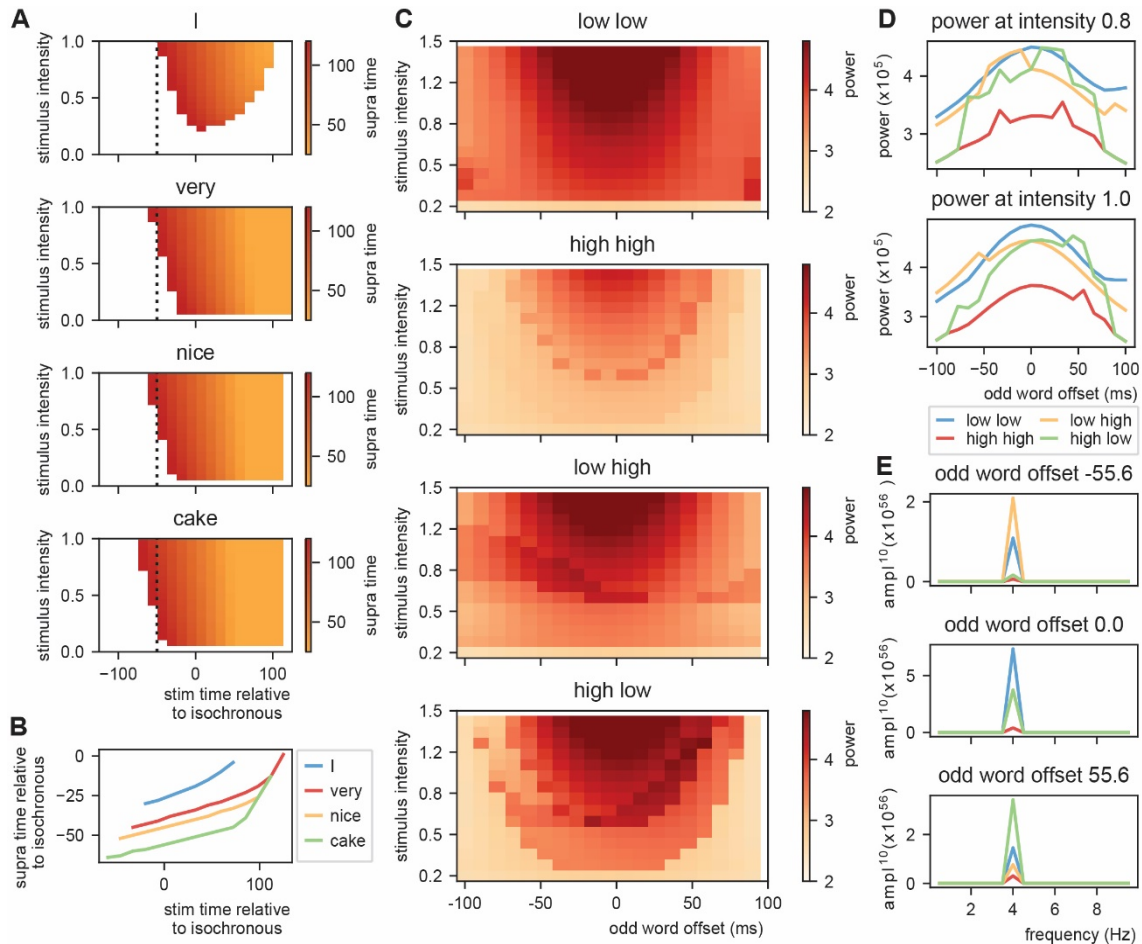


Figure 5. Model output on processing efficiency and rhythmicity. A) Time of presentation influences efficiency. Outcome variable is the time at which the node reached threshold activation (supra-time). The dashed line is presented to ease comparison between the four content types. White indicates that threshold is never reached. B) Same as A, but estimated at a threshold of 0.53 showing that oscillations regulate feedforward timing. Panel A shows that the earlier the stimuli are presented (on a weaker point of the ongoing oscillation), the longer it takes until supra-threshold activation is reached. This figure shows that timing relative to the ongoing oscillation is regulated such that the stimulus activation timing is closer to isochronous. Line discontinuities are a consequence of stimuli never reaching threshold for a specific node. C) Strength of 4 Hz power depends on predictability in the stream. When predictability is alternated between low and high, activation is more rhythmic when the predictable odd stimulus arrives earlier and vice versa. D) Slice of D at intensity of 0.8 and 1.0. E) Magnitude spectra at three different odd word offsets at 1.0 intensity. To more clearly illustrate the differences the magnitude to the power of 20 is plotted.

321 /cake/ at 0.5 probability, then /very/ at 0.2 probability). Additionally, the activation function is
322 asymmetric. This is a consequence of the interaction between the supra-activation caused by the
323 feedback and the sensory input. As soon as supra-activation is reached due to the feedback, sensory
324 input at any intensity will reach supra-activity (thus at early stages of the linearly increasing
325 confidence of the input). This is why for the /very/ stimulus activation is still reached at later delays
326 compared to /nice/ and /cake/ as the /very/-node reaches supra-activation due to feedback at a later
327 time point.

328 When we investigate timing differences in stimulus presentation it is important to also
329 consider what this means for the timing in the brain. Before, we showed that the amount of
330 prediction can influence timing in our model. It is also evident that the earlier a stimulus was
331 presented the more time it took (relative to the stimulus) for the nodes to reach threshold (more
332 yellow colors for earlier delays). This is a consequence of the oscillation still being at a relatively
333 low excitability point at stimulus onset for stimuli that are presented early during the cycle.
334 However, when we translate these activation threshold timing to the timing of the ongoing
335 oscillation, the variation is strongly reduced (Figure 5B). A stimulus timing that varies between 130
336 milliseconds (e.g. from -59 to +72 in the /cake/ line; excluding the non-linear section of the line)
337 only reaches the first supra-threshold response with 19 milliseconds variation in the model
338 (translating to a reduction of 53% to 8% of the cycle of the ongoing oscillation, i.e. a 1:6.9 ratio).
339 This means that within this model (and any oscillating model) the activation of nodes is robust to
340 some timing variation in the environment. This effect seemed weaker when no prediction was
341 present (for the /I/ stimulus this ratio was around 1:3.5. Note that when determining the /cake/
342 range using the full line the ratio would be 1:3.4).

343

344 **Modeling speech tracking: top-down interactions can provide rhythmic processing for non-** 345 **isochronous stimulus input**

346 The previous simulation demonstrate that oscillations provide a temporal filter and the processing
347 itself can actually be closer to isochronous than what can be solely extracted from the stimulus
348 input. Next, we investigated whether dependent on changes in top-down prediction, processing
349 within the model will be more or less rhythmic. To do this, we create stimulus input of 10 sequential
350 words at a base rate of 4 Hz to the model with constant (low at 0 and high at 0.8 predictability) or

351 alternating word-to-word predictability. For the alternating conditions word-to-word
352 predictability alternates between low to high (sequences which word are predicted at 0 or 0.8
353 predictability, respectively) or shift from high to low. For this simulation we used Gaussian sensory
354 input (with a standard deviation of 42 ms aligning the mean at the peak of the ongoing oscillation;
355 see Supporting Figure 5 for output with linear sensory input). Then, we vary the onset time of the
356 odd words in the sequence (shifting from -100 up to +100 ms) and the stimulus intensity (from 0.2
357 to 1.5). We extracted the overall activity of the model and computed the Fast Fourier transform of
358 the created time course (using a Hanning taper only including data from 0.5 – 2.5 seconds to exclude
359 the onset responses).

360 The first thing that is evident is that the model with no content predictions has overall
361 stronger power, and specifically around isochronous presentation (odd word offset of 0 ms) at high
362 stimulus intensities (Figure 5C-E). Adding overall high predictability drops the power, but also here
363 the power seems symmetric around zero. The spectra of the alternating predictability conditions
364 look different. For the low to high predictability condition the curve seems to be shifted to the left
365 such that 4 Hz power is strongest when the predictable odd stimulus is shifted to an earlier time
366 point (low-high condition). This is reversed for the high-low condition. At middle stimulus
367 intensities there is a specific temporal specificity window at which the 4 Hz power is particularly
368 strong. This window is earlier for the low-high than the high-low alternation (Figure 5D, Figure
369 5E, and Supporting Figure 6). The effect only occurs at specific middle intensity combination as at
370 high intensities the stimulus dominates the responses and at low intensities the stimulus does not
371 reach threshold activation. These results show that even though stimulus input is non-isochronous,
372 the interaction with the internal model can still create a potential rhythmic structure in the brain
373 (see [61, 62]). Note that the direction in which the brain response is more rhythmic matches with
374 the natural onset delays in speech (shorter onset delays for more predictable stimuli).

375

376 **Model validation: STiMCON's sinusoidal modulations of RNN predictions is optimally sensitive to**
377 **natural onset delays**

378 Next, we aimed to investigated whether STiMCON would be optimally sensitive to speech input
379 timings found naturally in speech. Therefore, we tried to fit STiMCON's expected word-to-word
380 onset differences to the word-to-word onset differences we found in the CGN. At a stable level of

381 intensity of the input and inhibition, the only aspect that changes the timing of the interaction
 382 between top-down predictions and bottom-up input within STiMCON is the ongoing oscillation.
 383 Considering that we only want to model for individual words how much the prediction ($C_{l-1 \rightarrow l} * A_{l-1, T}$)
 384 influences the expected timing we can set the contribution of the other factors from
 385 equation (1) to zero remaining with the relative contribution of prediction:

$$386 \quad C_{l+1 \rightarrow l} * A_{l+1, T} = \text{topdown influence} = -\text{osc}(T) \quad (4)$$

387 We can solve this formula in order to investigate the expected relative time shift (T) in processing
 388 that is a consequence of the strength of the prediction (ignoring that in the exact timing will also
 389 depend on the strength of the input and inhibition):

$$390 \quad \text{relative time shift} = \frac{1}{2\pi\omega} \left(\arcsin \left(\frac{C_{l+1 \rightarrow l} * A_{l+1, T}}{-Am} \right) - \varphi \right) \quad (5)$$

391 ω was set as the syllable rate for each sentence, Am and φ were systematically varied. We fitted a
 392 linear model between the STiMCON's expected time and the actual word-to-word onset
 393 differences. This model was similar to the model described in the section *word-by-word*
 394 *predictability predicts word onset differences* and included the predictor syllabrate and duration
 395 of the previous word. However, as we were interested in how well non-transformed data matches
 396 the natural onset timings we did not perform any normalization besides equation (5). As this might
 397 involve violating some of the assumptions of the ordinary least square fit, we estimate model
 398 performance by repeating the regression 1000 times fitting it on 90% of the data (only including
 399 the test data from the

400 RNN) and extracting R^2
 401 from the remaining 10%.

402 Results show a
 403 modulation of the R^2
 404 dependent on the
 405 amplitude and phase
 406 offset of the oscillation
 407 (Figure 6A) which was
 408 stronger than the non-
 409 transformed R^2 (which
 410 was 0.389). This suggests

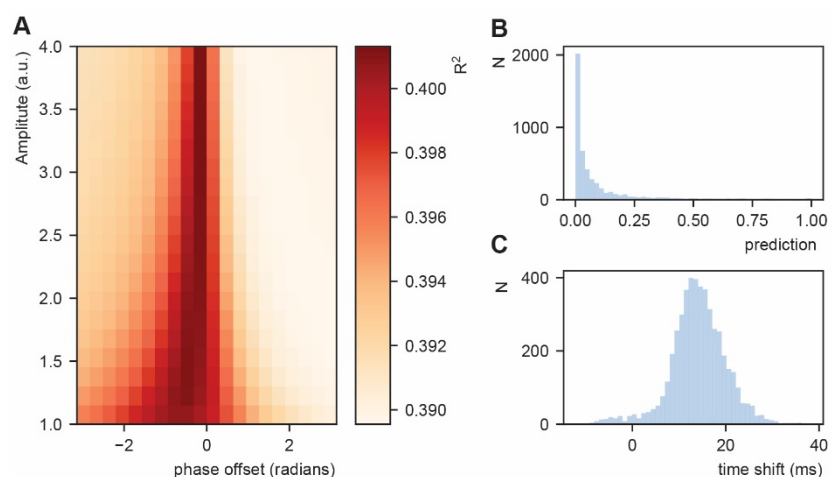


Figure 6. Fit between real and expected time shift dependent on predictability. A) Phase offset and amplitude of the oscillation modulate the fit to the word-to-word onset durations. B) Histogram of the predictions created by the deep neural net. C) Histogram of the relative time shift transformation at phase of -0.15π and amplitude of 1.5.

411 that STiMCON expected time durations matches the actual word-by-word duration. This was even
412 more strongly so for specific oscillatory alignments (around -0.25π offset) suggesting an optimal
413 alignment phase relative to the ongoing oscillation is needed for optimal tracking [3, 8].
414 Interestingly, the optimal transformation seemed to automatically alter a highly skewed prediction
415 distribution (Figure 6B) towards a more normal distribution of relative time shifts (Figure 6C).

416

417 **Model validation: STiMCON can explain perceptual effects in speech processing**

418 Due to the differential feedback strength and the inhibition after suprathreshold feedback
419 stimulation, STiMCON is more sensitive to lower predictable stimuli at phases later in the
420 oscillatory cycle. This property can explain two illusions that have been reported in the literature,
421 specifically, the observation that the interpretation of ambiguous input depends on the phase of
422 presentation [31, 32, 63] and on speech rate [46]. The only assumption that has to be made is that
423 there is an uneven base prediction balance between the ways the ambiguous stimulus can be
424 interpreted.

425 The empirical data we aim to model comprises an experiment in which ambiguous syllables,
426 that could either be interpreted as /da/ or /ga/, were presented [31]. In one of the experiments in
427 this study, broadband stimuli were presented at specific rates to entrain ongoing oscillations. After

428 the last entrainment stimulus an
 429 ambiguous /daga/ stimulus was
 430 presented at different delays
 431 (covering two cycles of the
 432 presentation rate at 12 different
 433 steps), putatively reflecting
 434 different oscillatory phases.
 435 Dependent on the delay of
 436 stimulation participants perceived
 437 either /da/ or /ga/ suggesting that
 438 phase modulates the percept of the
 439 participants. Besides this
 440 behavioral experiment, the authors
 441 also demonstrated that the same
 442 temporal dynamics were present
 443 when looking at ongoing EEG data
 444 showing that the phase of ongoing
 445 oscillations at the onset of
 446 ambiguous stimulus presentation
 447 determined the percept [31].

448 To illustrate that
 449 STiMCON is capable of showing a

450 phase (or delay) dependent effect, we use an internal language model similar to our original model
 451 (Table 2). The model consists of four nodes (N1, N2, Nda, and Nga) at which N1 activation predicts
 452 a second unspecific stimulus (S2) represented by N2 at a predictability of 1. N2 activation predicts
 453 either da or ga at 0.2 and 0.1 probability respectively. Then, we present STiMCON (same parameters
 454 as before) with /S1 S2 XXX/. XXX is varied to have different proportion of the stimulus /da/ and /ga/
 455 (ranging from 0% /da/ to 100% /ga/ in 12 times steps; these reflects relative proportions that sum up
 456 to 1 such that at 30% the intensity of /da/ would be at max 0.3. and of /ga/ 0.7) and is the onset is
 457 varied relate to the second to last word. We extract the time that a node reaches suprathreshold
 458 activity after stimulus onset. If both nodes were active at the same time the node with the highest
 459 total activates was chosen. Results showed that for some ambiguous stimuli, the delay determines

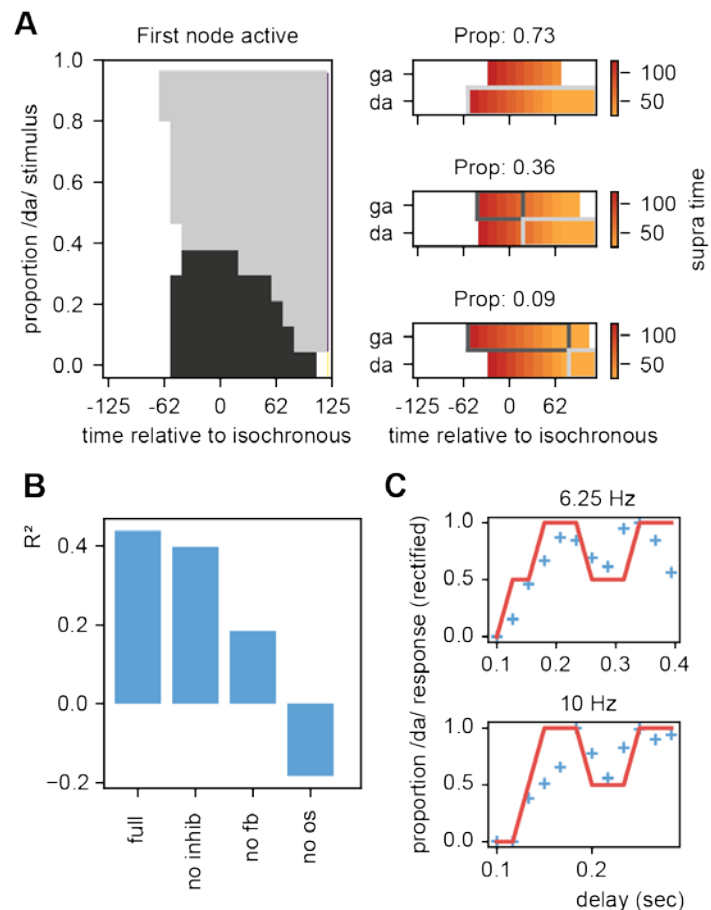


Figure 7. Results for /daga/ illusions. A) Modulations due to ambiguous input at different times. Illustration of the node that is active first. Different proportions of the /da/ stimulus show activation timing modulations at different delays (B). R² for the grid search fit of the full model, a model without inhibition (no inhib), without uneven feedback (no fb), or without an oscillation (no os). C) Fit of the full model on the rectified behavioral data of [31]. Blue crosses indicate rectified data and red lines indicate the fit.

460 which node is activated first, modulating the ultimate percept of the participant (Figure 7A, also
461 see supplementary Figure 7A). The same type of simulation can explain how speech rate can
462 influence perception (supplementary Figure 7B; but see [47].).

463 To further scrutinize on this effect we fitted our model to the behavioral data of Ten Oever
464 & Sack [31]. As we used an iterative approach in the simulations of the model, we optimized the
465 model using a grid search. We varied the parameters of proportion of the stimulus being /da/ or /ga/
466 (ranging between 10:20:80%), the onset time of the feedback (0.1:0.2:1.0 cycle), the speech of the
467 feedback decay (0:0.02:0.1), and a temporal offset of the final sound to account for the time it takes
468 to interpret a specific ambiguous syllable (ranging between -0.05:0.02:0.05 sec). Our outcome
469 variable was the node that show the first suprathreshold activation ($N_{da} = 1$, $N_{ga} = 0$). If both nodes
470 were active at the same time the node with the highest total activates was chosen. If both nodes
471 had equal activation or never reached threshold activation we coded the outcome to 0.5 (i.e. fully
472 ambiguous). These outcomes were fitted to the behavioral data of the 6.25 Hz and 10 Hz presentation
473 rate (the two rates showing a significant modulation of the percept). This data was normalize to
474 have a range between 0-1 to account for the model outcomes being binary (0, 0.5 or 1).

475 We found that our model could fit the data at an average explained variance of 43% (30%
476 and 58% for 6.25 Hz and 10 Hz respectively; Figure 7B+C). This explained variance was higher than
477 the original sinus fit (40% for 3 parameter sinus fit [amplitude, phase offset, and mean]). Note that
478 our fit cannot account for variance ranging inbetween 0-0.5 and 0.5-1, while the sinus fit can do
479 this. If we correct for this (by setting the sinus fit to the closest 0, 0.5 or 1 value and doing a grid
480 search to optimize the fitting) the average fit of the sinus is 21%. The average AIC of the model and
481 sinus fit are -27.0 and -24.1 respectively suggesting that the STiMCON model has the better fit.
482 Thus, STiMCON does better than a fixed-frequency sinus fit. This is a likely consequence of the
483 sinus fit not being able to explain the dampening of the oscillation later (i.e. the perception bias is
484 stronger for shorter compared to longer delays).

485 Finally, we investigated the relevance of the three key features of our model for this fit:
486 inhibition, feedback, and oscillations. We repeated the grid search fit but set either the inhibition
487 to zero, the feedback matrix equal for both /da/ and /ga/ (both 0.15), or the oscillation at an
488 amplitude of zero. Results showed that especially the oscillation and the differential feedback were
489 essential to reach a good fit (Figure 7B). Without the oscillation the model could not even fit better
490 than the mean of the model ($R^2 < 0$). Removing the inhibition had the least influence on the fit.

491 This suggest that all features (with a lesser extend the inhibition) are required to model the data
492 suggesting that oscillatory tracking is dependent on linguistic constrains flowing from the internal
493 language model.

494

495 Discussion

496 In the current paper, we combined an oscillatory model with a proxy for linguistic knowledge, an
497 internal language model, in order to investigate the model's processing capacity for onset timing
498 differences in natural speech. We show that word-to-word speech onset differences in natural
499 speech are indeed related to predictions flowing from the internal language model (estimated
500 through an RNN). Fixed oscillations aligned to the mean speech rate are robust against natural
501 temporal variations and even optimized for temporal variations that match the predictions flowing
502 from the internal model. Strikingly, when the pseudo-rhythmicity in speech matches the
503 predictions of the internal model, responses were more rhythmic for matched pseudo-rhythmic
504 compared to isochronous speech input. Our model is optimally sensitive to natural speech
505 variations, can explain phase dependent speech categorization behavior [31, 35, 44, 63], and
506 naturally comprises a neural phase code [40, 42, 43]. These results show that part of the pseudo-
507 rhythmicity of speech is expected by the brain and it is even optimized to process it in this manner,
508 but only when it follows the internal model.

509 Speech timing is variable and in order to understand how the brain tracks this pseudo-
510 rhythmic signal we need a better understanding of how this variability arises. Here, we isolated one
511 of the components explaining speech time variation, namely, constraints that are posed by an
512 internal language model. This goes beyond extracting the average speech rate [5, 19, 58], and might
513 be key to understanding how a predictive brain uses temporal cues. We show that speech timing
514 depends on the predictions made from an internal language model, even when those predictions
515 are highly reduced to be as simple as word predictability. While syllables generally follow a theta
516 rhythm, there is a systematic increase in syllabic rate as soon as more syllables are in a word. This
517 is likely a consequence of the higher close probability of syllables within a word which reduces the
518 onset differences of the later uttered syllables [59]. However, an oscillatory model constrained by
519 an internal language model is sensitive to these temporal variations, it is actually capable of
520 processing them optimally.

521 The oscillatory model we here pose has three components: oscillations, feedback, and
522 inhibition. The oscillations allow for the parsing of speech and provide windows in which
523 information is processed [3, 39, 64, 65]. Importantly, the oscillation acts as a temporal filter, such
524 that the activation time of any incoming signal will be confined to the high excitable window and
525 thereby is relatively robust against small temporal variations (Figure 5B). The feedback allows for

526 differential activation time dependent on the sensory input (Figure 5A). As a consequence, the
527 model is more sensitive to higher predictable speech input and therefore active earlier on the duty
528 cycle (this also means that oscillations are less robust against temporal variations when the feedback
529 is very strong). The inhibition allows for the network to be more sensitive to less predictable speech
530 units when they arrive later (the higher predictable nodes get inhibited at some point on the
531 oscillation; best illustrated by the simulation in Figure 7A). However, adding inhibition only
532 slightly improved the modeling fit (Figure 7B). In this way speech is ordered along the duty cycle
533 according to its predictability [43, 66]. The feedback in combination with an oscillatory model can
534 explain speech rate and phase dependent content effects. Moreover, it is an automatic temporal
535 code that can use time of activation as a cue for content [42]. The three components in the model
536 are common brain mechanisms [29, 42, 67-70] and follow many previously proposed organization
537 principles (e.g. temporal coding and parsing of information). While we implement these
538 components on an abstract level (not veridical to the exact parameters of neuronal interactions),
539 they illustrate how oscillations, feedback, and inhibition interact to optimize sensitivity to natural
540 pseudo-rhythmic speech.

541 The current model is not exhaustive and does not provide a complete explanation of all the
542 details of speech processing in the brain. For example, it is likely that the primary auditory cortex
543 is still mostly modulated by the acoustic pseudo-rhythmic input and only later brain areas follow
544 more closely the constraints posed by the language model of the brain. Therefore, more hierarchical
545 levels need to be added to the current model (but this is possible following equation (1)). Moreover,
546 the current model does not allow for phase or frequency shifts. This was intentional in order to
547 investigate how much a fixed oscillator could explain. We show that onset times matching the
548 predictions from the internal model can be explained by a fixed oscillator processing pseudo-
549 rhythmic input. However, when the internal model and the onset timings do not match the internal
550 model phase and/or frequency shift are still required and need to be incorporated (see e.g. [15, 19]).
551 Still, any coupling between brain oscillations and speech acoustics [19] needs to be extended with
552 the coupling of brain oscillations to brain activity patterns of internal models [71].

553 In the current paper we use an RNN to represent the internal model of the brain. However,
554 it is unlikely that the RNN captures the wide complexities of the language model in the brain. The
555 decades-long debates about the origin of a language model in the brain remains ongoing and
556 controversial. Utilizing the RNN as a proxy for our internal language model makes a tacit

557 assumption that language is fundamentally statistical or associative in nature, and does not posit the
 558 derivation or generation of knowledge of grammar from the input [72, 73]. In contrast, our brain
 559 could as well store knowledge of language that functions as fundamental interpretation principles
 560 to guide our understanding of language input [21, 24, 53, 65, 74]. Knowledge of language and
 561 linguistic structure could be acquired through an internal self-supervised comparison process
 562 extracted from environmental invariants and statistical regularities from the stimulus input [75-
 563 77]. Future research should investigate which language model can better account for the temporal
 564 variations found in speech.

565 A natural feature of our model is that time can act as a cue for content implemented as a
 566 phase code [43, 66]. This code unravels as an ordered list of predictability strength of the internal
 567 model. We predict that if speech nodes have a different base activity, ambiguous stimulus
 568 interpretation

569 should

570 dependent on
 571 the time/phase of
 572 presentation (see

573 [31, 63]). Indeed,
 574 we could model
 575 two temporal
 576 speech illusions
 577 (Figure 7). There

578 have also been
 579 null results

580 regarding the
 581 influence of
 582 phase on
 583 ambiguous

584 stimulus
 585 interpretation

586 [47, 78]. For the
 587 speech rate

588 effect, when

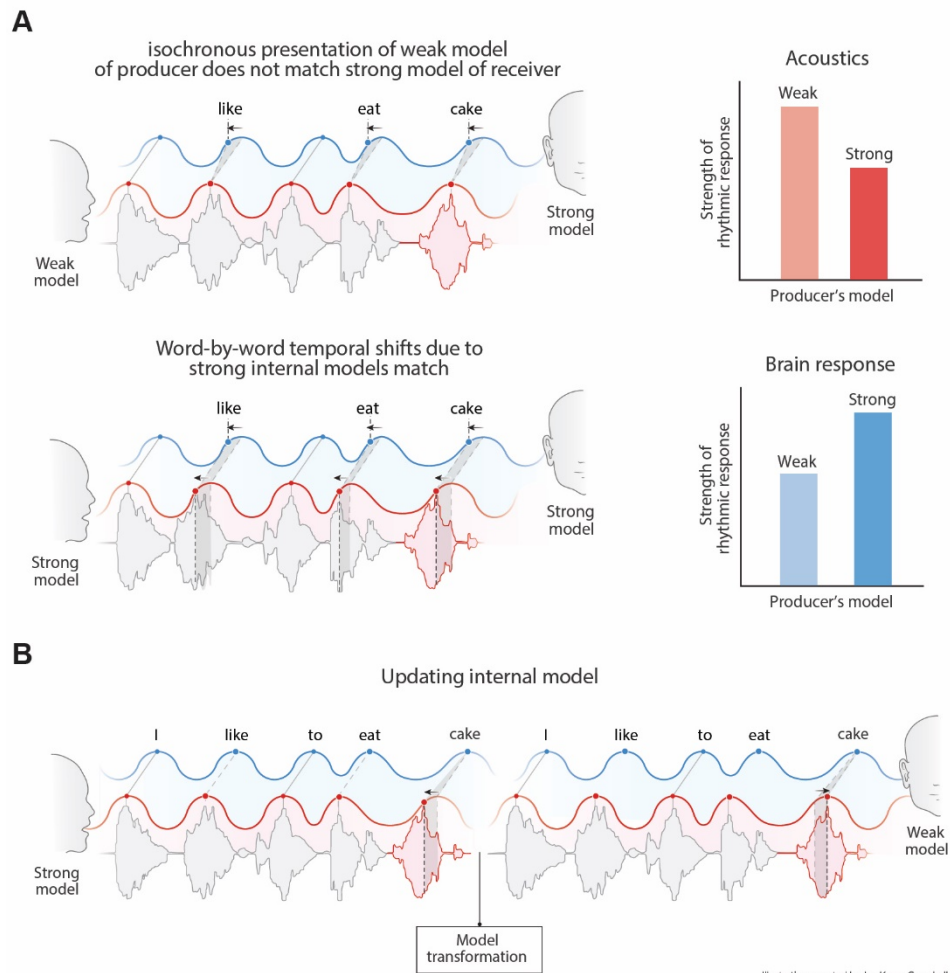


Figure 8. Predictions of the model. A) Acoustics signals will be more rhythmic when a producer has a weak versus a strong internal model (top right). When the producer's strong model matches the receiver's model the brain response will be more rhythmic for less rhythmic acoustic input. B) When a producer realizes the model of the receiver is weak it might transform its model and thereby their speech timing to match the receiver's expectations.

589 modifying the time of presentation with a neutral entrainer (summed sinusoidals with random
590 phase), no obvious phase effect was reported [47]. A second null result relates to a study where
591 participants were specifically instructed to maintain a specific perception in different blocks which
592 likely increases the pre-activation and thereby the phase [78]. Future studies need to investigate
593 the use of temporal/phase codes to disambiguate speech input and specifically use predictions in
594 their design.

595 The temporal dynamics of speech signals needs to be integrated with the temporal dynamics
596 of brain signals. However, it is unnecessary (and unlikely) that the exact duration of speech matches
597 with the exact duration of brain processes. Temporal expansion or squeezing of stimulus inputs
598 occur regularly in the brain [79, 80] and this temporal morphing also maps to duration [81-83] or
599 order illusions [84]. Our model predicts increased rhythmic responses for non-isochronous speech
600 matching the internal model. The perceived rhythmicity of speech could therefore also be an
601 illusion generated by a rhythmic brain signal somewhere in the brain.

602 When investigating the pseudo-rhythmicity in speech it is important to identify situations
603 where speech is actually more rhythmic. Two examples are the production of lists [85] and infant-
604 directed speech [86]. In both these examples it is clear that a strong internal predictive language
605 model is lacking either on the producer's or on the receiver's side, respectively. The infant-directed
606 speech also illustrates that a producer might proactively adapt its speech rhythm to the expectations

Table 3. Predictions from the current model

The more predictable a word, the earlier this word is uttered.
When there is a flat constraint distribution over an utterance (e.g., when probabilities are uniform over the utterance) the acoustics of speech should naturally be more rhythmic (Figure 8A).
If speech timing matches the internal language model, brain responses should be more rhythmic even if the acoustics are not (Figure 8A).
The more similar the internal language models of two speakers, the more effective they are in 'entraining' each other's brain.
If speakers suspect their listener to have a flatter constraint distribution than themselves (e.g., the environment is noisy, or the speakers are in a second language context), they adjust to the distribution by speaking more rhythmically (Figure 8B).
One adjusts the weight of the constraint distribution to a hierarchical level when needed. For example, when there is noise, participants adjust to the rhythm of primary auditory cortex instead of higher order language models. As a consequence, they speak more rhythmically.

The theoretical account provides various predictions that are listed in this table.

607 of the internal model of the receiver to align better with the predictions from the receiver’s model
608 (Figure 8B; similar to when you are speaking to somebody that is just learning a new language).
609 Other examples in which speech is more isochronous is during poems, during emotional
610 conversation [87], and in noisy situations [88]. While speculative, it is conceivable that in these
611 circumstances one puts more weight on a different level of hierarchy than the internal linguistic
612 model. In the case of poems and emotional conversation an emotional route might get more weight
613 in processing. In the case of noisy situations, stimulus input has to pass the first hierarchical level
614 of the primary auditory cortex which effectively gets more weight than the internal model.

615

616 **Conclusions**

617 We argued that pseudo-rhythmicity in speech is in part a consequence of top-down predictions
618 flowing from an internal model of language. This pseudo-rhythmicity is created by a speaker and
619 expected by a receiver if they have overlapping internal language models. Oscillatory tracking of
620 this signal does not need to be hampered by the pseudo-rhythmicity, but can use temporal
621 variations as a cue to extract content information since the phase of activation parametrically relates
622 to the likelihood of an input relative to the internal model. Brain responses can even be more
623 rhythmic to pseudo-rhythmic compared to isochronous speech if they follow the temporal delays
624 imposed by the internal model. This account provides various testable predictions which we list in
625 Table 3 and Figure 8. We believe that by integrating neuroscientific explanations of speech tracking
626 with linguistic models of language processing [21, 24], we can improve to explain temporal speech
627 dynamics. This will ultimately aid our understanding of language in the brain and provide a means
628 to improve temporal properties in speech synthesis.

629

630 **Acknowledgments**

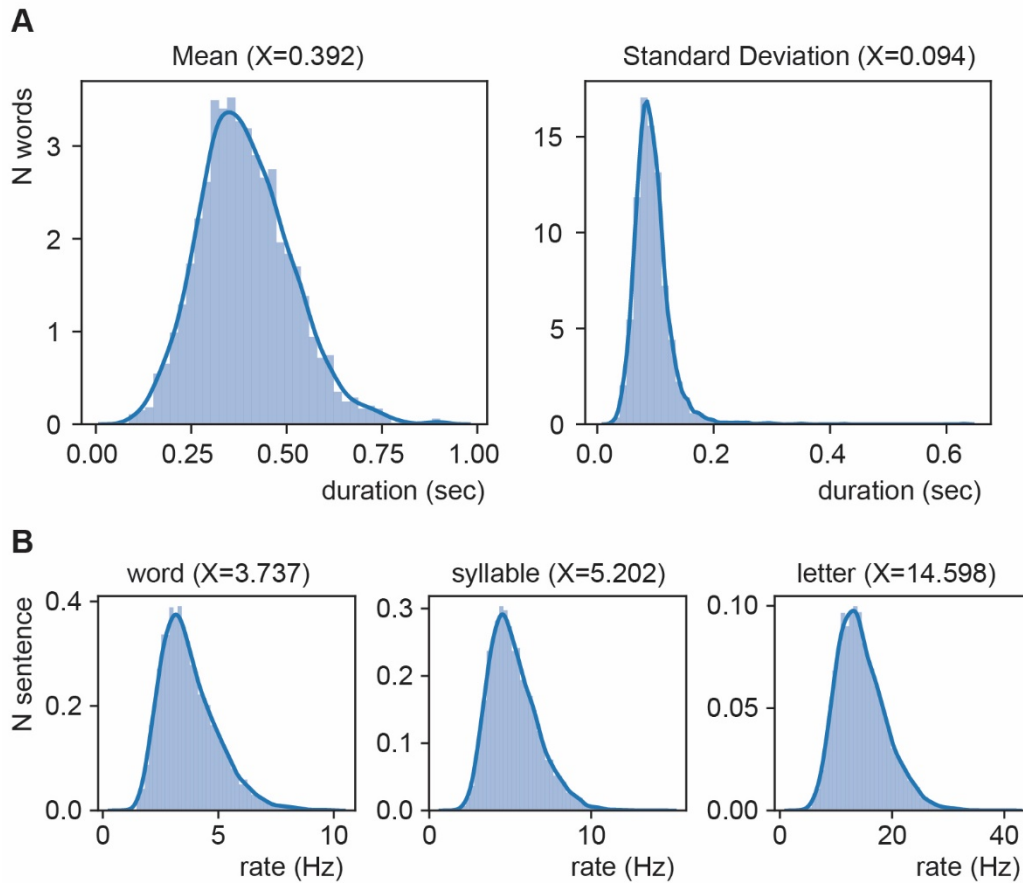
631 AEM was supported by the Lise Meitner Research Group “Language and Computation in Neural
632 Systems” from the Max Planck Society, and by the Netherlands Organization for Scientific Research
633 (grant 016.Vidi.188.029). Figure 1 and 8 were created in collaboration with scientific illustrator Jan-
634 Karen Campbell (www.jankaren.com).

635

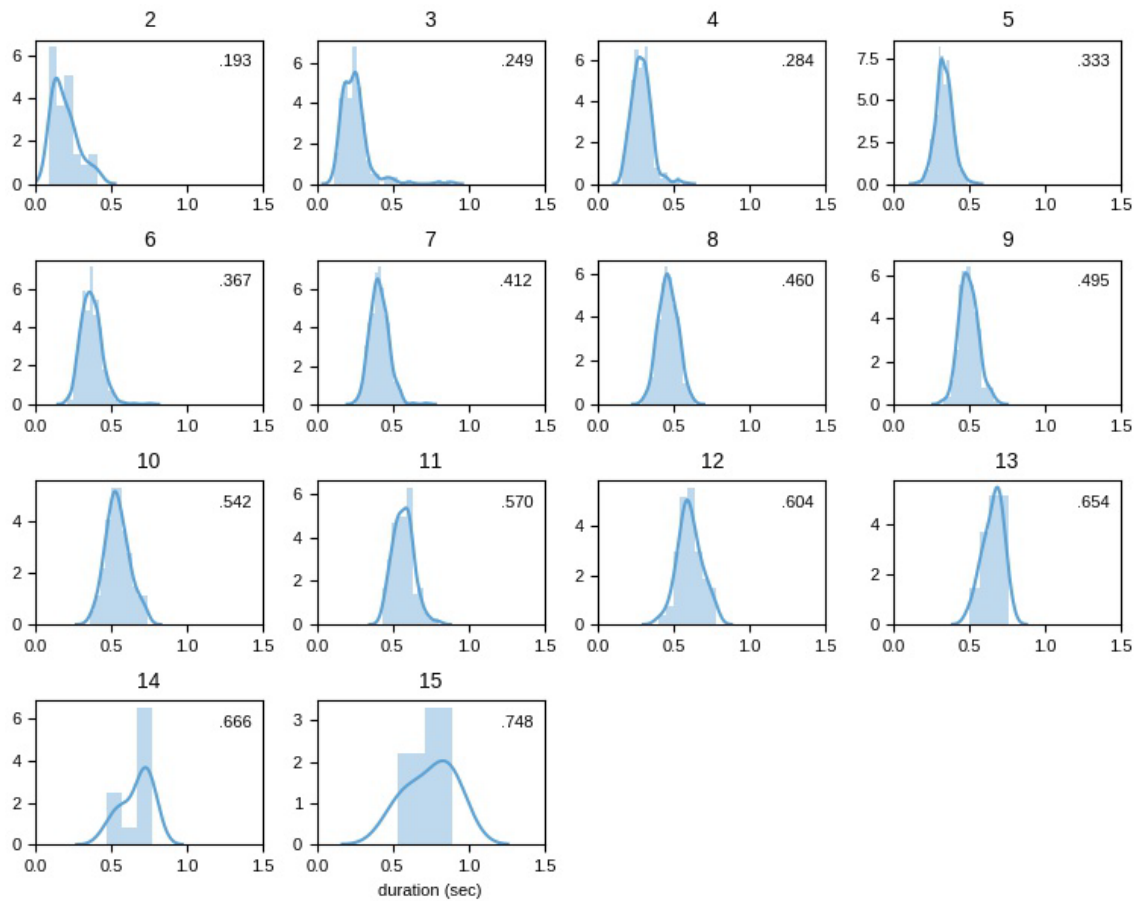
636 **Competing interests**

637 The authors declare no competing interests.

638 **Supporting Figures**



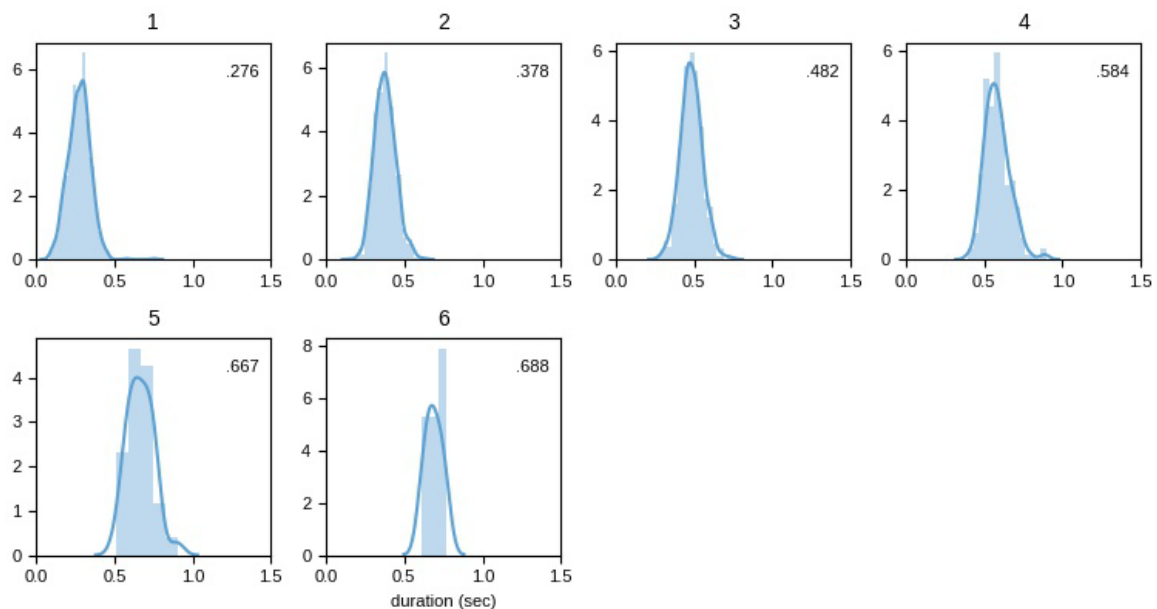
640 **Supporting figure 1.** Distribution of mean duration (A) and of average rate (B).



641

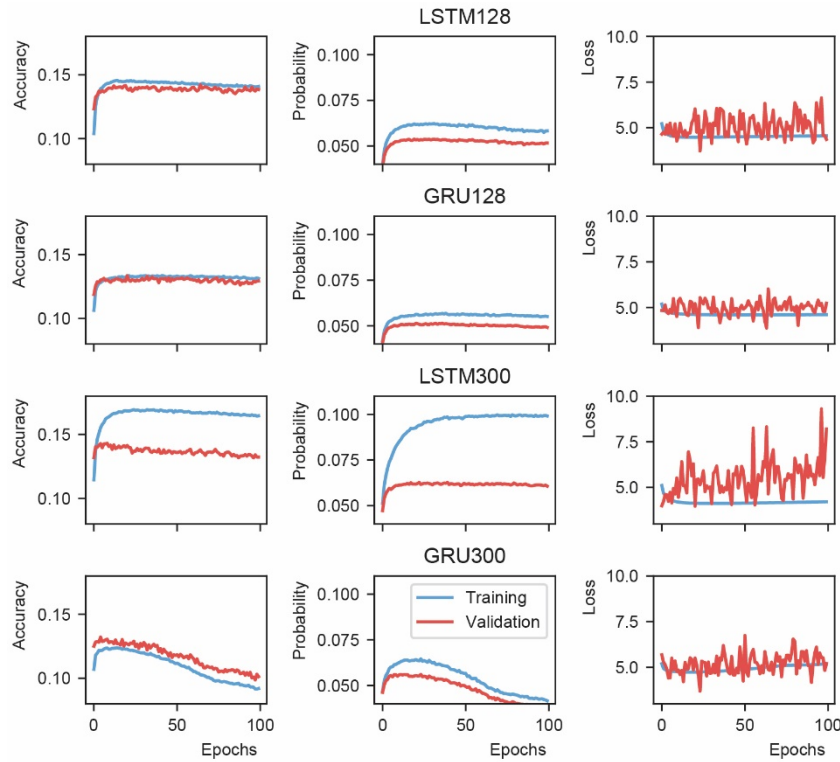
642 **Supporting figure 2.** Distribution of mean duration split up for word length (in characters).

643



644

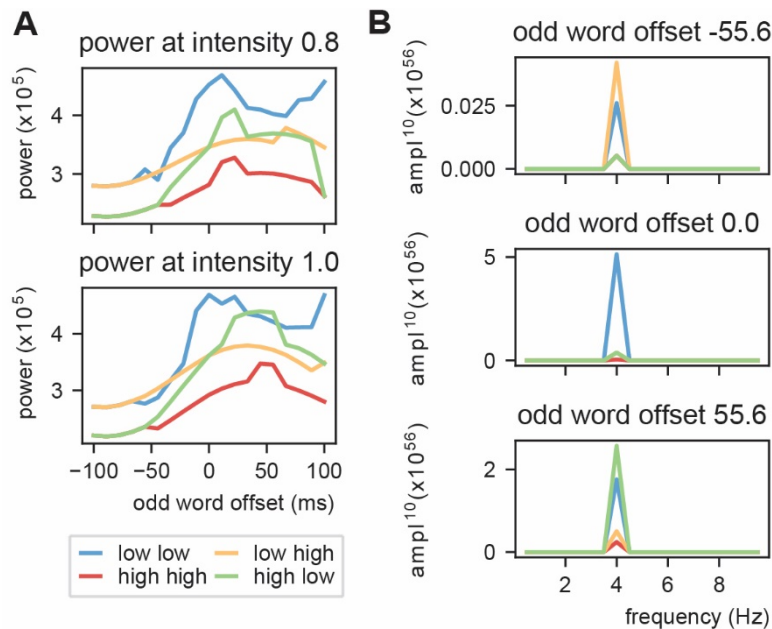
645 **Supporting figure 3.** Distribution of mean duration split up for syllable length.



646

647 **Supporting figure 4.** Recurrent neural network evaluation. Probability is defined as the mean of the
648 model output value at the node representing the next word.

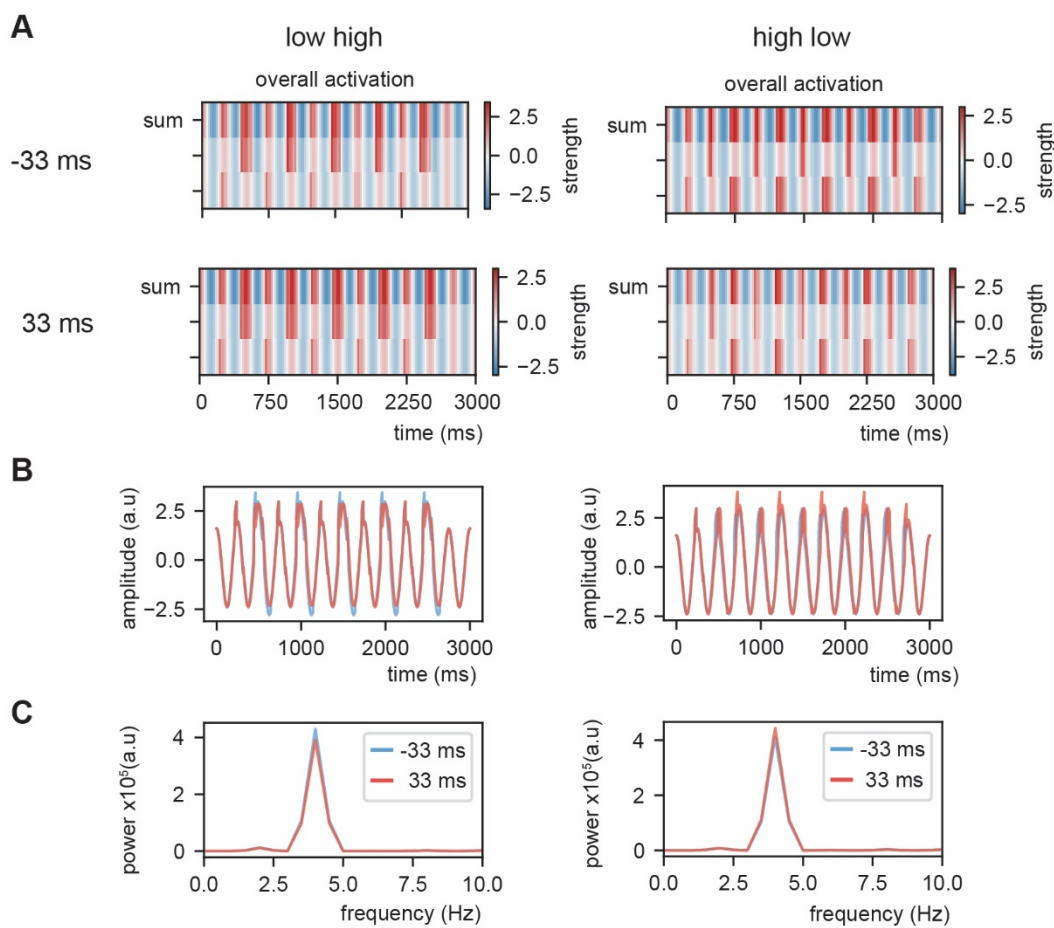
649



650

651 **Supporting figure 5.** Power at 4 Hz using linearly increasing sensory input. Conventions are the
652 same as in Figure 5D and E.

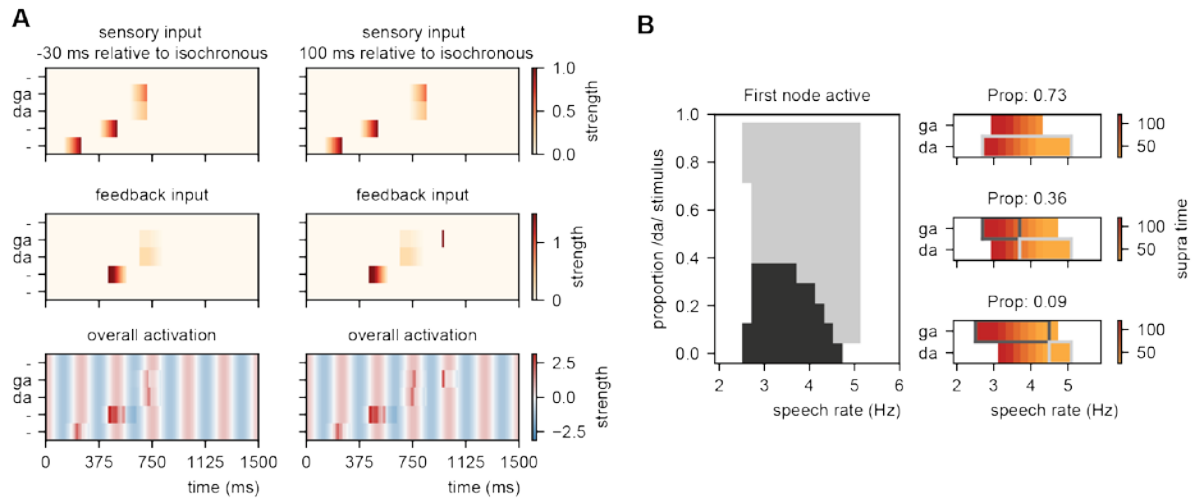
653



654

655 **Supporting figure 6.** Example of overall activation at threshold 0.8 (gaussian shaped input).

656



657

658 **Supporting figure 7.** Explaining speech timing illusions. A) Model activation of two example delays
659 for the fitting (figure 7A). B) Modulations due to ambiguous input at different speech rates.
660 Illustration of the node that is active first. Different proportions of the /da/ stimulus show activation
661 timing modulations at different speech rates. Conventions are the same as figure 7A.

662

663

664 References

- 665 1. Jones MR, Boltz M. Dynamic attending and responses to time. *Psychological Review*.
666 1989;96(3):459.
- 667 2. Large EW, Jones MR. The dynamics of attending: How people track time-varying events.
668 *Psychological Review*. 1999;106(1):119.
- 669 3. Giraud AL, Poeppel D. Cortical oscillations and speech processing: emerging computational
670 principles and operations. *Nat Neurosci*. 2012;15(4):511-7.
- 671 4. Ghitza O, Greenberg S. On the possible role of brain rhythms in speech perception:
672 intelligibility of time-compressed speech with periodic and aperiodic insertions of silence. *Phonetica*.
673 2009;66(1-2):113-26.
- 674 5. Ding N, Patel AD, Chen L, Butler H, Luo C, Poeppel D. Temporal modulations in speech and
675 music. *Neurosci Biobehav Rev*. 2017;81:181-7.
- 676 6. Arvaniti A. Rhythm, timing and the timing of rhythm. *Phonetica*. 2009;66(1-2):46-63.
- 677 7. Poeppel D. The analysis of speech in different temporal integration windows: cerebral
678 lateralization as 'asymmetric sampling in time'. *Speech Communication*. 2003;41(1):245-55.
- 679 8. Schroeder CE, Lakatos P. Low-frequency neuronal oscillations as instruments of sensory
680 selection. *Trends Neurosci*. 2009;32(1):9-18.
- 681 9. Luo H, Tian X, Song K, Zhou K, Poeppel D. Neural response phase tracks how listeners learn
682 new acoustic representations. *Curr Biol*. 2013;23(11):968-74.
- 683 10. Keitel A, Gross J, Kayser C. Perceptually relevant speech tracking in auditory and motor
684 cortex reflects distinct linguistic features. *PLoS Biol*. 2018;16(3):e2004473.
- 685 11. Lakatos P, Karmos G, Mehta AD, Ulbert I, Schroeder CE. Entrainment of neuronal oscillations
686 as a mechanism of attentional selection. *science*. 2008;320(5872):110-3.
- 687 12. Henry MJ, Obleser J. Frequency modulation entrains slow neural oscillations and optimizes
688 human listening behavior. *Proc Natl Acad Sci*. 2012;109(49):20095-100.
- 689 13. Herrmann B, Henry MJ, Grigutsch M, Obleser J. Oscillatory phase dynamics in neural
690 entrainment underpin illusory percepts of time. *J Neurosci*. 2013;33(40):15799-809.
- 691 14. Obleser J, Kayser C. Neural entrainment and attentional selection in the listening brain.
692 *Trends Cogn Sci*. 2019;23(11):913-26.
- 693 15. Rimmele JM, Morillon B, Poeppel D, Arnal LH. Proactive sensing of periodic and aperiodic
694 auditory patterns. *Trends Cogn Sci*. 2018;22(10):870-82.
- 695 16. Nolan F, Jeon H-S. Speech rhythm: a metaphor? *Philosophical Transactions of the Royal*
696 *Society B: Biological Sciences*. 2014;369(1658):20130396.
- 697 17. Jadoul Y, Ravignani A, Thompson B, Filippi P, de Boer B. Seeking temporal predictability in
698 speech: comparing statistical approaches on 18 world languages. *Front Hum Neurosci*. 2016;10:586.
- 699 18. Meyer L. The neural oscillations of speech processing and language comprehension: state of
700 the art and emerging mechanisms. *Eur J Neurosci*. 2018;48(7):2609-21.
- 701 19. Poeppel D, Assaneo MF. Speech rhythms and their neural foundations. *Nature Reviews*
702 *Neuroscience*. 2020:1-13.
- 703 20. Ten Oever S, Sack AT, Wheat KL, Bien N, Van Atteveldt N. Audio-visual onset differences are
704 used to determine syllable identity for ambiguous audio-visual stimulus pairs. *Frontiers in*
705 *Psychology*. 2013;4.
- 706 21. Martin AE. Language processing as cue integration: Grounding the psychology of language in
707 perception and neurophysiology. *Frontiers in psychology*. 2016;7:120.
- 708 22. Rosen S. Temporal information in speech: acoustic, auditory and linguistic aspects.
709 *Philosophical Transactions of the Royal Society of London Series B: Biological Sciences*.
710 1992;336(1278):367-73.
- 711 23. van de Ven V, Kochs S, Smulders F, De Weerd P. Learned interval time facilitates associate
712 memory retrieval. *Learn Memory*. 2017;24(4):158-61.
- 713 24. Martin AE. A compositional neural architecture for language. *J Cognit Neurosci*. 2020:1-20.

- 714 25. Marslen-Wilson WD. Functional parallelism in spoken word-recognition. *Cognition*.
715 1987;25(1-2):71-102.
- 716 26. Lau EF, Phillips C, Poeppel D. A cortical network for semantics:(de) constructing the N400.
717 *Nature Reviews Neuroscience*. 2008;9(12):920-33.
- 718 27. Nieuwland MS. Do 'early' brain responses reveal word form prediction during language
719 comprehension? A critical review. *Neurosci Biobehav Rev*. 2019;96:367-400.
- 720 28. O'Keefe J, Recce ML. Phase relationship between hippocampal place units and the EEG theta
721 rhythm. *Hippocampus*. 1993;3(3):317-30.
- 722 29. Malhotra S, Cross RW, van der Meer MA. Theta phase precession beyond the hippocampus.
723 *Reviews in the neurosciences*. 2012;23(1):39-65.
- 724 30. Bahramisharif A, Jensen O, Jacobs J, Lisman J. Serial representation of items during working
725 memory maintenance at letter-selective cortical sites. *PLoS Biol*. 2018;16(8):e2003805.
- 726 31. Ten Oever S, Sack AT. Oscillatory phase shapes syllable perception. *Proc Natl Acad Sci*.
727 2015;112(52):15833-7.
- 728 32. Kayser SJ, McNair SW, Kayser C. Prestimulus influences on auditory perception from sensory
729 representations and decision processes. *Proc Natl Acad Sci*. 2016;113(17):4842-7.
- 730 33. Di Liberto GM, O'Sullivan JA, Lalor EC. Low-frequency cortical entrainment to speech reflects
731 phoneme-level processing. *Curr Biol*. 2015;25(19):2457-65.
- 732 34. Ten Oever S, Hausfeld L, Correia J, Van Atteveldt N, Formisano E, Sack A. A 7T fMRI study
733 investigating the influence of oscillatory phase on syllable representations. *NeuroImage*. 2016;141:1-
734 9.
- 735 35. Thézé R, Giraud A-L, Mégevand P. The phase of cortical oscillations determines the
736 perceptual fate of visual cues in naturalistic audiovisual speech. *Science advances*.
737 2020;6(45):eabc6348.
- 738 36. Brennan JR, Martin AE. Phase synchronization varies systematically with linguistic structure
739 composition. *Philosophical Transactions of the Royal Society B*. 2020;375(1791):20190305.
- 740 37. Kaufeld G, Bosker HR, Alday PM, Meyer AS, Martin AE. Linguistic structure and meaning
741 organize neural oscillations into a content-specific hierarchy. *BioRxiv*. 2020.
- 742 38. Ghitza O. The theta-syllable: a unit of speech information defined by cortical function.
743 *Frontiers in psychology*. 2013;4:138.
- 744 39. Ghitza O. On the role of theta-driven syllabic parsing in decoding speech: intelligibility of
745 speech with a manipulated modulation spectrum. *Frontiers in Psychology*. 2012;3.
- 746 40. Panzeri S, Macke JH, Gross J, Kayser C. Neural population coding: combining insights from
747 microscopic and mass signals. *Trends Cogn Sci*. 2015;19(3):162-72.
- 748 41. Kayser C, Montemurro MA, Logothetis NK, Panzeri S. Spike-phase coding boosts and
749 stabilizes information carried by spatial and temporal spike patterns. *Neuron*. 2009;61(4):597-608.
- 750 42. Mehta M, Lee A, Wilson M. Role of experience and oscillations in transforming a rate code
751 into a temporal code. *Nature*. 2002;417(6890):741-6.
- 752 43. Lisman JE, Jensen O. The theta-gamma neural code. *Neuron*. 2013;77(6):1002-16.
- 753 44. Reinisch E, Sjerps MJ. The uptake of spectral and temporal cues in vowel perception is
754 rapidly influenced by context. *Journal of Phonetics*. 2013;41(2):101-16.
- 755 45. Kösem A, Bosker HR, Takashima A, Meyer A, Jensen O, Hagoort P. Neural entrainment
756 determines the words we hear. *Curr Biol*. 2018;28(18):2867-75. e3.
- 757 46. Bosker HR, Reinisch E, editors. Normalization for speechrate in native and nonnative speech.
758 18th International Congress of Phonetic Sciences (ICPhS 2015); 2015: International Phonetic
759 Association.
- 760 47. Bosker HR, Kösem A, editors. An entrained rhythm's frequency, not phase, influences
761 temporal sampling of speech. *Interspeech 2017*; 2017.
- 762 48. O'Malley S, Besner D. Reading aloud: Qualitative differences in the relation between
763 stimulus quality and word frequency as a function of context. *Journal of Experimental Psychology*:
764 *Learning, Memory, and Cognition*. 2008;34(6):1400.

- 765 49. Monsell S. The nature and locus of word frequency effects in reading. 1991.
766 50. Monsell S, Doyle MC, Haggard PN. Effects of frequency on visual word recognition tasks:
767 Where are they? *Journal of Experimental Psychology: General*. 1989;118(1):43.
768 51. Powers DM, editor Applications and explanations of Zipf's law. *New methods in language*
769 *processing and computational natural language learning*; 1998.
770 52. Piantadosi ST. Zipf's word frequency law in natural language: A critical review and future
771 directions. *Psychonomic bulletin & review*. 2014;21(5):1112-30.
772 53. Hagoort P. The core and beyond in the language-ready brain. *Neurosci Biobehav Rev*.
773 2017;81:194-204.
774 54. Beattie GW, Butterworth BL. Contextual probability and word frequency as determinants of
775 pauses and errors in spontaneous speech. *Language and speech*. 1979;22(3):201-11.
776 55. Gwilliams L, King J-R, Marantz A, Poeppel D. Neural dynamics of phoneme sequencing in real
777 speech jointly encode order and invariant content. *bioRxiv*. 2020.
778 56. Deacon D, Mehta A, Tinsley C, Nousak JM. Variation in the latencies and amplitudes of N400
779 and NA as a function of semantic priming. *Psychophysiology*. 1995;32(6):560-70.
780 57. Aubanel V, Schwartz J-L. The role of isochrony in speech perception in noise. *Scientific*
781 *reports*. 2020;10(1):1-12.
782 58. Pellegrino F, Coupé C, Marsico E. A cross-language perspective on speech information rate.
783 *Language*. 2011:539-58.
784 59. Thompson SP, Newport EL. Statistical learning of syntax: The role of transitional probability.
785 *Language learning and development*. 2007;3(1):1-42.
786 60. Guest O, Martin AE. How computational modeling can force theory building in psychological
787 science. *Perspectives on Psychological Science*. in press.
788 61. Meyer L, Sun Y, Martin AE. Synchronous, but not entrained: Exogenous and endogenous
789 cortical rhythms of speech and language processing. *Language, Cognition and Neuroscience*. 2019:1-
790 11.
791 62. Meyer L, Sun Y, Martin AE. "Entraining" to speech, generating language? *Language,*
792 *Cognition and Neuroscience*. in press.
793 63. Ten Oever S, Meierdierks T, Duecker F, De Graaf TA, Sack AT. Phase-coded oscillatory
794 ordering promotes the separation of closely matched representations to optimize perceptual
795 discrimination. *iScience*. 2020:101282.
796 64. Peelle JE, Davis MH. Neural oscillations carry speech rhythm through to comprehension.
797 *Frontiers in Psychology*. 2012;3.
798 65. Martin AE, Doumas LA. A mechanism for the cortical computation of hierarchical linguistic
799 structure. *PLoS Biol*. 2017;15(3):e2000663.
800 66. Jensen O, Bonnefond M, VanRullen R. An oscillatory mechanism for prioritizing salient
801 unattended stimuli. *Trends in cognitive sciences*. 2012;16(4):200-6.
802 67. Buzsáki G, Draguhn A. Neuronal oscillations in cortical networks. *science*.
803 2004;304(5679):1926-9.
804 68. Bastos AM, Usrey WM, Adams RA, Mangun GR, Fries P, Friston KJ. Canonical microcircuits for
805 predictive coding. *Neuron*. 2012;76(4):695-711.
806 69. Michalareas G, Vezoli J, Van Pelt S, Schoffelen J-M, Kennedy H, Fries P. Alpha-beta and
807 gamma rhythms subserve feedback and feedforward influences among human visual cortical areas.
808 *Neuron*. 2016;89(2):384-97.
809 70. Lisman JE. The theta/gamma discrete phase code occurring during the hippocampal phase
810 precession may be a more general brain coding scheme. *Hippocampus*. 2005;15(7):913-22.
811 71. Cumin D, Unsworth C. Generalising the Kuramoto model for the study of neuronal
812 synchronisation in the brain. *Physica D: Nonlinear Phenomena*. 2007;226(2):181-96.
813 72. Chater MHCN. *Connectionist psycholinguistics*: Greenwood Publishing Group; 2001.
814 73. McClelland JL, Elman JL. The TRACE model of speech perception. *Cognitive psychology*.
815 1986;18(1):1-86.

- 816 74. Friederici AD. The brain basis of language processing: from structure to function.
817 *Physiological reviews*. 2011;91(4):1357-92.
- 818 75. Martin AE, Doumas LA. Predicate learning in neural systems: using oscillations to discover
819 latent structure. *Current Opinion in Behavioral Sciences*. 2019;29:77-83.
- 820 76. Doumas LA, Hummel JE, Sandhofer CM. A theory of the discovery and predication of
821 relational concepts. *Psychological review*. 2008;115(1):1.
- 822 77. Doumas LA, Martin AE. Learning structured representations from experience. *Psychology of*
823 *Learning and Motivation*. 69: Elsevier; 2018. p. 165-203.
- 824 78. Kösem A, Basirat A, Azizi L, van Wassenhove V. High-frequency neural activity predicts word
825 parsing in ambiguous speech streams. *J Neurophysiol*. 2016;116(6):2497-512.
- 826 79. Eagleman DM, Peter UT, Buonomano D, Janssen P, Nobre AC, Holcombe AO. Time and the
827 brain: how subjective time relates to neural time. *J Neurosci*. 2005;25(45):10369-71.
- 828 80. Pariyadath V, Eagleman D. The effect of predictability on subjective duration. *PloS one*.
829 2007;2(11):e1264.
- 830 81. Eagleman DM. Human time perception and its illusions. *Curr Opin Neurobiol*.
831 2008;18(2):131-6.
- 832 82. Terao M, Watanabe J, Yagi A, Nishida Sy. Reduction of stimulus visibility compresses
833 apparent time intervals. *Nat Neurosci*. 2008;11(5):541-2.
- 834 83. Ulrich R, Nitschke J, Rammsayer T. Perceived duration of expected and unexpected stimuli.
835 *Psychological research*. 2006;70(2):77-87.
- 836 84. Vroomen J, Keetels M. Perception of intersensory synchrony: A tutorial review. *Attention,*
837 *Perception, & Psychophysics*. 2010;72(4):871-84.
- 838 85. Jefferson G. List construction as a task and resource. *Interaction competence*. 1990;63:92.
- 839 86. Fernald A. Speech to infants as hyperspeech: Knowledge-driven processes in early word
840 recognition. *Phonetica*. 2000;57(2-4):242-54.
- 841 87. Hawkins S. Situational influences on rhythmicity in speech, music, and their interaction.
842 *Philosophical Transactions of the Royal Society B: Biological Sciences*. 2014;369(1658):20130398.
- 843 88. Bosker HR, Cooke M. Talkers produce more pronounced amplitude modulations when
844 speaking in noise. *The Journal of the Acoustical Society of America*. 2018;143(2):EL121-EL6.

845