

# Speech understanding oppositely affects acoustic and linguistic neural tracking in a speech rate manipulation paradigm.

Abbreviated title: Untangling acoustic and linguistic neural tracking

**Eline Verschueren**<sup>a,\*, $\Delta$</sup>  & **Marlies Gillis**<sup>a, $\Delta$</sup> , **Lien Decruy**<sup>b</sup>, **Jonas**

**Vanthornhout**<sup>a</sup> and **Tom Francart**<sup>a</sup>

<sup>a</sup>: *Research Group Experimental Oto-rhino-laryngology (ExpORL), Department of Neurosciences, KU Leuven - University of Leuven, 3000 Leuven, Belgium*

<sup>b</sup>: *Institute for Systems Research, University of Maryland, College Park, MD 20742, USA*

$\Delta$ : *shared first authorship*

Correspondence\*:

Eline Verschueren

Herestraat 49, bus 721, 3000 Leuven, Belgium

[eline.verschueren@kuleuven.be](mailto:eline.verschueren@kuleuven.be)

1 **ABSTRACT**

2 When listening to continuous speech, the human brain can track features of the  
3 presented speech signal. It has been shown that neural tracking of acoustic features  
4 is a prerequisite for speech understanding and can predict speech understanding  
5 in controlled circumstances. However, the brain also tracks linguistic features  
6 of speech, which may be more directly related to speech understanding. We  
7 investigated acoustic and linguistic speech processing as a function of varying  
8 speech understanding by manipulating the speech rate. In this paradigm, acoustic  
9 and linguistic speech processing are affected simultaneously but in opposite  
10 directions: When the speech rate increases, more acoustic information per second  
11 is present. In contrast, linguistic information decreases as speech becomes less  
12 intelligible at higher speech rates. We measured the EEG of 18 participants who  
13 listened to speech at various speech rates. As expected and confirmed by the  
14 behavioral results, speech understanding decreased with increasing speech rate.  
15 Accordingly, linguistic neural tracking decreased with increasing speech rate, but  
16 acoustic neural tracking increased. This indicates that neural tracking of linguistic  
17 representations can capture the gradual effect of decreasing speech understanding.  
18 In addition, increased acoustic neural tracking does not necessarily imply better  
19 speech understanding. This suggests that, although more challenging to measure  
20 due to the low signal-to-noise ratio, linguistic neural tracking may be a more direct  
21 predictor of speech understanding.

22 **Keywords:** neural coding, natural speech, speech rate, EEG, acoustic hearing, linguistic  
23 representations

24

25 **Significance statement:** An increasingly popular method to investigate neural speech  
26 processing is to measure neural speech tracking. Although much research has been done  
27 on how the brain tracks acoustic speech features, linguistic speech features have received

28 less attention. In this study, we disentangled acoustic and linguistic characteristics of neural  
29 speech tracking via manipulating the speech rate. A proper way of objectively measuring  
30 auditory and language processing paves the way towards clinical applications: An objective  
31 measure of speech understanding would allow for behavioral-free evaluation of speech  
32 understanding, which allows to evaluate hearing loss and adjust hearing aids based on  
33 brain responses. This objective measure would benefit populations from whom obtaining  
34 behavioral measures may be complex, such as young children or people with cognitive  
35 impairments.

## 1 INTRODUCTION

36 Understanding speech relies on the integration of different acoustic and linguistic properties  
37 of the speech signal. The acoustic properties are mainly related to sound perception,  
38 while the linguistic properties are linked to the content and understanding of speech. When  
39 listening to continuous speech, our brain can track both the acoustic and linguistic properties  
40 of the presented speech signal.

41 Neural tracking of acoustic properties of natural speech has been the subject of many  
42 studies. Particular emphasis has been placed on recovering the temporal envelope, i.e., the  
43 slow modulations of the speech signal, from the brain responses, so-called neural envelope  
44 tracking. The temporal envelope is essential for speech understanding (Shannon et al.,  
45 1995), and neural envelope tracking can be linked to speech intelligibility (e.g. Ding and  
46 Simon, 2013; Vanthornhout et al., 2018; Lesenfants et al., 2019; Iotzov and Parra, 2019;  
47 Verschueren et al., 2020). However, only taking acoustic speech properties into account to  
48 investigate neural speech tracking would underestimate the complexity of the human brain,  
49 where linguistic properties also play an essential part, as reviewed in detail by Brodbeck  
50 and Simon (2020).

51 In addition to acoustic properties, there is growing interest in retrieving linguistic  
52 properties from brain responses to speech. Broderick et al. (2018) used semantic  
53 dissimilarity to quantify the meaning carried by words based on their preceding context.  
54 They report that the brain responds in a time-locked way to the semantic context of each  
55 content word. Additionally, neural tracking is also observed to linguistic properties derived  
56 from the probability of a given word or phoneme, i.e., word or phoneme surprisal (Brodbeck  
57 et al., 2018; Weissbart et al., 2019; Koskinen et al., 2020). Recently Gillis et al. (2021b)  
58 combined several linguistic neural tracking measures and evaluated the potential of each  
59 measure as a neural marker of speech intelligibility. After controlling for acoustic properties,  
60 phoneme surprisal, cohort entropy, word surprisal, and word frequency were significantly  
61 tracked. These results show the potential of linguistic representations as a neural marker of

62 speech intelligibility. In addition, this underlines the importance of controlling for acoustic  
63 features when investigating linguistic neural processing, as acoustic and linguistic features  
64 are often correlated (Brodbeck and Simon, 2020).

65 We investigated whether neural speech processing can capture the effect of gradually  
66 decreasing speech understanding by manipulating the speech rate. In this study, we focused  
67 on acoustic and linguistic speech processing. By changing the speech rate, we manipulate  
68 acoustic and linguistic speech processing simultaneously but in opposite directions: When  
69 increasing the speech rate, more phonemes, words, and sentences, and thus more acoustic  
70 information per second is present. In contrast, linguistic information decreases because it  
71 becomes more challenging to identify the individual phonemes or words at high speech rates,  
72 causing decreased speech understanding. We hypothesize that neural tracking of acoustic  
73 features will increase with increasing speech rate because more acoustic information will  
74 be present. However, linguistic speech tracking will decrease with increasing speech rate  
75 because of decreasing speech understanding. The effect of speech rate on neural responses  
76 to speech has already been investigated. However, all these studies only investigated brain  
77 responses to the acoustic properties of the speech signal (Ahissar et al., 2001; Nourski  
78 et al., 2009; Hertrich et al., 2012; Müller et al., 2019; Casas et al., 2021). No study, to  
79 our knowledge, reported on how speech rate affects linguistic speech processing and the  
80 potential interaction between both. In addition, no consensus has been reached on the effect  
81 of speech rate on acoustic neural tracking. For example, Nourski et al. (2009) reported  
82 that phase-locked responses decrease with increasing speech rate, similar to Ahissar et al.  
83 (2001) and Hertrich et al. (2012). However, in the same data, Nourski et al. (2009) also  
84 reported that time-locked responses to the envelope (70-250 Hz) could still be found at very  
85 high speech rates where speech is no longer understood.

86 We investigated how linguistic and acoustic speech tracking are affected when  
87 speech understanding gradually decreases. Analyzing neural speech tracking to different

88 characteristics of the presented speech allows us to identify neural patterns associated with  
89 speech understanding.

## 90 **2 MATERIAL AND METHODS**

### 91 **2.1 Participants**

91 Eighteen participants aged between 19 and 24 years (4 men and 14 women) took part in the  
92 experiment after having provided informed consent. Participants had Dutch as their mother  
93 tongue and were all normal-hearing, confirmed with pure tone audiometry (thresholds  $\leq$   
94 25 dB HL at all octave frequencies from 125 Hz to 8 kHz). The study was approved by the  
95 Medical Ethics Committee UZ Leuven / Research (KU Leuven) with reference S57102.

### 96 **2.2 Speech material**

97 The story presented during the EEG measurement was ‘A casual vacancy’ by J.K. Rowling,  
98 narrated in Dutch by Nelleke Noordervliet. The story was manually cut into 12 blocks  
99 of varying length randomly selected from the following list: 4 min, 5 min, 8.5 min, 12.5  
100 min, 18 min, and 23 min. After cutting the story, the story was time-compressed with  
101 the Pitch Synchronous Overlap and Add algorithm (PSOLA) from PRAAT (Boursma  
102 and Weenink, 2018) to manipulate the speech rate. Six different compression ratio’s (CR)  
103 were used: 1.4, 1.0, 0.6, 0.4, 0.28, 0.22 with corresponding speech rates varying from  
104  $\approx 2.6$  syllables/second (CR=1.4) to  $\approx 16.2$  syllables/second (CR=0.22). The fastest CR  
105 (CR=0.22) was applied to the longest part (23 min), the one but fastest CR (CR=0.28) to  
106 the one but longest part (18 min), and so on. This way, all story parts were compressed or  
107 expanded to  $\approx 5$  minutes. These blocks had slightly different lengths because word and  
108 sentence boundaries were taken into account while cutting the story, which is important for  
109 the linguistic analysis. Every speech rate was presented twice to obtain 10 minutes of speech  
110 at the same rate. The story was presented in chronological order. For each stimulus block,  
111 we determined the number of syllables using the forced aligner of the speech alignment

112 component of the reading tutor (Duchateau et al., 2009) and CELEX database (Baayen  
113 et al., 1996). The number of syllables uttered for each speech block was then divided by  
114 the duration of the speech block in seconds to obtain the speech rate.

115 After each part of the story, content questions were asked to maximize the participants'  
116 attention and motivation. In addition, speech intelligibility was measured after each block  
117 by asking the participants to rate their speech understanding on a scale from 0 to 100%  
118 following the question 'Which percentage of the story did you understand?'. A short  
119 summary of the story was shown in the beginning of the experiment to enhance intelligibility  
120 as some participants started with more difficult speech rates.

## 121 **2.3 Experimental setup**

### 122 2.3.1 EEG recording

123 EEG was recorded with a 64-channel BioSemi ActiveTwo EEG recording system at a  
124 sample rate of 8192 Hz. Participants sat in a comfortable chair and were asked to move as  
125 little as possible during the EEG recordings. All stimuli were presented bilaterally using  
126 APEX 4 (Francart et al., 2008), an RME Multiface II sound card (Haimhausen, Germany),  
127 and Etymotic ER-3A insert phones (Illinois, USA). The setup was calibrated using a 2 cm<sup>3</sup>  
128 coupler of the artificial ear (Brüel & Kjør 4152, Denmark). Recordings were made in a  
129 soundproof and electromagnetically shielded room.

## 130 **2.4 Signal processing**

### 131 2.4.1 EEG processing

132 We processed the EEG in 5 consecutive steps. Firstly, we drift-corrected the EEG signals  
133 by applying a first-order highpass Butterworth filter with a cutoff frequency of 0.5 Hz  
134 in the forward and backward direction. Then, we reduced the sampling frequency of the

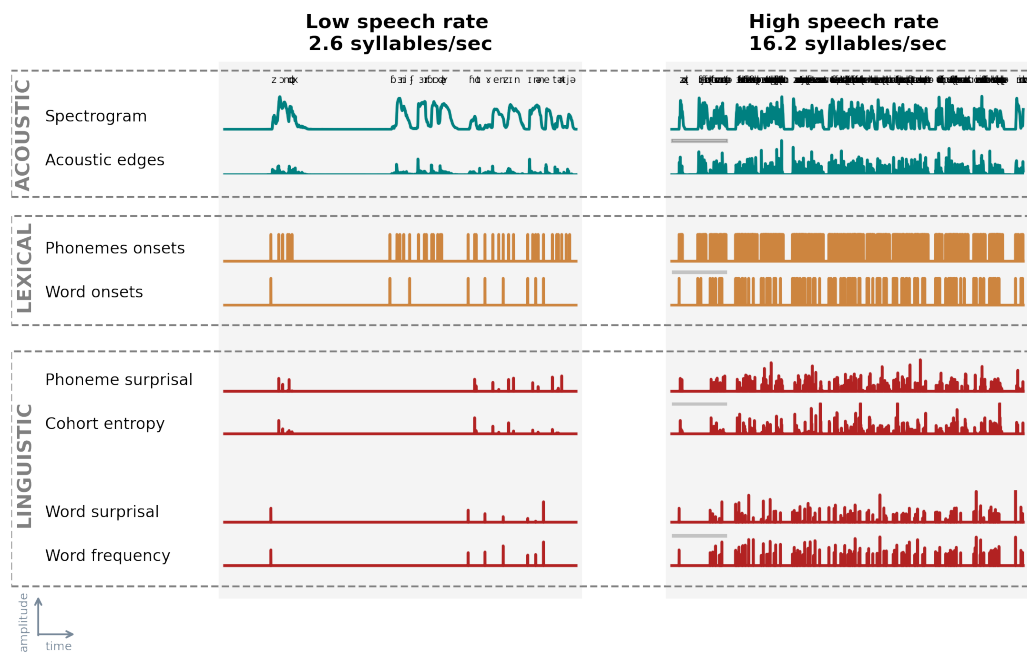
135 EEG from 8192 Hz to 256 Hz to reduce computation time. Artifacts related to eyeblinks  
136 were removed with a multichannel Wiener filter (Somers et al., 2018). Subsequently, we  
137 referenced the EEG signals to the common average signal. Lastly, we removed the power  
138 line frequency of 50 Hz by using a second-order IRR notch filter at this frequency with a  
139 quality factor of 35 to determine the filter's bandwidth at the -3dB.

#### 140 2.4.2 Stimuli Representations

141 This study aims to investigate acoustic and linguistic neural tracking at different speech  
142 rates. To examine acoustic tracking, we estimated neural tracking based solely on acoustic  
143 representations of the stimulus, namely the spectrogram and acoustic edges. To investigate  
144 linguistic neural tracking, we created two models: a model to control speech acoustics,  
145 which consisted of acoustic and lexical segmentation representations, and a model that  
146 included linguistic representations on top of these acoustic and lexical segmentation  
147 representations. All speech representations used in the analysis are visualized in Figure 1.

148 The spectrogram representations were calculated based on the low-pass filtered speech  
149 stimulus (zero-phase low-pass FIR filter with a hamming window of 159 samples). We  
150 low-pass filtered the stimulus at a cut-off frequency of 4 kHz as the insert earphones also  
151 low-pass filter at this frequency. Subsequently, we calculated the spectrogram representation  
152 from this filtered stimulus using the Gammatone Filterbank Toolkit (Heeris, 2014, center  
153 frequencies between 70 and 4000 Hz with 256 filter channels and an integration window of  
154 0.01 second). By using a Gammatone Filterbank, the estimated filter outputs are closer to  
155 the human auditory response (Slaney, 1998). We combined the filter outputs by averaging  
156 them into eight frequency bands with center frequencies of 124 Hz, 262 Hz, 455 Hz, 723 Hz,  
157 1098 Hz, 1618 Hz, 2343 Hz, and 3352 Hz. To calculate the acoustic edges representations,  
158 we took the derivative of the spectrogram's response in each frequency band and mapped all  
159 its negative values to 0. Lastly, we reduced the sampling frequency of these representations  
160 to the same sampling frequency as the EEG, namely 256 Hz.





**Figure 1. Speech representations at acoustic, lexical and linguistic level** We visualized the speech representations used in this study for all three levels: acoustic (averaged across frequency bands; top box; blue), lexical (middle box; orange), and linguistic (bottom box; red) for the lowest speech rate (SR = 2.6 syllables/sec, left) and highest speech rate (SR = 16.2 syllables/sec; right) for the first 10 seconds of the speech material.

161 To determine linguistic neural tracking, we carefully controlled for neural responses  
162 related to acoustic and lexical characteristics of the speech. As pointed out by Brodbeck and  
163 Simon (2020); Gillis et al. (2021b), it is important to control for these characteristics  
164 when investigating linguistic neural tracking, as otherwise spurious linguistic neural  
165 tracking can be observed due to the high correlation between linguistic and acoustic  
166 representations. To evaluate linguistic neural tracking, we determined the added value of  
167 linguistic representations by subtracting the performance of the model containing acoustic  
168 and lexical segmentation characteristics of the speech from the performance of the model  
169 that included the same representations together with the linguistic representations. We used  
170 four linguistic representations: phoneme surprisal, cohort entropy, word surprisal, and word

171 frequency, which according to Gillis et al. (2021b), have an added value over and beyond  
172 acoustic representations.

173 All linguistic representations are one-dimensional arrays with impulses at the onsets  
174 of phonemes or words. The amplitude of an impulse represents the amount of linguistic  
175 information conveyed by the phoneme or word. To obtain the timing of the phonemes  
176 and words, we used the forced aligner of the speech alignment component of the reading  
177 tutor (Duchateau et al., 2009). Similar to linguistic representations, lexical segmentation  
178 representations are one-dimensional arrays. However, the impulses' amplitudes are one and  
179 thus independent of the amount of linguistic information.

180 Phoneme surprisal and cohort entropy are two linguistic representations that describe the  
181 linguistic content of a phoneme. Phoneme surprisal is thought to be a measure of phoneme  
182 prediction error as it represents how surprising a phoneme is given the previously uttered  
183 phonemes. It is calculated as the negative logarithm of the inverse conditional probability of  
184 the phoneme given the preceding phonemes of the word. Another linguistic representation  
185 at the phoneme level is cohort entropy derived from the cohort of words congruent with  
186 the already uttered phonemes. More specifically, it is calculated as the Shannon entropy of  
187 this active cohort of words, reflecting the degree of competition between them. To calculate  
188 both representations, we used a custom pronunciation dictionary that maps a word to its  
189 phoneme representation. This dictionary was created by manual and grapheme-to-phoneme  
190 conversion and contained the segmentation of 9157 words. The word probabilities were  
191 derived from the SUBTLEX-NL database (Keuleers et al., 2010). The linguistic information  
192 of the initial phoneme was not modeled in these representations. More details regarding  
193 phoneme surprisal and cohort entropy, as well as the mathematical determinations, can be  
194 found in Brodbeck et al. (2018).

195 The linguistic information conveyed by a word is described by word surprisal and word  
196 frequency. Similar to phoneme surprisal, word surprisal is thought to model a word's  
197 prediction error. It reflects how surprising a word is given its preceding words. We used

198 a 5-gram model to determine the negative logarithm of the conditional probability of the  
199 word given the preceding words. Therefore, a word's surprisal is estimated given its four  
200 preceding words. Word frequency was derived from the same 5-gram model but without  
201 including previous words, describing the word's unigram probability.

### 202 2.4.3 Determination of Neural Tracking

203 To determine neural tracking, we used a forward modeling approach, estimating how the  
204 brain responds to specific speech characteristics. The temporal response function (TRF)  
205 describes the relationship between the presented stimulus and measured EEG. It also  
206 allows us to predict the EEG responses associated with the speech stimulus. By correlating  
207 the predicted EEG responses with the measured EEG responses, we obtain a prediction  
208 accuracy per EEG channel. This prediction accuracy is a measure of neural tracking.

209 We used the boosting algorithm (David et al., 2007) implemented by the Eelbrain Toolbox  
210 (Brodbeck, 2020) to estimate the TRF and obtain the prediction accuracy. We used an  
211 integration window of -100 to 600 ms, i.e., the neural response is estimated ranging from  
212 100 ms before activation of the stimulus characteristic to 600 ms after its activation. We  
213 use a broad integration window to ensure that the model captures the brain responses to the  
214 linguistic representations, which occur at longer latencies. As each speech rate condition  
215 was presented twice, we estimated the TRF on the concatenation of these two blocks per  
216 speech rate, i.e., ten minutes of data. Before the TRF estimation, the data is normalized by  
217 dividing by the Euclidean norm per channel. We applied this normalization for the stimulus  
218 and EEG data individually. Then the boosting algorithm estimates the associated response  
219 behavior using a fixed step size of 0.005. We derived the TRF and prediction accuracy per  
220 channel using a cross-validation scheme: the TRF was estimated and validated on partitions  
221 unseen during testing of the TRF to obtain the prediction accuracy. More specifically, we  
222 used 10-fold cross-validation, implying the data was split into ten equally long folds, of  
223 which eight folds are used for estimating the TRF, one fold for validation, and one fold for

224 testing. The obtained TRFs and prediction accuracies are then averaged across the different  
225 folds. Note that the validation fold is required to determine the stopping criterion: we used  
226 an early stopping based on the  $\ell_2$ -norm, i.e., estimation of the TRF is stopped when the  
227 Euclidian distance between the actual and predicted EEG data on the validation partition  
228 stops decreasing. The resulting TRFs are sparse. Therefore, to account for the inter-subject  
229 variability and obtain a meaningful average TRF response across subjects, we smoothed  
230 the TRFs across time by convolving the estimated response with a hamming window of 50  
231 ms in the time dimension.

232 To determine the acoustic tracking of the speech, we purely used acoustic representations.  
233 Therefore, we determined the prediction accuracy and TRFs based on the spectrogram  
234 and acoustic edges. Regarding the linguistic tracking of speech, we investigated the added  
235 value of these linguistic representations. To determine the added value, we subtracted  
236 the prediction accuracies of two different models. Firstly, we estimated a baseline model  
237 consisting of acoustic and lexical segmentation representations. Secondly, we estimated  
238 a combined model which included linguistic representations on top of the acoustic and  
239 lexical segmentation representations. By subtracting the prediction accuracy obtained with  
240 the baseline model from the prediction accuracy of the combined model, we can examine  
241 the added value of the linguistic representations after controlling for the acoustic and lexical  
242 segmentation representations.

243 We used two predetermined channel selections to investigate the effect of acoustic and  
244 linguistic tracking. The neural responses to acoustics are significantly different from those  
245 to linguistic content and therefore require a different channel selection. We used a frontal  
246 channel selection for acoustic neural tracking based on Lesenfants et al. (2019) and a central  
247 channel selection for linguistic neural tracking as reported by Gillis et al. (2021b).

248 These channel selections were used to visualize the TRFs and to determine associated peak  
249 latency and amplitudes. To determine the peak characteristics, we set a preset time window  
250 based on the TRF averaged across subjects (see Table 1). Within this time window, we

**Table 1.** Time windows selected per speech representation to determine the peak characteristics

Speech representation	Time window(s)	Channel selection
Spectrogram Acoustic edges	0 to 90 ms, 110 to 200 ms	frontocentral
Phoneme surprisal Cohort entropy	200 to 300 ms	central
Word surprisal Word frequency	350 to 450 ms	central

251 normalized the TRF per channel by dividing the TRF by its  $\ell_2$ -norm over time to decrease  
 252 across subject variability and averaged the TRF across the channel selection. Depending  
 253 on a positive or negative peak, we determined the maximal or minimal amplitude and its  
 254 corresponding latency to obtain the peak amplitude and latency. If the peak latency was the  
 255 same as the beginning of the window, indicating the end of the previous peak, we discarded  
 256 the peak from the analysis (see Table 2).

**Table 2.** Number of peaks detected per speech representation per speech rate with  $n_{\max} = 18$  (= amount of participants).

	2.6 syll/sec	3.6 syll/sec	6.2 syll/sec	9.0 syll/sec	12.9 syll/sec	16.2 syll/sec
Spectrogram - peak 1	n = 15	n = 17	n = 18	n = 18	n = 18	n = 18
Spectrogram - peak 2	n = 14	n = 16	n = 18	n = 15	n = 13	n = 11
Acoustic edges - peak 1	n = 17	n = 17	n = 18	n = 18	n = 18	n = 18
Acoustic edges - peak 2	n = 18	n = 17	n = 16	n = 10	n = 8	n = 8
Phoneme surprisal	n = 17	n = 18	n = 17	n = 16	n = 15	n = 18
Cohort entropy	n = 17	n = 16	n = 16	n = 15	n = 18	n = 17
Word surprisal	n = 15	n = 16	n = 15	n = 18	n = 16	n = 17
Word frequency	n = 16	n = 15	n = 13	n = 14	n = 15	n = 17

## 257 **2.5 Statistics**

258 Statistical analysis was performed using MATLAB (version R2018a) and R (version 3.4.4)  
259 software. The significance level was set at  $\alpha=0.05$  unless otherwise stated.

260 To evaluate the subjectively rated speech understanding results we calculated the  
261 correlation between speech rate and rated speech understanding using a Spearman rank  
262 correlation. In addition, we fitted a sigmoid function on the data to address the relation  
263 between rated speech understanding and speech rate using the `minpack.lm` package (Elzhov  
264 et al., 2016) in R. For further statistical analysis, we selected speech rate (and not  
265 subjectively rated speech understanding) as a main predictor. We opted for this because a  
266 subjective rating is very subject-dependent: some participants will give a higher estimate of  
267 their speech understanding than others at the same level of speech understanding. Thirdly,  
268 we investigated a homogeneous group of participants' neural responses: all normal-hearing  
269 participants between 19 and 24 years old. Therefore we do not expect large differences in  
270 speech understanding between participants at a particular speech rate.

271 To determine whether the topographies or the TRFs were significantly different from zero,  
272 we performed non-parametric permutation tests (Maris and Oostenveld, 2007). For the  
273 analysis of the acoustic TRFs, we limited the window of interest to the time region between  
274 0 and 200 ms. As speech is more difficult to understand, the latency of the neural responses  
275 to acoustic representation increases. These effects are most prominent in a time region of 0  
276 to 200 ms (Verschuere et al., 2020; Mirkovic et al., 2019; Kraus et al., 2020), explaining  
277 the rationale to limit the time window of interest. However, no time window of interest was  
278 set to determine the significance of the linguistic TRFs. We are not aware of any studies that  
279 assess the effect of linguistic tracking when speech comprehension becomes challenging.  
280 Therefore we chose not to specify a time window of interest when investigating the neural  
281 responses to linguistic representations. As observed in previous literature, linguistic TRFs

282 are associated with negative responses in central areas. Therefore, we applied this test in a  
283 one-sided fashion, i.e., we determined where the TRF was significantly negative.

284 To assess the relationship between speech rate, neural speech tracking and speech  
285 understanding, we created a linear mixed effect (LME) model using the LME4 package  
286 (Bates et al., 2015) in R with the following general formula:

$$287 \quad \text{neuralMeasure} \sim \text{rate}(+\text{rate}^2)(+\text{understanding}) + \text{random} = \text{participant}$$

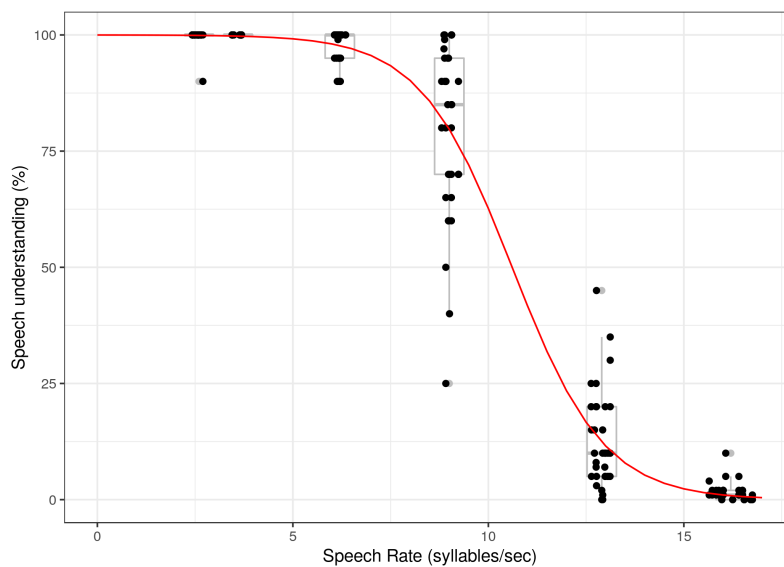
288 where “neuralMeasure” refers to neural speech tracking, TRF amplitude or TRF latency,  
289 depending on the model being investigated and “rate” refers to the speech rate the speech  
290 was presented at. Speech rate was also added as a quadratic effect, “rate<sup>2</sup>”, as we do not  
291 expect neural speech tracking will decrease or increase linearly indefinitely with increasing  
292 speech rate. Lastly, “understanding”, referring to rated speech understanding, was added to  
293 the model to investigate whether speech understanding is able to explain additional variance  
294 on top of speech rate. An additional random intercept per participant was included in the  
295 model to account for the multiple observations per participant. “Rate<sup>2</sup>” and “understanding”  
296 are added between brackets to the general formula because these factors were only included  
297 if they benefited the model. We controlled this by calculating the Akaike Information  
298 Criterion (AIC) for the model with and without “Rate<sup>2</sup>” and “understanding”. The model  
299 with the lowest AIC was selected and its residuals plot was analyzed to assess the normality  
300 assumption of the LME residuals. Unstandardized regression coefficients (beta) with 95%  
301 confidence intervals and p-value of the factors included in the model are reported in the  
302 results section.

### 3 RESULTS

#### 303 3.1 Effect of speech rate on speech understanding

304 Figure 2 shows that when speech rate increases, rated speech understanding decreases  
305 ( $r=-0.91$ ,  $p<0.001$ , Spearman rank correlation). To model the data, we fitted a sigmoid

306 function between speech understanding and speech rate. The function shows a plateau  
307 until 6 syllables/sec ( $\pm 1.5$  times the rate of normal speech). When the speech rate further  
308 increases, speech understanding drops. Because speech understanding and speech rate are  
309 highly correlated, we select speech rate for further analysis in function of neural speech  
310 tracking. As mentioned above, speech rate is more reliable than the subjectively rated  
311 speech understanding scores since it is objectively derived from the acoustic stimulus.



**Figure 2.** Rated speech understanding in function of speech rate. The dots show speech understanding per participant for every participant specific speech rate. The boxplots show the participants' results for the averaged speech rates (based on compression ratio). The red line is the sigmoid function fitted on the data over participant.

### 312 **3.2 Effect of speech rate on neural processing of speech**

313 To obtain the results in this section, we created two models: an acoustic model and a  
314 linguistic model as explained in detail in section 2.4.3.



### 315 3.2.1 Acoustic neural tracking

316 First, we investigated the acoustic model containing acoustic edges and the spectrogram.  
 317 Figure 3.A shows how accurately the acoustic model can predict the speech signal for every  
 318 electrode used. To better quantify this result, we selected the frontocentral channels based  
 319 on Lesenfants et al. (2019) (channels are highlighted in red in the inset of figure 3.B) and  
 320 averaged them per subject. This resulted in one neural tracking value per speech rate per  
 321 subject. Neural tracking of this frontocentral channel selection increased with increasing  
 322 speech rate ( $p < 0.001$ ,  $b = 7.61 \times 10^{-3}$ ,  $CI(95\%) : \pm 1.59 \times 10^{-3}$ , LME, table 3). However,  
 323 as visualized on the right in figure 3.A this increase is not monotonous, but quadratic  
 324 ( $p < 0.001$ ,  $b = -3.28 \times 10^{-4}$ ,  $CI(95\%) : \pm 8.41 \times 10^{-5}$ , LME, table 3). Finally, adding speech  
 325 understanding as an extra predictor to the model does not improve the model ( $AIC_{SR} =$   
 326  $-623$ ,  $AIC_{SR + understanding} = -605$ ).

**Table 3.** Linear Mixed Effect Model of prediction accuracies and amplitude and latency of the TRF peaks in function of speech rate for the acoustic model

Acoustic model			Rate			Rate <sup>2</sup>		
			$\beta$	CI(95%)	p	$\beta$	CI(95%)	p
Prediction accuracy			$7.61 \times 10^{-3}$	$\pm 1.59 \times 10^{-3}$	<0.001	$-3.28 \times 10^{-4}$	$\pm 8.41 \times 10^{-5}$	<0.001
Spectrogram	p1	ampl	$3.95 \times 10^{-2}$	$\pm 8.00 \times 10^{-3}$	<0.001	$-1.45 \times 10^{-3}$	$\pm 4.20 \times 10^{-4}$	<0.001
		lat	$1.02 \times 10^{-3}$	$\pm 5.35 \times 10^{-4}$	<0.001	does not improve AIC		
	p2	ampl	$-6.64 \times 10^{-3}$	$\pm 1.77 \times 10^{-3}$	<0.001	does not improve AIC		
		lat	$-9.81 \times 10^{-5}$	$\pm 7.55 \times 10^{-4}$	NS	does not improve AIC		
Acoustic edges	p1	ampl	$2.33 \times 10^{-2}$	$\pm 1.05 \times 10^{-2}$	<0.001	$-1.33 \times 10^{-3}$	$\pm 5.54 \times 10^{-4}$	<0.001
		lat	$1.05 \times 10^{-3}$	$\pm 4.58 \times 10^{-4}$	<0.001	does not improve AIC		
	p2	ampl	$-3.92 \times 10^{-3}$	$\pm 1.83 \times 10^{-3}$	<0.001	does not improve AIC		
		lat	$2.37 \times 10^{-3}$	$\pm 2.84 \times 10^{-3}$	NS	does not improve AIC		

Every line represents a different model. NS = not significant, ampl = TRF amplitude, lat = TRF latency, p1 = peak 1, p2 = peak 2

327 To better understand the obtained quadratic tendency of neural tracking as a function of  
 328 speech rate, we analyzed the acoustic features separately using TRFs. Figure 3.B visualizes

329 the averaged TRF of the frontocentral channels for the spectrogram (left panel) and for the  
330 acoustic edges (right panel). For both speech features, two significant positive peaks appear  
331 around 70 and 150 ms (horizontal bars show the TRF parts significantly different from zero).  
332 The topographies of these peaks are shown underneath the TRFs. More detailed analysis on  
333 both peaks is done by calculating the maximum value for every participant per speech rate  
334 between 0 and 90 ms (= peak value 1) and between 110 and 200 ms (= peak value 2). We  
335 investigated the amplitude and the latency of these peak values as a function of speech rate as  
336 shown in Figure 4. The amplitude of peak 1 increases quadratically with increasing speech  
337 rate for the spectrogram feature (SR:  $p < 0.001$ ,  $b = 3.95 \times 10^{-2}$ ,  $CI(95\%) : \pm 8.00 \times 10^{-3}$ ;  $SR^2$ :  
338  $p < 0.001$ ,  $b = -1.45 \times 10^{-3}$ ,  $CI(95\%) : \pm 4.20 \times 10^{-4}$ ; LME; table 3) and acoustic edges (SR:  
339  $p < 0.001$ ,  $b = 2.33 \times 10^{-2}$ ,  $CI(95\%) : \pm 1.05 \times 10^{-2}$ ;  $SR^2$ :  $p < 0.001$ ,  $b = -1.33 \times 10^{-3}$ ,  $CI(95\%)$ :  
340  $\pm 5.54 \times 10^{-4}$ ; LME; table 3). In contrast, the amplitude of peak 2 decreases with increasing  
341 speech rate (spectrogram:  $p < 0.001$ ,  $b = -6.64 \times 10^{-3}$ ,  $CI(95\%) : \pm 1.77 \times 10^{-3}$ ; acoustic edges:  
342  $p < 0.001$ ,  $b = -3.92 \times 10^{-3}$ ,  $CI(95\%) : \pm 1.83 \times 10^{-3}$ ; LME; table 3). Interestingly, the second  
343 peak for the acoustic edges even disappears when the speech rate is 9 syllables/sec or higher  
344 and speech understanding drops below 80% (Figure 3.B, horizontal bars show the TRF  
345 parts significantly different from zero). For the latency analysis of peak 2 for acoustic edges,  
346 we thus only include the latency of the peaks in the 3 easiest speech rate conditions, as no  
347 peaks (and latencies) can be found anymore at higher speech rates. The latency of peak  
348 1, for both speech features, increases with increasing speech rate (spectrogram:  $p < 0.001$ ,  
349  $b = 1.02 \times 10^{-3}$ ,  $CI(95\%) : \pm 5.35 \times 10^{-4}$ ; acoustic edges:  $p < 0.001$ ,  $b = 1.05 \times 10^{-3}$ ,  $CI(95\%)$ :  
350  $\pm 4.58 \times 10^{-4}$ ; LME; table 3), while the latency of peak 2 shows no significant relation with  
351 speech rate (spectrogram:  $p = 0.80$ ,  $b = -9.81 \times 10^{-5}$ ,  $CI(95\%) : \pm 7.55 \times 10^{-4}$ ; acoustic edges:  
352  $p = 0.11$ ,  $b = 2.37 \times 10^{-3}$ ,  $CI(95\%) : \pm 2.84 \times 10^{-3}$ ; LME; table 3).

### 353 3.2.2 Linguistic neural tracking

354 Next to acoustic neural tracking, we also investigated the effect of speech rate on linguistic  
355 neural tracking (see section 2.4.3 for more details). Figure 5.A (left panel) shows how

356 accurately the linguistic model can predict the speech signal over subjects per speech rate per  
357 channel. The channels in the cluster which drives the topography from 0 are annotated with  
358 grey markers. The higher the speech rate, the fewer channels have significant neural tracking.  
359 To quantify this, we averaged the prediction accuracy over a central channel selection based  
360 on Gillis et al. (2021b) (channels are highlighted in red in the inset of figure 5.B), resulting  
361 in one neural tracking value per speech rate per subject. As shown in figure 5.A (right),  
362 neural tracking significantly drops monotonically with increasing speech rate ( $p=0.008$ ,  
363  $b=-4.59 \times 10^{-5}$ ,  $CI(95\%): \pm 3.31 \times 10^{-5}$ , LME, table 4). Interestingly, this is the opposite  
364 trend from the acoustic model in section 3.2.1. Finally, adding speech understanding as  
365 a predictor does not improve the linguistic model ( $AIC_{SR} = -1175$ ,  $AIC_{SR + understanding} =$   
366  $-1151$ ).

367 To thoroughly investigate the neural responses to the linguistic features, we examined the  
368 TRFs of the central channel selection. Figure 5.B visualizes the averaged normalized TRF  
369 in the central channel selection for the different linguistic features per speech rate. The grey  
370 zone is where, based on Gillis et al. (2021b), we would expect a neural response. Significant  
371 responses can be found in the lower speech rates when speech can still be understood for  
372 all features. In the higher speech rates, where speech understanding is worse or absent,  
373 the linguistic neural response also disappears (Figure 5.B, horizontal bars show the TRF  
374 parts significantly different from zero). The topographies of these responses are shown  
375 in Figure 5.B underneath the TRFs. Interestingly, the topographies switch from central  
376 negativity when speech is understood to frontal negativity when speech understanding is  
377 worse or absent. To investigate whether the amplitude or latency of these peaks is related to  
378 speech rate, we calculated the minimum value for every participant within the grey zone  
379 (= peak value). For all linguistic features the peak amplitude shrinks significantly with  
380 increasing speech rate as shown in Figure 6 (Phoneme surprisal:  $p < 0.001$ ,  $b = 6.81 \times 10^{-3}$ ,  
381  $CI(95\%): \pm 2.82 \times 10^{-3}$ ; Cohort entropy:  $p = 0.008$ ,  $b = 4.21 \times 10^{-3}$ ,  $CI(95\%): \pm 3.07 \times 10^{-3}$ ;  
382 Word surprisal:  $p < 0.001$ ,  $b = 6.46 \times 10^{-3}$ ,  $CI(95\%): \pm 3.09 \times 10^{-3}$ ; Word Frequency:  $p < 0.001$ ,  
383  $b = 6.44 \times 10^{-3}$ ,  $CI(95\%): \pm 3.00 \times 10^{-3}$ ; LME; table 4). In other words, when speech becomes

384 faster and more difficult to understand, the peak amplitude of the linguistic features  
 385 decreases until it finally disappears when the speech rate is 9.0 or 12.9 syllables/sec,  
 386 or higher, and speech understanding is dropping below 90%. Similar to the analysis of the  
 387 second peak for acoustic features, we only include the latency of significant peaks. Phoneme  
 388 surprisal shows a significant increase of peak latency with increasing speech rate ( $p=0.015$ ,  
 389  $b=5.47 \times 10^{-3}$ ,  $CI(95\%): \pm 4.20 \times 10^{-3}$ ; LME; table 4). Cohort entropy, word surprisal and  
 390 word frequency, on the other hand, reveal no significant effect of speech rate on peak latency  
 391 (Cohort entropy:  $p=0.18$ ,  $b=3.59 \times 10^{-3}$ ,  $CI(95\%): \pm 5.11 \times 10^{-3}$ ; Word surprisal:  $p=0.23$ ,  
 392  $b=-3.33 \times 10^{-3}$ ,  $CI(95\%): \pm 5.38 \times 10^{-3}$ ; Word frequency:  $p=0.95$ ,  $b=-1.09 \times 10^{-4}$ ,  $CI(95\%):$   
 393  $\pm 3.06 \times 10^{-3}$ ; LME; table 4).

**Table 4.** Linear Mixed Effect Model of prediction accuracy and amplitude and latency of the TRF peaks in function of speech rate for the linguistic model

Linguistic model		rate			rate <sup>2</sup>		
		$\beta$	CI(95%)	p	$\beta$	CI(95%)	p
Prediction accuracy		$-4.59 \times 10^{-5}$	$\pm 3.31 \times 10^{-5}$	0.0081	does not improve AIC		
Phoneme surprisal	ampl	$6.81 \times 10^{-3}$	$\pm 2.82 \times 10^{-3}$	<0.001	does not improve AIC		
	lat	$5.47 \times 10^{-3}$	$\pm 4.20 \times 10^{-3}$	0.015	does not improve AIC		
Cohort entropy	ampl	$4.21 \times 10^{-3}$	$\pm 3.07 \times 10^{-3}$	0.008	does not improve AIC		
	lat	$3.59 \times 10^{-3}$	$\pm 5.11 \times 10^{-3}$	NS	does not improve AIC		
Word surprisal	ampl	$6.46 \times 10^{-3}$	$\pm 3.09 \times 10^{-3}$	<0.001	does not improve AIC		
	lat	$-3.33 \times 10^{-3}$	$\pm 5.38 \times 10^{-3}$	NS	does not improve AIC		
Word frequency	ampl	$6.44 \times 10^{-3}$	$\pm 3.00 \times 10^{-3}$	<0.001	does not improve AIC		
	lat	$-1.09 \times 10^{-4}$	$\pm 3.06 \times 10^{-3}$	NS	does not improve AIC		

Every line represents a different model. NS = not significant, ampl = TRF amplitude, lat = TRF latency

## 4 DISCUSSION

394 We aimed to investigate whether neural speech processing can capture the effect of gradually  
 395 decreasing speech understanding by manipulating the speech rate. With increasing speech

396 rate, we found that neural tracking of the acoustic features increased, while neural tracking  
397 of the linguistic features decreased.

#### 398 **4.1 Effect of speech rate on acoustic neural processing of speech**

399 We found an increase of acoustic neural tracking with increasing speech rate and thus  
400 decreasing speech understanding, confirming our hypothesis. When speech becomes faster,  
401 the model is better at predicting acoustic speech features.

402 This increase of acoustic neural tracking with decreasing speech understanding seems  
403 discrepant with previous research trying to link acoustic neural tracking to speech  
404 understanding using, for example, speech-in-noise paradigms (Vanthornhout et al., 2018;  
405 Verschueren et al., 2020; Ding and Simon, 2013; Iotzov and Parra, 2019; Etard and  
406 Reichenbach, 2019). The experimental paradigm could explain this discrepancy. Previous  
407 studies used, for example, noise to manipulate speech understanding. In those cases,  
408 decreased neural tracking was accompanied by a decrease in speech understanding and an  
409 acoustically degraded speech signal. Because speech understanding and signal-to-noise  
410 ratio are highly correlated, it is challenging to unravel to what extent the decreased neural  
411 tracking is driven by decreased speech understanding or the signal-to-noise ratio used to  
412 vary speech understanding. In this study, we manipulated speech understanding by speeding  
413 up the speech signal and preserving its signal-to-noise ratio, in contrast to the speech in  
414 noise studies. We hypothesize that the brain mainly responds to acoustic boundaries, i.e.  
415 onsets of sounds, which are more prominent in the faster speech presented in this study,  
416 explaining the increasing tendency. When presenting speech in noise, acoustic boundaries  
417 can be masked and, therefore, more challenging to observe. Therefore, it is difficult to  
418 attribute this decrease in acoustic neural tracking: a decrease in speech understanding or a  
419 decrease in neural detection of acoustic boundaries, or a combination of both?

420 In addition, because the speech is sped up, the duration of the silences in between words  
421 or sentences inherently decreases, which increases the amount of speech data allowing

422 the model to improve its estimate of the TRF and obtain higher prediction accuracies.  
423 However, this increase of acoustic neural tracking and speech rate is not linear but quadratic,  
424 saturating, and even decreasing at very high speech rates. This may be due to the stimulus  
425 characteristics. When increasing the speech rate, the spectrogram and acoustic edges contain  
426 more and more peaks. Possibly these peaks are occurring so fast after each other making it  
427 difficult for the brain to perceive them still (see Figure 1). A different hypothesis is related  
428 to the motor cortex. When participants listen to the speech, they tend to mimic the speech  
429 in their brain, activating neural activity in motor areas (Casas et al., 2021). However, most  
430 speakers cannot produce speech as fast as 16.2 syllables/sec. Hence, the corresponding  
431 mouth movements are unnatural, which implies that the listener cannot mimic the speech  
432 in their brain anymore, decreasing the related responses in the motor areas and thus brain  
433 responses to the acoustic speech features in general.

434 To better understand the observed quadratic tendency of acoustic neural tracking with  
435 increasing speech rate, we investigated the TRFs of the speech features separately. Two  
436 significant peaks with opposite behavior could be observed for both acoustic features. The  
437 first peak is the largest, and its amplitude increases quadratically with increasing speech  
438 rate, similar to the previously discussed neural acoustic tracking results. On the other  
439 hand, the second peak amplitude decreases with increasing speech rate. This discrepancy  
440 is intriguing as it suggests that both peaks have different underlying brain processes as  
441 confirmed by literature (Picton, 2011; Brodbeck and Simon, 2020). Peak 1 occurs relatively  
442 fast, around 50 ms, and is probably mainly related to the acoustics of the incoming  
443 speech and thus benefits from an increased speech rate. Peak 2, on the other hand, occurs  
444 somewhat later, around 150 ms, and could, in addition to the acoustics, be influenced by  
445 top-down processing related to speech understanding and attention (Ding and Simon, 2012;  
446 Vanthornhout et al., 2019). Besides the amplitude, we also investigated the latencies. The  
447 latency of the first peak increases with increasing speech rate. Increased latencies are often  
448 observed in more complex conditions with a higher task demand, like for example lower  
449 stimulus intensity, vocoded speech or speech in noise (Mirkovic et al., 2019; Verschueren

450 et al., 2021; Kraus et al., 2020). The latency of the neural responses can also be related to  
451 neural processing efficiency (Bidelman et al., 2019; Gillis et al., 2021a). In more detail, a  
452 larger latency indicates that more processing time is required to process the same speech  
453 characteristics, showing reduced neural processing efficiency. More words and phonemes  
454 need to be processed as the speech rate increases, resulting in a more challenging condition  
455 to process the incoming speech.

## 456 **4.2 Effect of speech rate on linguistic neural processing of speech**

457 When speech becomes faster, speech understanding drops. Interestingly, this same decrease  
458 can be observed in linguistic neural tracking (in contrast to acoustic neural tracking, section  
459 4.1). To the best of our knowledge, this is the first study that evaluates linguistic neural  
460 tracking when manipulating the level of speech understanding as a gradual effect. The  
461 studies of Brodbeck et al. (2018) and Broderick et al. (2018) using a two-talker paradigm  
462 are most comparable. They compared two conditions, i.e., intelligible and attended speech  
463 versus unintelligible and ignored speech, but not the spectrum in between. Nevertheless,  
464 their findings converge with our results and support our hypothesis of linguistic neural  
465 tracking as a neural marker of speech understanding. When the speech is not understood or  
466 ignored, the brain does not track the linguistic aspects of the speech, while for intelligible  
467 speech linguistic tracking is present.

468 To better understand the observed decrease of linguistic neural tracking with increasing  
469 speech rate, we investigated the TRFs of the speech representations separately. We observed  
470 a characteristic negative peak for each linguistic representation as observed in previous  
471 literature (e.g. Brodbeck et al., 2018; Gillis et al., 2021b; Weissbart et al., 2019). For  
472 the phoneme-related features, phoneme surprisal and cohort entropy, this peak occurs  
473 around 250 ms. For the word-related features, word surprisal and word frequency, this peak  
474 occurs somewhat later, around 350 ms. The difference in timescale between both feature  
475 groups could be linked to the different speech processing stages (phonemes versus words)



476 they represent (Van Canneyt et al., 2021). Regarding the topographies of these peaks, the  
477 understandable speech conditions are associated with a typical topography, similar to the  
478 classical N400 responses characterized by central negative channels. As speech becomes  
479 less understandable, i.e., the speech rate increases, the associated topography disappears.

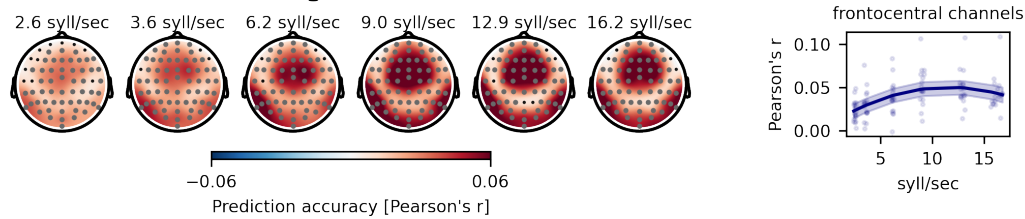
480 For all speech features, the amplitude of this negative peak decreases with increasing  
481 speech rate until it disappears at speech rates as high as 9.0 tot 12.9 syllables/sec. Gillis  
482 et al. (2021b) already showed that these linguistic representations have an added value  
483 above and beyond acoustic and lexical representations. However, the authors did not  
484 compare intelligible to unintelligible speech. Here, we elegantly showed that as the speech  
485 becomes less understandable but remains audible and acoustically intact (in contrast to  
486 speech in noise studies or vocoder studies), the characteristic negative peak decreases and  
487 finally disappears. Altogether, our results suggest that these characteristic negative peaks to  
488 linguistic representations could be neural correlates of speech understanding.

### 489 **4.3 Conclusion**

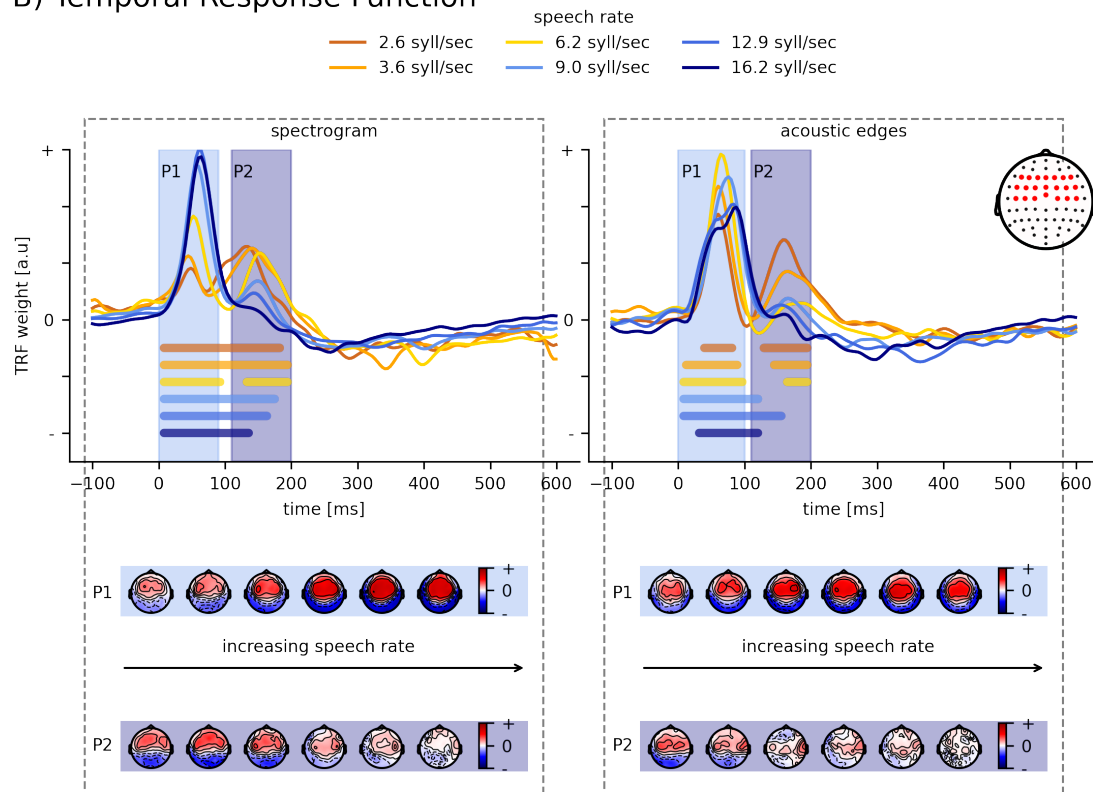
490 Using a speech rate paradigm, we map how the level of speech understanding affects  
491 acoustic and linguistic neural speech processing. When speech rate increases, acoustic  
492 neural tracking increases, although speech understanding drops. However, the amplitude  
493 of the later acoustic neural response decreases with increasing speech rate, suggesting  
494 influence of top-down processing related to speech understanding and attention. In contrast,  
495 linguistic neural tracking decreases with increasing speech rate and even disappears when  
496 speech is no longer understood. Altogether, this suggests that linguistic neural tracking  
497 could possibly be a more direct predictor of speech understanding compared to acoustic  
498 neural tracking.



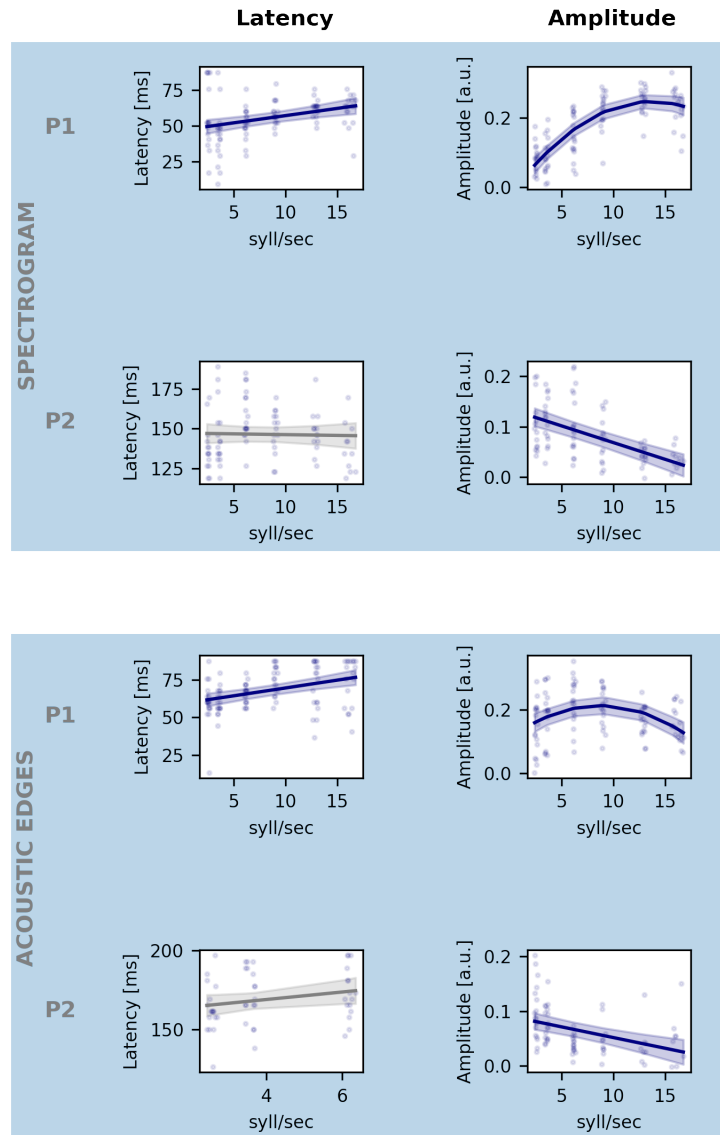
### A) Acoustic neural tracking



### B) Temporal Response Function



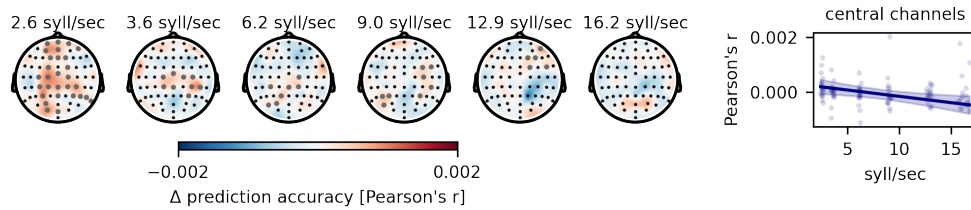
**Figure 3. Effect of speech rate on acoustic tracking** Panel A: visualization of the average prediction accuracy across participants for each speech rate. The annotated grey channels indicate the cluster which drives the significant difference from 0. How acoustic tracking, averaged across frontocentral channels, changes according to the speech rate is shown on the right. Panel B: Normalized TRFs of the spectrogram and acoustic edges. The bold horizontal lines indicate where the TRFs are significantly different from 0 (the same color as the TRF of the considered speech rate). The topographies below show the associated peak topographies in the TRF.



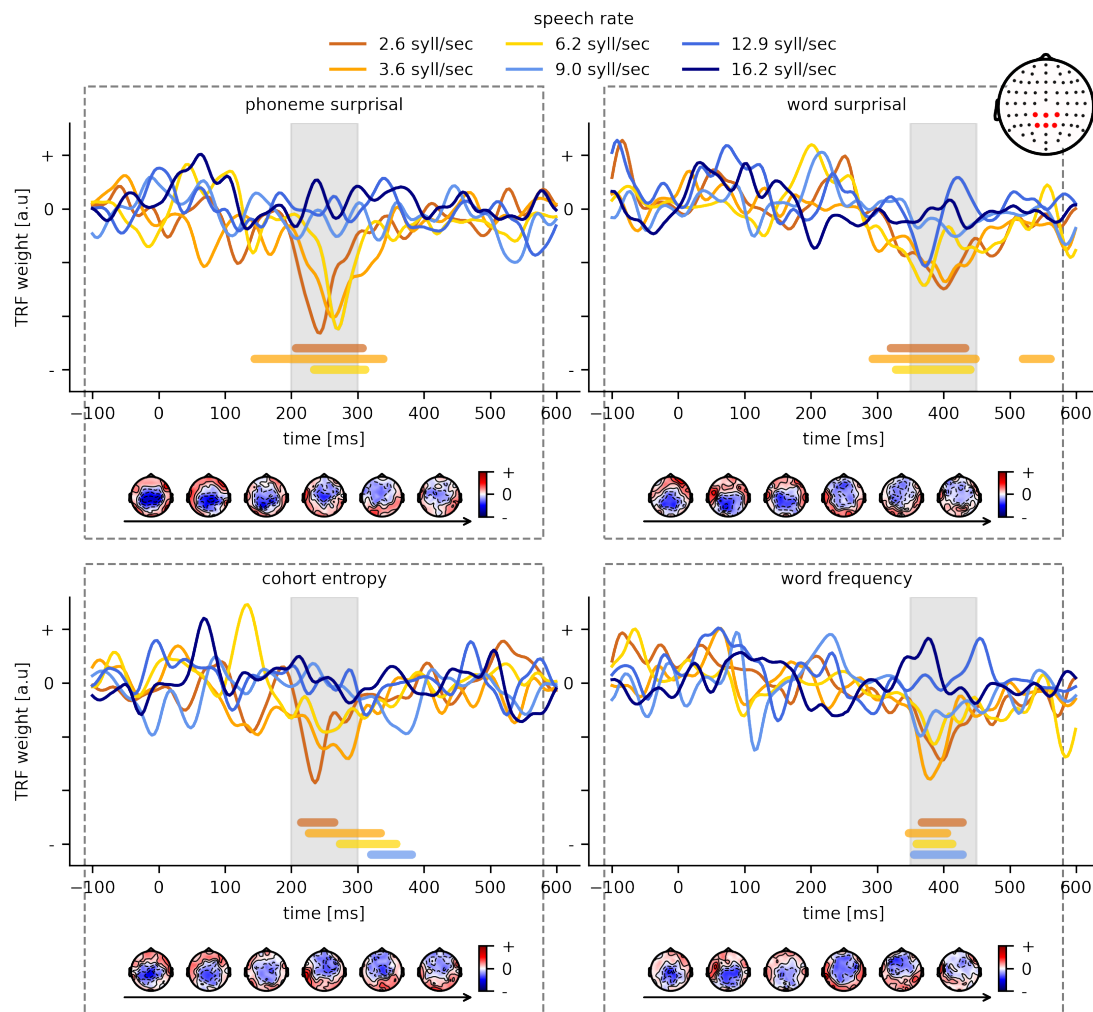
**Figure 4. Effect of speech rate on the amplitude and latency of acoustic representation** Each row shows the effect of speech rate on latency (left plot) and the amplitude (right plot) of the neural response to spectrogram (top row) and acoustic edges (bottom row) for respectively the first and second identified peak as indicated on Figure 3. The blue line shows the model's prediction for each speech rate; the shaded area indicates the confidence interval of the model's prediction. The non-significant models are shown in grey. Remark that we only include the latency of significant peaks for the latency analysis.

26

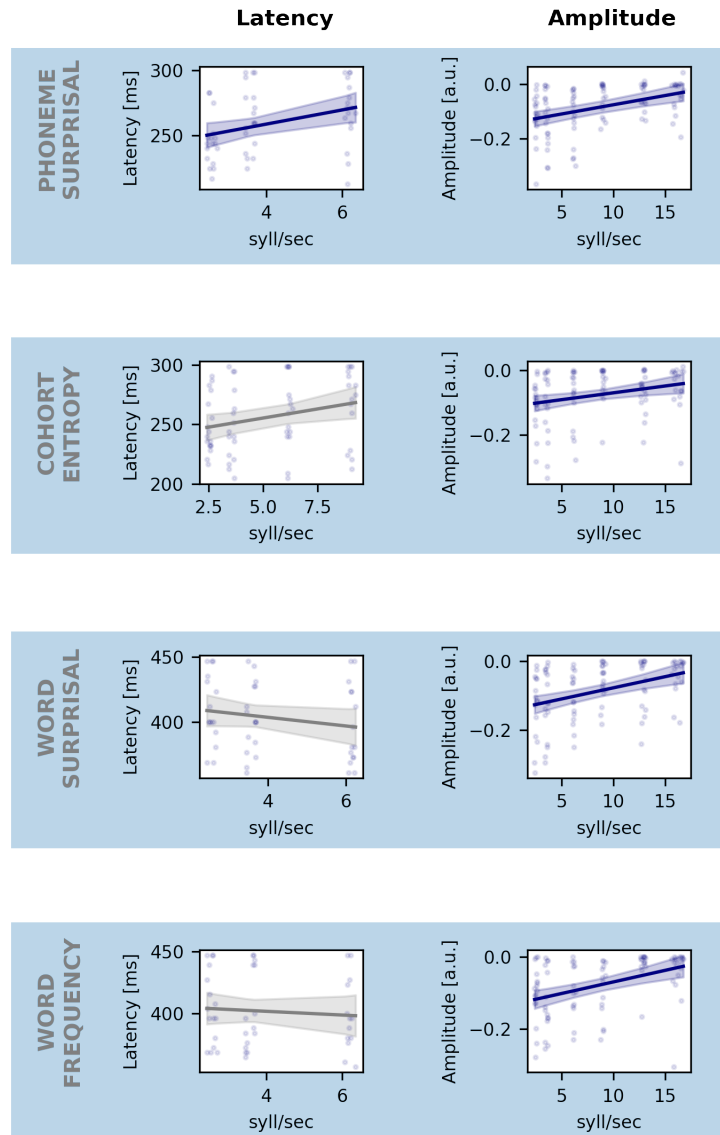
### A) Linguistic neural tracking



### B) Temporal Response Function



**Figure 5. Effect of speech rate on linguistic tracking** Panel A: visualization of the added value of linguistic representations across participants for each speech rate. The annotated grey channels indicate the cluster which drives the significant difference from 0. **H27** linguistic tracking, averaged across central channels, changes according to the speech rate is shown on the right. Panel B: The associated normalized TRFs for the linguistic representations. The bold horizontal lines indicate where the TRFs are significantly different from 0 (the same color as the TRF of the considered speech rate). The topographies below indicate the associated peak topographies to the TRF in the grey shaded area. The horizontal arrow underneath the topographies indicates the increasing speech rate.



**Figure 6. Effect of speech rate on the amplitude and latency of linguistic representation** Each row shows the effect of speech rate on latency (left plot) and the amplitude (right plot) of the neural response to phoneme surprisal (top row), cohort entropy (second row), word surprisal (third row) and word frequency (bottom row). The blue line shows the model's prediction for each speech rate; the shaded area indicates the confidence interval of the model's prediction. The non-significant models are shown in grey. Remark that we only include the latency of significant peaks for the latency analysis.

499 Acknowledgment: The authors would like to thank Sofie Keunen and Elise Verwaerde for  
500 their help in data acquisition. Funding: The presented study received funding from the  
501 European Research Council (ERC) under the European Union’s Horizon 2020 research and  
502 innovation programme (Tom Francart; grant agreement No. 637424). Research of Eline  
503 Verschueren (PhD grant: SB 1S86118N), Marlies Gillis (PhD grant: SB 1SA0620N) and  
504 Jonas Vanthornout (postdoctoral grant: 1290821N) was funded by the Research Foundation  
505 Flanders (FWO).

## REFERENCES

- 506 Ahissar, E., Nagarajan, S., Ahissar, M., Protopapas, A., Mahncke, H., and Merzenich,  
507 M. M. (2001). Speech comprehension is correlated with temporal response patterns  
508 recorded from auditory cortex. *Proceedings of the National Academy of Sciences of the*  
509 *United States of America* 98, 13367–72. doi:10.1073/pnas.201400998
- 510 Baayen, R. H., Piepenbrock, R., and Gulikers, L. (1996). The celex lexical database  
511 (cd-rom)
- 512 Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects  
513 models using lme4. *Journal of Statistical Software* 67, 1–48. doi:10.18637/jss.v067.i01
- 514 Bidelman, G. M., Price, C. N., Shen, D., Arnott, S. R., and Alain, C. (2019). Afferent-  
515 efferent connectivity between auditory brainstem and cortex accounts for poorer speech-  
516 in-noise comprehension in older adults. *Hearing research* 382, 107795
- 517 Boursma, P. and Weenink, D. (2018). Praat: doing phonetics by computer [computer  
518 program]. *version 6.0.37, retrieved 3 Februari 2018 from <http://www.praat.org/>*
- 519 [Dataset] Brodbeck, C. (2020). Eelbrain 0.34. [http://doi.org/10.5281/zenodo.](http://doi.org/10.5281/zenodo.3923991)  
520 3923991
- 521 Brodbeck, C., Hong, L. E., and Simon, J. Z. (2018). Rapid transformation from auditory to  
522 linguistic representations of continuous speech. *Current Biology* 28, 3976–3983
- 523 Brodbeck, C. and Simon, J. Z. (2020). Continuous speech processing. *Current Opinion in*  
524 *Physiology* 18, 25–31. doi:10.1016/j.cophys.2020.07.014

- 525 Broderick, M. P., Anderson, A. J., Liberto, G. M. D., Crosse, M. J., and Lalor, E. C. (2018).  
526 Electrophysiological Correlates of Semantic Dissimilarity Reflect the Comprehension of  
527 Natural , Report Electrophysiological Correlates of Semantic Dissimilarity Reflect the  
528 Comprehension of Natural , Narrative Speech. *Current Biology* 28, 803–809. doi:10.  
529 1016/j.cub.2018.01.080
- 530 Casas, A. S. H., Lajnef, T., Pascarella, A., Guiraud-vinatea, H., Laaksonen, H., Bayle,  
531 D., et al. (2021). Neural oscillations track natural but not artificial fast speech: Novel  
532 insights from speech-brain coupling using meg. *NeuroImage* 244, 118577
- 533 David, S. V., Mesgarani, N., and Shamma, S. A. (2007). Estimating sparse spectro-  
534 temporal receptive fields with natural stimuli. *Network: Computation in Neural Systems*  
535 18, 191–212. doi:10.1080/09548980701609235
- 536 Ding, N. and Simon, J. Z. (2012). Emergence of neural encoding of auditory objects while  
537 listening to competing speakers. *Proceedings of the National Academy of Sciences of the*  
538 *United States of America* 109, 11854–9. doi:10.1073/pnas.1205381109
- 539 Ding, N. and Simon, J. Z. (2013). Adaptive Temporal Encoding Leads to a Background-  
540 Insensitive Cortical Representation of Speech. *Journal of Neuroscience* 33, 5728–5735.  
541 doi:10.1523/JNEUROSCI.5297-12.2013
- 542 Duchateau, J., Kong, Y. O., Cleuren, L., Latacz, L., Roelens, J., Samir, A., et al. (2009).  
543 Developing a reading tutor: Design and evaluation of dedicated speech recognition and  
544 synthesis modules. *Speech Communication* 51, 985–994
- 545 Elzhov, T. V., Mullen, K. M., Spiess, A.-N., and Bolker, B. (2016). *minpack.lm: R Interface*  
546 *to the Levenberg-Marquardt Nonlinear Least-Squares Algorithm Found in MINPACK,*  
547 *Plus Support for Bounds.* R package version 1.2-1
- 548 Etard, O. and Reichenbach, T. (2019). Neural speech tracking in the theta and in the  
549 delta frequency band differentially encode clarity and comprehension of speech in noise.  
550 *Journal of Neuroscience* 39, 5750–5759
- 551 Francart, T., van Wieringen, A., and Wouters, J. (2008). APEX 3: a multi-purpose test  
552 platform for auditory psychophysical experiments. *Journal of Neuroscience Methods*

- 553 172, 283–293
- 554 Gillis, M., Decrui, L., Vanthornhout, J., and Francart, T. (2021a). Hearing loss is associated  
555 with delayed neural responses to continuous speech. *bioRxiv*
- 556 Gillis, M., Vanthornhout, J., Simon, J. Z., Francart, T., and Brodbeck, C. (2021b).  
557 Neural markers of speech comprehension: measuring EEG tracking of linguistic speech  
558 representations, controlling the speech acoustics. *Journal of neuroscience* , [Accepted for  
559 publication]
- 560 [Dataset] Heeris, J. (2014). Gammatone filterbank toolkit 1.0. [https://github.com/](https://github.com/detly/gammatone)  
561 [detly/gammatone](https://github.com/detly/gammatone)
- 562 Hertrich, I., Dietrich, S., Trouvain, J., Moos, A., and Ackermann, H. (2012). Magnetic  
563 brain activity phase-locked to the envelope, the syllable onsets, and the fundamental  
564 frequency of a perceived speech signal. *Psychophysiology* 49, 322–334. doi:10.1111/j.  
565 1469-8986.2011.01314.x
- 566 Iotzov, I. and Parra, L. (2019). EEG can predict speech intelligibility To. *J. Neural Eng.* 16,  
567 036008 (11p). doi:10.1088/1741-2552/ab07fe
- 568 Keuleers, E., Brysbaert, M., and New, B. (2010). Subtlex-nl: A new measure for dutch  
569 word frequency based on film subtitles. *Behavior research methods* 42, 643–650
- 570 Koskinen, M., Kurimo, M., Gross, J., Hyv, A., and Hari, R. (2020). NeuroImage Brain  
571 activity reflects the predictability of word sequences in listened continuous speech.  
572 *NeuroImage* 219, 116936. doi:10.1016/j.neuroimage.2020.116936
- 573 Kraus, F., Tune, S., Ruhe, A., Obleser, J., and Woestmann, M. (2020). Unilateral acoustic  
574 degradation delays attentional separation of competing speech. *bioRxiv*
- 575 Lesenfants, D., Vanthornhout, J., Verschueren, E., Decrui, L., and Francart, T. (2019).  
576 Predicting individual speech intelligibility from the neural tracking of acoustic- and  
577 phonetic-level speech representations. *Hearing Research* 380, 1–9. doi:10.1016/j.heares.  
578 2019.05.006
- 579 Maris, E. and Oostenveld, R. (2007). Nonparametric statistical testing of eeg-and meg-data.  
580 *Journal of neuroscience methods* 164, 177–190



- 581 Mirkovic, B., Debener, S., Schmidt, J., Jaeger, M., and Neher, T. (2019). Effects of  
582 directional sound processing and listener's motivation on eeg responses to continuous  
583 noisy speech: Do normal-hearing and aided hearing-impaired listeners differ? *Hearing*  
584 *Research* 377, 260–270
- 585 Müller, J. A., Wendt, D., Kollmeier, B., Debener, S., Brand, T., and Hunter, C. R. (2019).  
586 Effect of Speech Rate on Neural Tracking of Speech. *Frontiers in psychology* 10, 1–15.  
587 doi:10.3389/fpsyg.2019.00449
- 588 Nourski, K. V., Reale, R. A., Oya, H., Kawasaki, H., Kovach, C. K., Chen, H., et al. (2009).  
589 Temporal Envelope of Time-Compressed Speech Represented in the Human Auditory  
590 Cortex. *The Journal of Neuroscience* 29, 15564–15574. doi:10.1523/JNEUROSCI.  
591 3065-09.2009
- 592 Picton, T. W. (2011). *Human Auditory Evoked Potentials* (San Diego: Plural Publishing  
593 inc.)
- 594 Shannon, R. V., Zeng, F. G., Kamath, V., Wygonski, J., and Ekelid, M. (1995). Speech  
595 recognition with primarily temporal cues. *Science* 270, 303–304
- 596 Slaney, M. (1998). Auditory toolbox. *Interval Research Corporation, Tech. Rep* 10
- 597 Somers, B., Francart, T., and Bertrand, A. (2018). A generic EEG artifact removal  
598 algorithm based on the multi-channel Wiener filter. *Journal of neural engineering* 15.  
599 doi:10.1088/1741-2552/aaac92
- 600 Van Canneyt, J., Gillis, M., Vanthornhout, J., and Francart, T. (2021). Neural tracking as an  
601 objective measure of auditory perception and speech intelligibility. *bioRxiv*
- 602 Vanthornhout, J., Decruy, L., and Francart, T. (2019). Effect of task and attention on neural  
603 tracking of speech. *BioRxiv* , doi: <http://dx.doi.org/10.1101/568204>doi:10.3389/fpsyg.  
604 2019.00449
- 605 Vanthornhout, J., Decruy, L., Wouters, J., Simon, J. Z., and Francart, T. (2018). Speech  
606 intelligibility predicted from neural entrainment of the speech envelope. *JARO* 19,  
607 181–191. doi:10.1007/s10162-018-0654-z



- 608 Verschueren, E., Vanthornhout, J., and Francart, T. (2020). The effect of stimulus choice  
609 on an EEG-based objective measure of speech intelligibility. *Ear & Hearing*, Publish  
610 Ahead of preprintdoi:<https://doi.org/10.1101/421727>
- 611 Verschueren, E., Vanthornhout, J., and Francart, T. (2021). The effect of stimulus intensity  
612 on neural envelope tracking. *Hearing Research* 403, 108175
- 613 Weissbart, H., Kandylaki, K. D., and Reichenbach, T. (2019). Cortical Tracking of  
614 Surprisal during Continuous Speech Comprehension. *Journal of cognitive neuroscience*  
615 32, 155–166. doi:10.1162/jocn.a\_01467