

1

Neural tracking as a diagnostic tool to assess the auditory pathway

2

Marlies Gillis^{1A}, Jana Van Canneyt^{1A}, Tom Francart^{1B}, Jonas Vanthornhout^{1B}

3

¹: Experimental Oto-Rhino-Laryngology, Department of Neurosciences, Leuven Brain Institute, KU Leuven,

4

Belgium

5

^A: Shared first authorship

6

^B: Shared last authorship

7

Corresponding authors:

8

Marlies Gillis (marlies.gillis@kuleuven.be)

9 **Abstract**

When a person listens to sound, the brain time-locks to specific aspects of the sound. This is called neural tracking and it can be investigated by analysing neural responses (e.g., measured by electroencephalography) to continuous natural speech. Measures of neural tracking allow for an objective investigation of a range of auditory and linguistic processes in the brain during natural speech perception. This approach is more ecologically valid than traditional auditory evoked responses and has great potential for research and clinical applications. This article reviews the neural tracking framework and highlights three prominent examples of neural tracking analyses: neural tracking of the fundamental frequency of the voice (f_0), the speech envelope and linguistic features. Each of these analyses provides a unique point of view into the human brain's hierarchical stages of speech processing. F_0 -tracking assesses the encoding of fine temporal information in the early stages of the auditory pathway, i.e., from the auditory periphery up to early processing in the primary auditory cortex. Envelope tracking reflects bottom-up and top-down speech-related processes in the auditory cortex and is likely necessary but not sufficient for speech intelligibility. Linguistic feature tracking (e.g. word or phoneme surprisal) relates to neural processes more directly related to speech intelligibility. Together these analyses form a multi-faceted objective assessment of an individual's auditory and linguistic processing.

10 *Keywords:* Neural tracking, Speech intelligibility, EEG, f_0 tracking, envelope tracking, linguistic features

11

12 **1. Introduction**

13 Understanding speech is a complex process that relies on activation and cooperation between various brain regions.
14 Different characteristics of incoming speech are processed in different brain regions. Roughly, purely acoustic pro-
15 cessing of the speech occurs in subcortical areas and the primary auditory cortex. In contrast, segmentation of words
16 and phonemes occurs in temporal regions of the brain, and integration of words into their context occurs in language-
17 related brain regions, such as superior temporal gyrus and inferior frontal gyrus (Brodbeck et al., 2018c,a). However,
18 only if all stages in this neural pathway are successful speech understanding can be achieved.

19 Audiologists rely on an extensive test battery to assess a person's speech understanding. A commonly performed test
20 is to let a subject recall a list of sentences. The outcome of this test expresses speech understanding as a percentage
21 of correctly recalled words. However, such behavioural tests have some disadvantages. First, the subject must listen
22 actively to the stimulus and recall the words. Although this seems like an easy task, it can be challenging or impossible
23 for many populations: persons with locked-in syndrome, young children, persons with aphasia, etc. Although they
24 might understand the speech, they might not be able to recall the heard words. Second, the outcomes of these tests do

25 not pinpoint the origin of the deficit. Is the deficit situated cortically, indicating an issue with the higher-order language
26 processing, or peripherally, suggesting a hearing loss? Third, these behavioural tests rely on highly controlled stand-
27 alone sentences or words spoken by a professional speaker. Such speech material has limited contextual information.
28 Therefore, it does not resemble a typical day-to-day listening environment.

29 One can use a more objective approach, such as neurophysiological measures, i.e., metrics derived from brain signals
30 to overcome these issues. Traditional neurophysiological measures, like the auditory brainstem response (ABR),
31 the auditory steady-state response (ASSR) or the frequency following response (FFR), require EEG measurement.
32 During such a measurement, a participant listens to repetitive presentations of a short sound stimulus (for a review,
33 see Picton (2010)). Typical stimuli include clicks, tones, chirps and vowels. The repetitive stimulation is necessary as
34 response instances need to be averaged to reduce measurement noise, but it is highly unnatural and demotivating for
35 the listener (Theunissen et al., 2000; Hamilton and Huth, 2018). In recent years, technical advances in data analysis
36 have made it possible to analyse neural responses measured while a participant listens to continuous natural speech
37 without repetition (for a review, see Brodbeck and Simon, 2020). These neural responses time-lock to the presented
38 speech and this phenomenon is called neural tracking. Measuring neural responses to continuous natural speech was
39 originally proposed by Lalor et al. (Lalor et al., 2009; Lalor and Foxe, 2010) and the methods were further developed
40 by, amongst others, Ding and Simon (2012a,b), O’Sullivan et al. (2015) and Crosse et al. (2016a).

41 The possibility of investigating continuous speech processing by measuring neural tracking is an important innovation.
42 Humans do not communicate with the stimuli of traditional objective measures: repetitive tones or clicks. Context-
43 rich continuous speech better approximates natural language use, and as a result, research findings with these stimuli
44 are more relevant for auditory processing in day-to-day communication (Kei et al., 1999; Pichora-Fuller et al., 2016;
45 Hamilton and Huth, 2018; Keidser et al., 2020). Moreover, continuous speech is more comfortable and enjoyable for
46 the listener. The stimulus can even be targeted towards the population of interest: e.g. a fairy tale for young children
47 or a podcast for adults. When participants are interested in the content of the stimulus, they maintain attention for
48 longer, and as a result, the neural response measurement may be of higher quality. Finally, natural speech stimuli are
49 better suited for research with hearing aids. Hearing aid signal processing is designed specifically for natural speech
50 and may behave unpredictably with artificial sounds, corrupting the experiment.

51 In this article, we will give an overview of neural tracking of continuous speech with primary emphasis on neural
52 tracking of single-talker speech. However, a promising and emerging related field is the use of neural tracking to
53 determine the focus of attention, called auditory attention decoding (AAD). When listening to a target speaker in a
54 mixture of multiple speakers, higher neural tracking is observed for the attended speaker than ignored speakers. Au-
55 ditory attention decoding may enable smart hearing devices that can reinforce the attended speaker while attenuating
56 the unattended speakers (Geirnaert et al., 2021).

57 In this article, we evaluated neural tracking as a diagnostic tool to assess multiple levels of the auditory pathway. Such

58 a tool would be based on one EEG recording of a person listening to natural speech, which is analysed in various ways
59 to assess whether there is a speech understanding deficit and if so, where it is situated in the auditory pathway (e.g.,
60 peripherally or cortically).

61 **2. Methods to measure neural tracking**

62 The most common approach for measuring neural tracking are linear, decoding or encoding models. These models
63 measure the amount of neural tracking, i.e., how strongly the neural responses time-lock to a stimulus feature and are
64 discussed in more detail in the following subsections.

65 Measuring neural tracking requires two inputs: neural responses in the form of single-channel or multi-channel EEG
66 (or MEG) and one or more features representing the stimulus (see section 3). A linear relation between the EEG and
67 the stimulus feature is modelled to investigate how well the stimulus information is encoded in the neural activity.
68 Linear modelling is possible in two ways: reconstructing the feature from the EEG (decoding, section 2.1) and,
69 conversely, predicting the EEG from the feature (encoding, section 2.2) or a combined approach (CCA (de Cheveigné
70 et al., 2018) as discussed in section 2.5). As discussed below, the decoding and encoding analyses provide different
71 but complementary information about neural tracking.

72 *2.1. Decoding modelling*

73 In decoding modelling, one reconstructs the stimulus feature from a weighted sum of the EEG signals of the different
74 channels and their time-shifted versions. The time-shifted versions are included to account for the time difference
75 between the sound and the associated neural response.

76 The decoding modelling procedure is visualised in panel A of Figure 1. First, the weights that provide the optimal
77 reconstruction are determined based on the time-shifted EEG and its corresponding stimulus feature. Then those
78 weights are applied to the EEG, resulting in a reconstructed stimulus feature. The reconstructed feature is correlated
79 with the actual stimulus feature of the test data to determine the reconstruction accuracy. This reconstruction accuracy
80 is a measure of neural tracking as it indicates how well the stimulus information is time-locked with the EEG. A
81 higher reconstruction accuracy will therefore reflect higher neural tracking of the speech.

82 The decoding modelling approach is a powerful analysis tool since the information of multiple EEG channels (often 32
83 or more) can be combined to reconstruct a stimulus feature with often only one dimension (although multi-dimensional
84 features are possible).

85 A disadvantage of decoding modelling is that the weights of the model are extraction filters which cannot and should
86 not be interpreted to investigate the spatial pattern of the response (Haufe et al., 2014). Extraction filters do not
87 always have large weights when the corresponding EEG channels contain a lot of response information. When an
88 EEG channel captures information about a noise component, it can be used in the modelling process to remove the

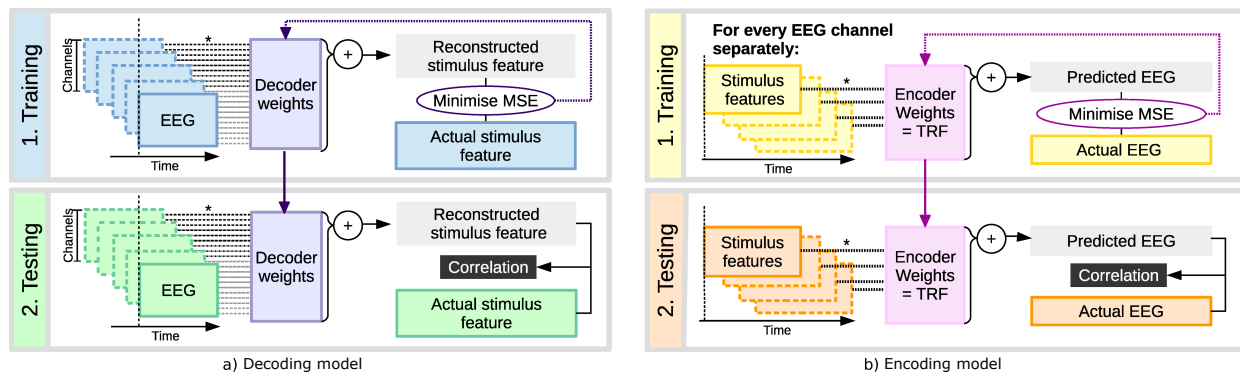


Figure 1: *A. Schematic representation of decoding modelling.* The stimulus feature is reconstructed in decoding modelling based on a linear combination of time-shifted EEG data. In the training phase, the model is estimated by optimising the decoder weights to minimise the MSE (mean squared error) between the reconstructed stimulus feature and the actual stimulus feature for a training data set. Then, the weights are applied in the testing phase to reconstruct the stimulus feature for the testing dataset. The final output is the correlation between the reconstructed and actual stimulus features for the testing dataset. The division in training and testing dataset is done according to the cross-validation technique (described in 2.4) *B. Schematic representation of encoding modelling.* In encoding modelling, the EEG data in each EEG channel is predicted based on a linear combination of time-shifted stimulus features. Again, the encoder weights or TRFs (temporal response functions) are estimated by minimising the reconstruction MSE for a training data set. Then the TRFs can be studied as is or used to predict the EEG for a testing data set. The output of the testing phase is the correlation between the predicted EEG and the actual EEG. The division in training and testing dataset is done according to the cross-validation technique (described in 2.4)

89 noise component from other EEG channels. As a result, some channels may receive large weights because they are
 90 helpful for noise reduction purposes and not because they contain response information (Montoya-Martínez et al.,
 91 2021). Haufe et al. (2014) defines an inversion method to make the topography of the decoding model interpretable,
 92 or one can use an encoding model.

93 Another disadvantage of decoding modelling is that the evaluation of the model, i.e., the loss function, should be
 94 adapted when investigating the reconstruction accuracy of sparse, impulsive features. These sparse, impulsive features
 95 can, for example, code the onset of a phoneme or word. They consist of an array of zeroes with a given value at the
 96 onset of a word or phoneme. Such features cannot be reconstructed from the continuous EEG signals with a linear
 97 model. Therefore, if sparse, impulsive features are used, a correlation between the actual and reconstructed features
 98 is sub-optimal. The evaluation of a decoding model, i.e., taking a correlation, should be adapted when using sparse,
 99 impulsive features. Another option is to equalise the feature spectrum to the EEG spectrum by convolving the sparse,
 100 impulsive feature with an appropriately smooth kernel. However, these approaches remain to be investigated in detail,
 101 and typically encoding models are used with impulsive features.

102 2.2. Encoding modelling

103 Encoding modelling can be used to study the spatio-temporal properties of the response: the EEG signal in each
 104 channel is predicted from the speech features via a weighted sum of the different speech features and their time-

105 shifted versions. These weights form a temporal response function (TRF) for each speech feature and each EEG
106 channel. A TRF consists of the estimated weights at the different time-shifts of the speech feature, reflecting how the
107 EEG response is modulated by the stimulus at different time-shifts. The encoding model can deal to some extent with
108 autocorrelation of the speech stimulus. Autocorrelation denotes that speech is correlated with itself at different time
109 lags. Because the computation of the encoding model uses this autocorrelation, it can prevent the smearing over time
110 of the TRF (for more details, see Crosse et al., 2016b).

111 Panel B of figure 1 schematically presents the encoding modelling process. Note that for the encoding modelling,
112 the time-shifting occurs in the opposite direction than for decoding modelling. Each EEG channel is considered
113 separately, so encoding models cannot reduce noise in the EEG signal by combining information across channels.
114 The advantage of this approach is that the TRF weights are activation patterns and not extraction filters. Activation
115 patterns have large weights for the EEG channels containing a lot of response information and can therefore be
116 interpreted. The temporal aspects of these TRF weights are particularly interesting to investigate the neural response
117 latency. Depending on the considered brain area, different bottom-up neural response latencies can be expected:
118 about 5-10 ms for auditory processing in the upper brainstem and at least 12-30 ms for processes in the primary
119 auditory cortex (Tichko and Skoe, 2017; Brugge et al., 2009). Higher-order cortical processes that modulate the neural
120 response, like attention and interpretation of the speech, occur with delays of 200 ms or more (for a review, see Martin
121 et al., 2008). Similarly to decoding modelling, a measure of neural tracking can be extracted by correlating the actual
122 measured EEG responses with the predicted EEG responses. This correlation is called the prediction accuracy and is
123 a measure of neural tracking, i.e., the better the brain tracks the speech, the higher the prediction accuracy.

124 Compared to decoding modelling, encoding modelling results in a lower magnitude of neural tracking. The reasons
125 are twofold. Firstly, in encoding modelling, the actual EEG is correlated with the EEG predicted from the speech
126 features. However, the predicted EEG is a gross simplification of the content of the actual EEG signals, which contain
127 responses to the speech together with a plethora of non-speech-related EEG activity and noise. Secondly, encoding
128 modelling cannot use across-channel information to reduce noise in the EEG signals (Das et al., 2019). However,
129 a lower magnitude of neural tracking does not necessarily mean that the encoding model is less valid or reliable
130 than the decoding model. The only issue is that the metric to assess the quality of the model is noisy and has lower
131 values.

132 For each channel, the TRF can be interpreted as the impulse response of the measured auditory system: the information
133 in the input stimulus, i.e., the speech feature, is transformed with this impulse response to produce the output response,
134 i.e., the preprocessed EEG. Please note that the impulse response depends on the preprocessing (further described in
135 section 2.4). The channel-specific TRFs are noisy and therefore often averaged over a selection of EEG channels and
136 subjects. Based on the time-shifts that receive large weights for many of the EEG channels/subjects, the dominant
137 latencies of the response can be derived. These latencies (or delays) can then be used to estimate which stages of
138 neural processing along the auditory pathway contribute to the response. The spatial properties of the response can be

139 further investigated from the distribution of the magnitude of TRF weights over the scalp. This information is usually
140 visualised on a topographic map. Examples of TRFs and topographic maps are shown in figure 4, which is discussed
141 below. The TRFs and the corresponding topographic maps are similar to ERPs with the advantage that they can be
142 computed using a continuous signal. Moreover, the prediction accuracies of the encoding model can also be visualised
143 on a topographic map. Note that such topographic maps only divulge spatial information on scalp level, where the
144 electrodes were located. To study the actual sources of the neural responses within the head, the information from
145 electrode space should be transformed to neural source space (e.g. Brodbeck et al., 2018c).

146 Although the encoding modelling approach can deal to some extent with the autocorrelation of the stimulus, it can
147 still be problematic for certain features, like the fundamental frequency (further discussed in section 4).

148 *2.3. Algorithms to calculate decoding models and encoding models*

149 Different algorithms exist to acquire decoding and encoding models. In essence, the algorithms have the same goal:
150 minimizing the error between the reconstructed or predicted signal with the actual signal. However, due to the noisy
151 nature of EEG signals and the fact that only a small fraction of the EEG signal is auditory stimulus related, this
152 question is ill-posed meaning that multiple solutions are possible. Therefore, these algorithms may yield slightly
153 different outcomes as they might rely on different priors. In this overview, we will focus on two algorithms: ridge
154 regression (Machens et al., 2004) and boosting (David et al., 2007). Both algorithms are supported by a dedicated
155 toolbox, respectively the mTRF toolbox (Crosse et al., 2016b) and Eelbrain (Brodbeck et al., 2021b).

156 The ridge regression algorithm minimises the mean-squared error between the predicted or reconstructed signal and
157 the actual signal. It relies on the inverse of the autocorrelation matrix of the time-shifted input. The input depends on
158 the considered model, i.e. EEG for a decoding model and speech features for the encoding model. Taking the inverse
159 of this autocorrelation matrix is ill-posed as the rows of the matrix are mutually dependent, because the time-shifted
160 inputs are dependent. To solve this, a ridge parameter is added to each diagonal element of the autocorrelation matrix.
161 Using a cross-validation approach (described in 2.4), the ridge parameter can be determined. To do so, the measure of
162 neural tracking is calculated for multiple values of the ridge parameter (for example ranging from 10^{-2} to 10^3 in steps
163 of the powers of 10). The ridge parameter with the highest accuracy is selected and used for the analysis.

164 Like ridge regression, the boosting algorithm minimises the error between the predicted or reconstructed signal and
165 the actual signal, and aims for a maximally sparse solution. In contrast to ridge regression, which has a closed-form
166 solution, boosting relies on an iterative approach to determine the model's weights. Initially, all weights are set to
167 zero. Subsequently, each weight is changed by a specific value. The error between the predicted or reconstructed
168 signal and the actual signal is calculated for every change. The weight which resulted in the smallest error is selected.
169 Then, these weights are used to start the next iteration. These iterations are performed until a stopping criterion is
170 met, e.g. the error stops decreasing, or the number of iterations exceeds the limit.

171 Although both algorithms lead to a similar results, the outcome has some apparent differences (Kulasingham and
172 Simon, 2022). Ridge regression results in smooth solutions which are more spread across time and channels while
173 boosting results in a more sparse solution which is better defined in time and space. On the other hand, as boosting
174 relies on this iterative approach, it is more computationally expensive than ridge regression.

175 *2.4. Preprocessing and model evaluation*

176 Preprocessing of the EEG and the speech features affects the results and interpretation of the model. Especially
177 the interpretation of the patterns in the TRF should be made carefully with respect to the preprocessing characteris-
178 tics.

179 The filtering method is a crucial aspect to consider when preprocessing the EEG data. Every filter has specific
180 characteristics which affect the impulse response of the system. Two important filter characteristics should always be
181 considered: the causality and the filter's phase response. The causality relates to which data points the filter uses. For
182 causal filters, only past data points can affect the output of a specific data point, while for acausal filters, the output can
183 be affected by past and future data points. The phase of the filter is also important as it denotes the delay introduced
184 by the filter.

185 Especially when interpreting the TRF, the filter characteristics should be considered. Firstly, we want to emphasise
186 that using a causal filter makes more sense: the stimulus evokes a response in the brain, so only past data points can
187 influence the output. However, a causal filter cannot be zero-phase, and therefore, introduces a delay which should be
188 accounted for. An exception is when the EEG and stimulus features are filtered the same way; the delay will affect
189 both in the same way, and thus the time delay and time-locking between EEG and stimulus is preserved. It should
190 be emphasised that all filter characteristics should reported for a study, as pointed out by de Cheveigné and Nelken
191 (2019). The best practice is to filter the stimulus and the EEG similarly. Additionally, the EEG and stimulus spectrum
192 will become similar, leading to higher prediction or reconstruction accuracies. Regarding decoding modelling, the
193 filter causality and phase are less critical as the weights of the decoding model are not interpreted.

194 Another critical preprocessing step is the referencing of the EEG signals. Different choices can be made: the central
195 electrode Cz, mastoids or a common average. The reference choice does not affect the overall magnitude of the
196 reconstruction or prediction accuracy, but it does affect the weights of the model. This aspect is not essential for
197 decoding models as the weights cannot be interpreted. For encoding models, the TRF and the distribution of prediction
198 accuracies across the EEG channels are affected by the reference choice. If the EEG signals are referenced to one or
199 a combination of couple of electrodes such as the two mastoids, the obtained referenced EEG favours specific brain
200 regions. If this is not wanted, the common average referencing is a good choice which allows a broader picture of the
201 neural activity. However, when making claims about neural activity in dedicated regions in the brain, using source
202 localisation techniques is a better solution than deriving conclusions based on the TRF.

203 Preprocessing of EEG also incorporates artefact removal. The above-discussed linear models investigate the time-
204 locked relationship between stimulus features and the EEG responses. As artefacts due to eye blinks or movement
205 are not strictly time-locked to the speech, the models can cope with these artefacts if sufficient data is provided.
206 Nevertheless, artefact suppression should always be considered. Multiple options are possible: multi-channel Wiener
207 filtering (Somers et al., 2018), independent component analysis, denoising source separation (Särelä et al., 2005),
208 multiway canonical correlation analysis (de Cheveigné et al., 2019), etc. These techniques are useful to suppress
209 various artefacts in the EEG signal.

210 After preprocessing the data, the linear models can be estimated. We want to emphasise the necessity of cross-
211 validation to estimate and evaluate the models. If the model is estimated and evaluated on the same data, the results
212 are likely to be biased: inflated reconstruction or prediction accuracies and distorted TRF patterns. For example, more
213 peaks might be seen in the pattern because the model learnt the noise in the data. Altogether, estimation and evaluation
214 of the model on the same data leads to unreliable results specific to the used data and may not generalise well to new,
215 unseen data. This can be avoided by using the cross-validation technique (Crosse et al., 2021). This technique relies
216 on a training set, i.e., part of the data used for model estimation and a testing set, i.e., another part of the data, unseen
217 during the model estimation, to evaluate the model on (visualised in Figure 1). The cross-validation technique divides
218 the data into n folds of equal length. Subsequently, the model is estimated on $n - 1$ folds and evaluated on the left-out
219 fold. This is repeated until all folds have been used to evaluate the model. The TRFs and prediction or reconstruction
220 accuracies obtained by respectively model estimation and evaluation are then averaged across the different folds. This
221 technique allows identifying a robust model that generalises to new, unseen data.

222 2.5. *Other methods to extract neural tracking*

223 Decoding and encoding models imply a directionality: either the stimulus feature is reconstructed or the EEG re-
224 sponses are predicted. Another linear approach is canonical component analysis (CCA), which operates bidirection-
225 ally. CCA transforms both EEG and stimulus, so they are maximally correlated, thereby combining the advantages
226 of the decoding and encoding models. Stimulus dimensions irrelevant for measurable responses are removed, as are
227 EEG dimensions irrelevant for auditory perception. Although CCA is a flexible tool that can discover more complex
228 relations than a simple encoding or decoding model, it has more parameters and with this a higher risk of overfit-
229 ting. Therefore, an appropriate cross-validation strategy is needed, or one has to use dimensionality reduction or
230 regularisation. Furthermore, as with all techniques based on least-squared minimisation, it is prone to outliers. CCA
231 decomposes the signals into multiple components. Although this can ease the interpretation of the results, it should
232 be done with care. CCA orders the components based on their correlation, however a high correlation does not guar-
233 antee a physiological interpretation. For example, components that have lowpass or narrowband filtering can have
234 very high correlations. Finally, there is no exact match between the components and the actual neural sources which
235 can complicate the interpretation of the underlying neural basis of the component. This technique has been used in
236 de Cheveigné et al. (2018) and O’Sullivan et al. (2021)

237 Another linear analysis method is a cross-correlation (Kong et al., 2014; Aiken and Picton, 2008; Petersen et al., 2016;
238 Aljarboa et al., 2022). Here, the cross-correlation is computed between the EEG channels and the speech features.
239 The cross-correlation is computationally inexpensive and can give some insight into the neural responses. Similar
240 to the encoding model, it cannot integrate multiple channels. An disadvantage compared to encoding and decoding
241 modelling, is that the cross-correlation is more sensitive to the autocorrelation of the stimulus, leading to patterns that
242 are smeared out over time. Another disadvantage is that cross-correlation cannot be applied on an unseen test dataset.
243 A comparison of the TRF to the cross-correlation is nicely shown in Crosse et al. (2016b).

244 Non-linear techniques have been explored to overcome the inherent limitations of linear models. Mutual information
245 is a metric that uses information theory, which captures the shared information between the EEG responses and the
246 stimulus expressed in the unit ‘bits’. Therefore, higher mutual information indicates higher neural tracking. Because
247 it does not make explicit assumptions about the relationship between the stimulus and EEG, it can capture non-
248 linear aspects. Moreover, mutual information can be calculated between the EEG and the time-lagged stimulus to
249 understand how the mutual information metric behaves for multiple time lags. This results in a pattern similar to the
250 TRF. However, an important consideration is that autocorrelation of the stimulus is not taken into account. Note that
251 the mutual information method does not need to be applied between stimulus and EEG necessarily, but can also be
252 applied between different EEG signals to get additional insights. This technique has been used by Gross et al. (2014),
253 Zan et al. (2020), and Kaufeld et al. (2020).

254 Another measure of neural tracking can be obtained with a match-mismatch paradigm. In this case, a model is trained
255 to classify whether a given EEG segment is matched or not with a given stimulus segment. This can be done for a
256 single stimulus segment or N segments, in which case the model classifies which of the N segments is matched with
257 the EEG. In the case of decoding or encoding linear models, the stream with the highest prediction or reconstruction
258 accuracy can be identified as the matched stream. The accuracy, i.e., the percentage of correctly identified matched
259 speech streams, is a metric of neural tracking. Please note that auditory attention decoding (AAD) relies on the same
260 principle. Instead of presenting a match and mismatch speech stream, the attended and unattended speech streams are
261 given as input to the model (e.g. Fuglsang et al., 2020; Das et al., 2018; Deckers et al., 2018; O’Sullivan et al., 2015;
262 De Cheveigné et al., 2021).

263 This paradigm can also be solved in a non-linear fashion with neural networks (e.g. Accou et al., 2021; Monesi
264 et al., 2021; Bollens et al., 2022). Accou et al. (2021) showed that the accuracy of a neural network solving a
265 match-mismatch task could be used to estimate the speech reception threshold. Therefore, this neural network allows
266 for an evaluation of speech understanding based on the EEG responses. Although neural networks can model non-
267 linear relationships between the EEG and stimulus, they give limited insight into how the network handles the EEG
268 responses. Therefore, evaluating different levels of the auditory system becomes more challenging as it is difficult to
269 tell which information and how the neural network uses it.

270 Although other methods can quantify neural tracking, in the continuation of this overview, we focus on decoding and
271 encoding models.

272 **3. The stimulus feature**

273 The stimulus feature is derived from the presented speech and reflects how a particular speech characteristic varies over
274 time. Many stimulus features can be used, ranging from low-level acoustic characteristics (e.g. the acoustic envelope)
275 to high-level linguistic information (e.g. word surprisal). This flexibility makes the neural tracking framework highly
276 versatile and allows for evaluating multiple levels of the auditory system. It also underlies one of the most prominent
277 advantages of the framework: a single EEG measurement can be analysed with various features of the stimulus and
278 provides information on a range of auditory/language processes. Note that the feature choice is arbitrary, and thus
279 different features will reflect the different stages of the auditory pathway. In this manuscript, we focus on f0 tracking,
280 envelope tracking and linguistic tracking to target different stages in the auditory pathway. Other features are possible
281 and have been investigated in other studies, e.g. word category (Brennan and Hale, 2019), acoustic onsets (Brodbeck
282 et al., 2018a), the spectrogram (Di Liberto et al., 2015), phonetic features (Di Liberto et al., 2015), etc.

283 In the following sections, we will illustrate the use of neural tracking measures using three prominent (groups of)
284 stimulus features corresponding to three types of neural tracking analyses. We discuss these following the hierarchical
285 organisation of the auditory pathway: starting with auditory processing of the fundamental frequency (f0, section 4),
286 which happens mainly in the subcortical stages of the auditory pathway, then moving on to envelope processing
287 (section 5) which happens in the auditory cortex and ending with linguistic processing (section 6) which happens
288 in the language network of the brain. We focus on how these stimulus features can be used to investigate different
289 aspects of speech processing and different parts of the auditory pathway. Moreover, we provide example results and
290 review findings from relevant studies, including how the responses relate to important clinical measures like hearing
291 thresholds and speech perception. Please note that it is out of the scope of this paper to comprehensively review all
292 stimulus features and related analyses that have been published.

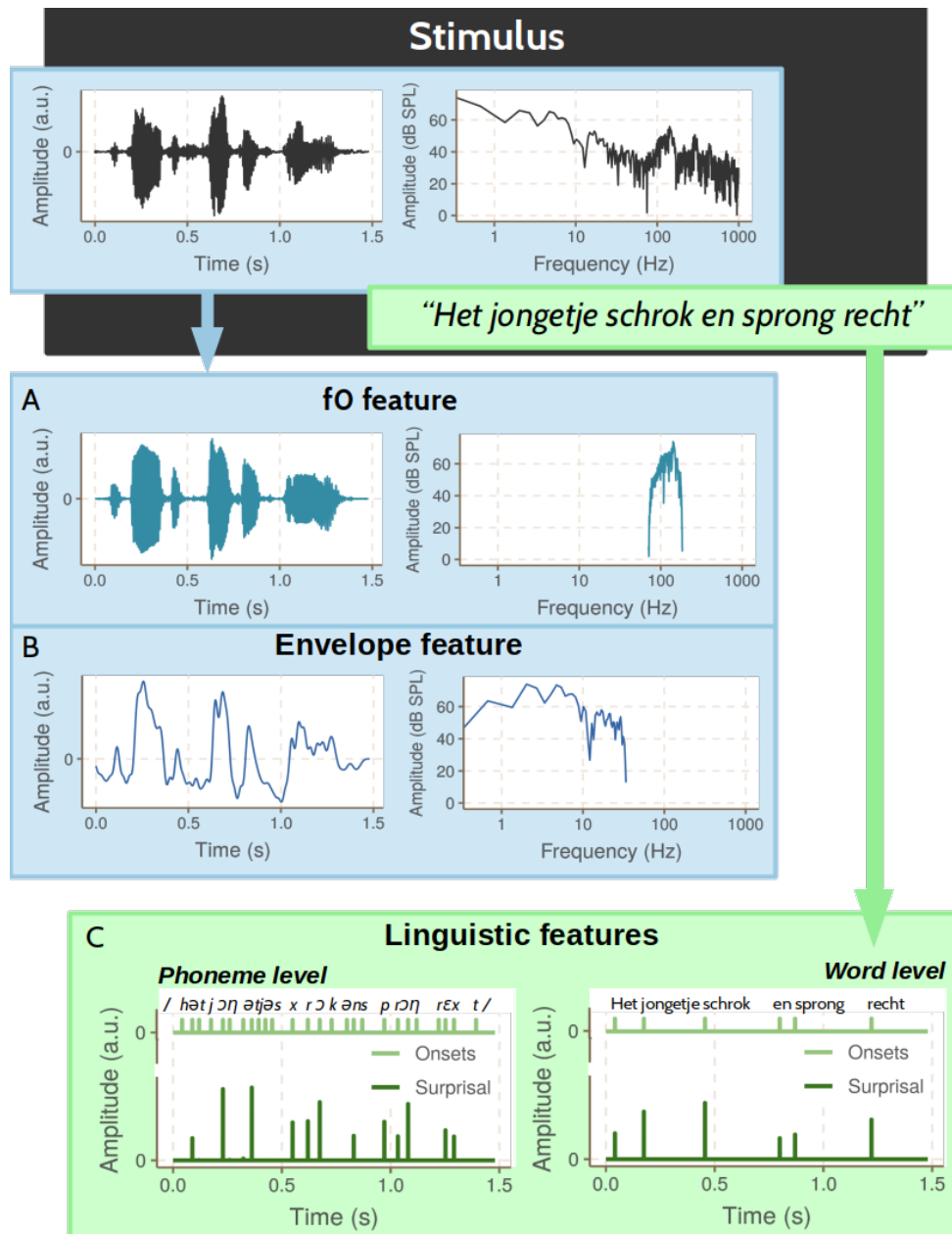


Figure 2: Example of stimulus and derived features for an example sentence by a male speaker. The f0 (panel A) and envelope feature (panel B) are derived from the stimulus waveform, whereas linguistic features (panel C) are derived from the stimulus transcription. The f0 and envelope features concern different spectral ranges, with the envelope focusing on low frequencies (< 50 Hz) and the f0 focusing on higher frequencies (~ 85 – 300 Hz). Linguistic features can focus on different segmentation levels, including phoneme level and word level. Panel C visualises an example onset and surprisal feature for each level.

293 *3.1. Inter-feature correlation and feature evaluation*

294 All the speech features are derived from the same speech stimulus. This leads to a high inter-feature correlation.
295 For example, at every impulse of a linguistic speech feature, there is a word or phoneme onset, which tends to be
296 associated with a high burst of acoustic energy. In Figure 3, we visualised the inter-feature correlation (panel A).
297 Moreover, we created a TRF to predict the envelope based on the remaining stimulus features (using the boosting
298 algorithm; panel B). Not surprisingly, the envelope feature can be predicted from the other speech features. Even
299 some aspects of linguistic features (such as word surprisal) explain some variance in the envelope feature. This
300 inter-feature correlation can affect the model performance and complicate the interpretation of the results.

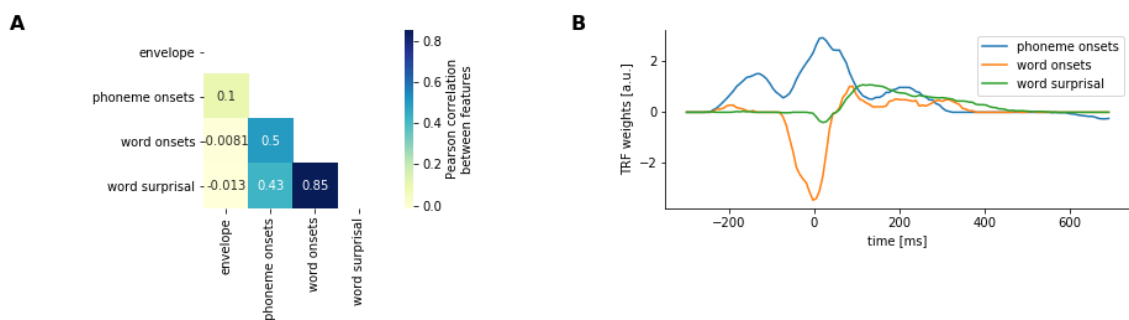


Figure 3: *Illustrative example of the inter-feature correlation for the stimulus used in Figure 2.* Panel A shows the correlation between the different features. Panel B shows the TRF for phoneme onsets, word onsets and word surprisal when trying to predict the envelope of the stimulus.

301 This inter-feature correlation can bias the interpretation of results, which is important to consider when investigating
302 the tracking of a linguistic feature using a model with only the linguistic feature of interest. When significant neural
303 tracking is observed, it cannot be attributed to solely the brain tracking linguistic aspects of speech due to this inter-
304 feature correlation (Daube et al., 2019).

305 To overcome this issue, a good approach is to assess the feature's contribution to the model performance. To do
306 so, the features of interest must be defined together with control features, i.e., features that are not of interest in the
307 study's goal but are correlated with the feature of interest. Here, we discuss three approaches: subtracting correlations,
308 residual fitting and feature shuffling. For the first method, two models are created, one with and one without the feature
309 of interest. If the reconstruction or prediction accuracies significantly increase when the feature of interest is included,
310 the feature of interest contributes unique information to the model. This method has been applied in previous studies
311 (e.g. Di Liberto et al., 2015; Brodbeck et al., 2018a; Gillis et al., 2021b). The disadvantage of this method is that it is
312 too conservative. Only the unique contribution of the feature of interest compared to the control features is captured.
313 As shown on Figure 3, some control features, like the envelope, can contain information which is also captured by
314 word surprisal. Therefore, the linguistic information common to both features will be attributed to the envelope.
315 Another disadvantage is that the degrees of freedom change between the two models. This is especially important

316 to consider when the cross-validation technique is not applied. If this approach is used with ridge regression, Crosse
317 et al. (2021) suggests using banded-ridge regression as different features might require different regularisation.

318 Another approach is to look at the model's residuals with the control features. In this case, first, a model is created
319 using the control features. Subsequently, the predicted EEG or reconstructed envelope is subtracted from the actual
320 signal, creating the models' residuals. Then a new model is created using the features of interest and these residuals.
321 If a significant prediction or reconstruction accuracy is achieved, the feature of interest explains variance in the EEG
322 responses, which is not explained by the control features. Similar to the subtracting correlation approach, this method
323 is very conservative as it only considers the unique contribution of the model.

324 The last approach is shuffling of the features. This approach is used in studies by Broderick et al. (2018, 2021). Now
325 two models include the control features combined with either the feature of interest or a shuffled version of the feature
326 of interest. If the model with the feature of interest performs significantly better than the model with its shuffled
327 version, the feature contributes significantly to the model. Note that shuffling of the feature should be done with care.
328 For impulsive features, shuffling the feature can be done by preserving the timing of the pulses, but changing the
329 amplitudes. For continuous features, the shuffled feature can be created by filtering noise with the same spectrum
330 as the feature of interest. In this method, the number of features is kept constant to investigate the feature's added
331 value. However, a disadvantage is that the inter-feature correlation is not preserved. For example, returning to Figure
332 3, if the feature of interest is word surprisal, the correlation between word surprisal and envelope gets lost, i.e., the
333 envelope captures no effect of the shuffled word surprisal. The loss of the inter-feature correlation might affect the
334 model performance, mainly if the cross-validation technique is not applied.

335 Although each approach has disadvantages, they are all valid and used in different studies. However, good control
336 conditions should always be considered to evaluate whether or not a feature captures the desired effect. For example,
337 for a linguistic feature, this might be a foreign language, vocoded speech or time reversed speech. However, each
338 of these control conditions also has its drawbacks. Depending on the choice of a foreign language, the speaker
339 and language structure can vary, which affects neural tracking. Although vocoded speech can preserve the speech
340 envelope, other speech cues are lost. Time reversed speech has a limitation that the onsets become unnatural, i.e.
341 acoustic boundaries occur at the end of the sound instead of the beginning. An overarching disadvantage is that the
342 listener's attention may drift over the course of the control stimulus. As the stimulus is not understandable, it is more
343 challenging to listen attentively to the stimulus.

344 **4. Neural tracking of the f0**

345 Neural tracking of the fundamental frequency of the voice, or f0-tracking, is used to investigate how the f0 is rep-
346 resented in the brain activity (Forte et al., 2017; Etard et al., 2019; Van Canneyt et al., 2021c). The f0 is a periodic
347 modulation in the speech signal generated by vocal fold vibration during speech production. It is related to the per-

348 ception of pitch. The f_0 of adult speakers typically ranges from 85 to 300 Hz, with male and female voices situated
349 respectively at the lower and higher ends of the range. The f_0 is an essential characteristic of the human voice, and it is
350 vital to convey intonation and emotion. However, proper perception of the f_0 is not required for speech intelligibility
351 (e.g. cochlear implant listeners). Nevertheless, f_0 -tracking can provide information on the quality of fine temporal
352 processing in the early stages of the auditory pathway, which is the foundation for proper speech processing in the
353 brain.

354 Temporal processing of the f_0 in the human auditory system happens through the synchronisation of the activity of
355 the neurons to the f_0 modulations, i.e. phase-locking. Due to the relatively high frequency of the f_0 modulations, this
356 phase-locking occurs mainly in peripheral and subcortical stages of the auditory pathway, up to the upper brainstem.
357 Neurons at cortical stages have poor phase-locking above 100 Hz and are therefore less likely to contribute to f_0 -
358 tracking (Joris et al., 2004). However, it has been shown that early cortical contributions to f_0 -tracking responses (and
359 FFRs) can occur for low-frequency stimuli (85-100 Hz, e.g. low male voices) (Coffey et al., 2016, 2017; Van Canneyt
360 et al., 2021c).

361 F_0 tracking analysis requires an f_0 feature that represents the f_0 modulations in the presented speech. The f_0 feature
362 can be extracted from the speech stimulus in various ways. A simple yet effective way is to band-pass filter the
363 stimulus in the range of the f_0 (Etard et al., 2019; Van Canneyt et al., 2021c). An example of this type of feature is
364 provided in panel A of figure 2. More complicated and computationally expensive techniques have also been explored,
365 including empirical mode decomposition (Etard et al., 2019; Forte et al., 2017) and auditory modelling (Van Canneyt
366 et al., 2021b). Constructing an f_0 feature that approximates the expected neural response using auditory modelling has
367 proven particularly effective, nearly doubling the reconstruction accuracies obtained with the neural tracking analysis
368 (Van Canneyt et al., 2021b). This is likely explained by the fact that the auditory model is more physiologically
369 valid. Moreover, the auditory model simulates the contribution of the higher harmonics to the f_0 response. Because
370 the neural response is also driven by the harmonics and not just by the f_0 , the level of measured f_0 tracking will
371 increase.

372 Section 1 of figure 4 shows the results of a typical encoding modelling analysis for f_0 -tracking. These results were
373 obtained by filtering the EEG responses between 75 and 175 Hz and referencing to the average EEG response (32
374 participants; Van Canneyt et al. (more details regarding the preprocessing are described in 2021c)). The data set
375 used for this visualisation (and all others in figure 4) contained 64-channel EEG data from 32 young normal-hearing
376 subjects measured in response to male-narrated speech (Accou et al. (dataset from 2021)). Panel A shows the mean
377 TRF across subjects for the channel selection indicated in pink on panel B. The TRF for each subject is plotted as
378 well to indicate the variance. The TRFs in this example are absolute value of the TRFs of the stimulus and the Hilbert
379 transformed EEG. This technique aids with interpretation as the brainstem response may occur at a phase that is
380 different from that of the f_0 (for more information, see Van Canneyt et al. (2021c); Forte et al. (2017)). The TRF
381 pattern indicates that the activity in the auditory system (\sim EEG) best reflects the f_0 information (\sim the feature) at

382 a latency of about 10-25 ms. Panel B of figure 4 presents an example f0 tracking topographic map with common-
383 average rereferencing at 15 ms latency. The topographic map indicates strong response activity in the centre of the
384 scalp and across the back of the head. The temporal and spatial response patterns are consistent with dominant f0-
385 related activity in the upper brainstem and early cortical regions. Saiz-Alía and Reichenbach (2020) has performed
386 detailed computational modelling of the subcortical sources of the f0 tracking response, demonstrating important
387 contributions from the cochlear nuclei and the inferior colliculus. Van Canneyt et al. (2021c) argues for additional
388 contributions from the right primary auditory cortex for f0 tracking of low-frequency voices.

389 Although f0-tracking was only recently developed, the technique has led to several interesting findings. Forte et al.
390 (2017) and Etard et al. (2019) have demonstrated that the f0 tracking response holds information on selective attention,
391 possibly indicating that neural mechanisms for attention influence the brainstem. Kulasingham et al. (2020) and
392 Van Canneyt et al. (2021a) have investigated how the age of the listener impacts f0 tracking. Kulasingham et al.
393 (2020) found no age effects using MEG, which is most sensitive to cortical sources. In contrast, Van Canneyt et al.
394 (2021a) found a significant reduction in response strength with advancing age using EEG (which is more sensitive
395 to subcortical sources than MEG and will capture both cortical and subcortical sources). This observation is in line
396 with an age-related decrease in the phase-locking ability of the subcortical (and early cortical) auditory system. Van
397 Canneyt et al. (2021a) also studied the effect of hearing loss and found increased f0-tracking responses in participants
398 with hearing impairment compared to age-matched controls. The response enhancement was due to additional cortical
399 activity phase-locked to the f0 (with a latency of ~40 ms), likely compensating for the reduced quality of bottom-up
400 auditory input due to diminished peripheral auditory sensitivity. Moreover, the amount of additional compensatory
401 cortical activity was significantly related to the pure tone average (PTA) hearing loss of the participant. As such, a
402 significant relationship exists between the degree of hearing loss of an individual and the strength of their f0 tracking
403 response.

404 At the moment, f0-tracking also has some limitations, which future advances may mitigate. One of the main issues
405 is auto-correlative smearing in TRFs and topographic maps because the f0 stays relatively steady over multiple f0
406 periods. This periodic smearing over latencies can be somewhat mitigated with Hilbert-transformed TRFs, which
407 disregard phase information. However, TRF and topographic map interpretation are still limited to the most dominant
408 peaks (see Van Canneyt et al. (2021c) for more details). A second limitation is that the f0 is only present in speech
409 during voiced sounds (~ 50-60 % of the time) and not during unvoiced sounds (~ 40 % of the time), including
410 silences. During analysis, these unvoiced sections in the speech stimulus (and corresponding sections in the EEG)
411 are disregarded. As a result, only about half of the measured data can be used to analyse f0-tracking, increasing
412 the required measurement time. Another limitation is that the f0 tracking response is reduced for voices with higher
413 and more variable f0, leading to weak and often non-significant responses for typical female voices. This occurs
414 because neural phase-locking ability is decreased for higher and more variable f0s, especially for cortical sources.
415 As such, the stimulus choice has a large impact on the f0 tracking response. A final limitation is that f0-tracking

416 requires careful interpretation: f0-tracking reflects the capability of the auditory system to phase-lock to the f0, but
417 it does not reflect the ability of a person to perceive pitch or speech in general. Neural tracking analyses with other
418 features help complete the picture. Lastly, the use of the f0 feature is still in its infancy with only a limited number of
419 published studies. Nevertheless, it seems to be a promising approach to evaluate neural responses at the level of the
420 brainstem.

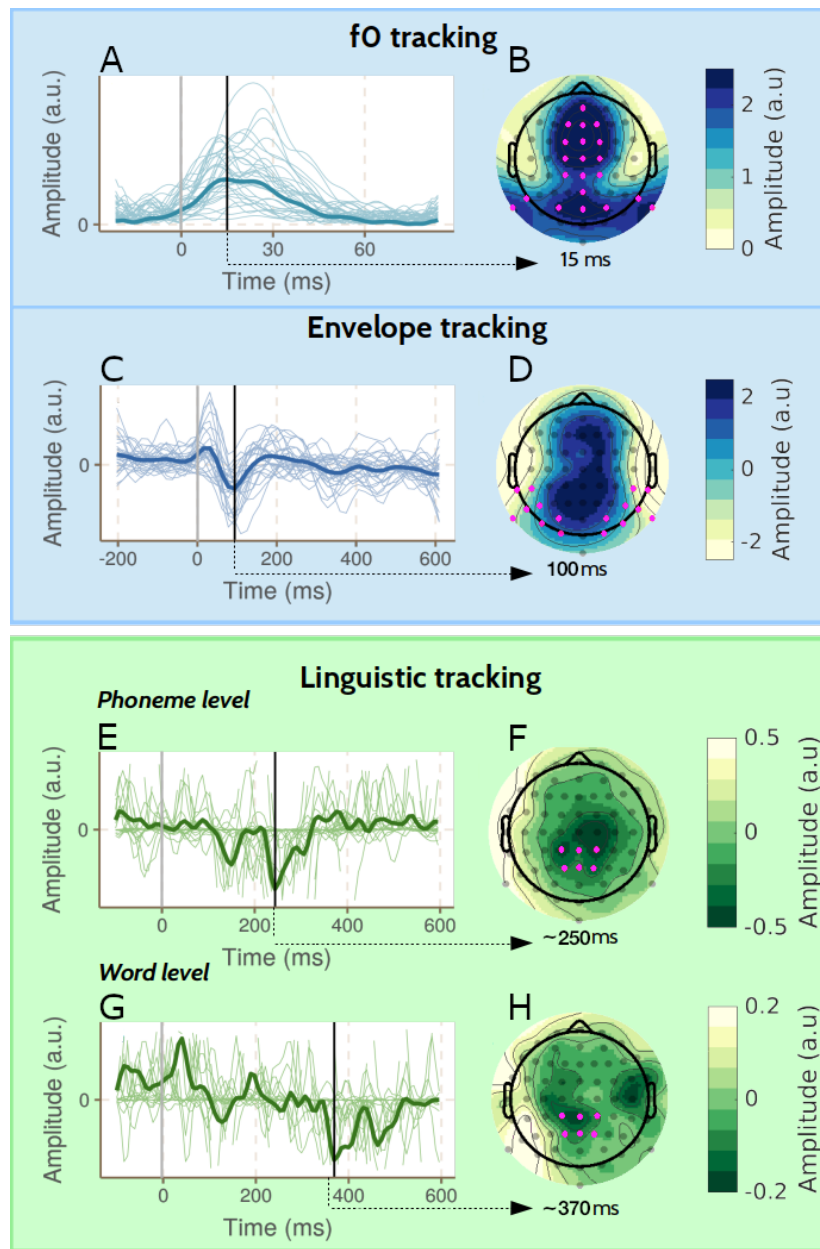


Figure 4: *Example of encoding modelling results: TRFs and topographic maps.* The figure is divided into three sections on f0-tracking, envelope tracking and linguistic tracking, respectively. For each type of tracking, an example mean TRF (+ individual TRFs) is presented (panel A, C, E and G), together with a corresponding topographic map at an important latency (panel B, D, F and H). These TRFs are estimated for the same participants and the same speech material of around 15 minutes long. The channels indicated with pink on the topographic map represent the channel selection used to obtain the corresponding TRF. Note the drastically different time scales in the TRFs, reflecting the presence of neural activity at different latencies for each feature. In panel A, the TRFs are visualized of the stimulus and the Hilbert transformed EEG, similarly to Van Canneyt et al. (2021c) and Forte et al. (2017).

421 **5. Neural tracking of the speech envelope**

422 The speech envelope consists of slow-varying modulations (< 50 Hz) in the speech signal. It contains acoustic
423 temporal information (Rosen, 1992) but also reflects phonemes, syllables and word transitions (Peelle and Davis,
424 2012). Moreover, it also correlates with the area of the mouth opening during articulation (Chandrasekaran et al.,
425 2009). Therefore it is not surprising that research indicates that the envelope is an essential acoustic cue for speech
426 intelligibility (Shannon et al., 1995; Drullman et al., 1994a,b).

427 Measuring envelope tracking can be used to analyse the neural encoding of the speech envelope during speech percep-
428 tion (Ding and Simon, 2012a; O'Sullivan et al., 2015; Vanthornhout et al., 2018). From animal studies (Wang et al.,
429 2008) and human studies with electrocochleography (ECoG), it is known that the speech envelope is processed in the
430 primary auditory cortex, specifically in Heschl's Gyrus (Nourski et al., 2009). A growing body of evidence demon-
431 strates that envelope tracking is a requirement for speech understanding. Multiple studies show that neural tracking of
432 the speech envelope is strongly correlated with behaviourally measured speech intelligibility (e.g. Ding et al. (2014);
433 Vanthornhout et al. (2018); Lesenfants et al. (2019); Iotzov and Parra (2019); Verschueren et al. (2021)). As a specific
434 example, Vanthornhout et al. (2018) found a significant correlation of 0.69 between the speech reception threshold
435 (SRT) estimated based on envelope tracking and the SRT measured with behavioural speech audiometry.

436 Although the broadband envelope can be used, it is also possible to study the neural tracking of specific frequency
437 bands of the envelope. Envelope tracking responses are most commonly investigated in the delta band (0.5-4 Hz),
438 theta band (4-8 Hz) and gamma band (> 30 Hz) (Ding and Simon, 2013; Verschueren et al., 2021; Molinaro and
439 Lizarazu, 2017). The lower envelope frequencies are often the main interest as they correspond with word onsets
440 and the syllabic rate of the speech, which is hypothesised to be crucial for speech intelligibility. Higher envelope
441 frequencies are typically related to the onsets of phonemes. Some studies suggest that speech intelligibility, i.e., how
442 well a person can understand speech, is specifically related to the theta band (4-8 Hz) and not the delta band (1-4 Hz)
443 (Ding and Simon, 2013; Peelle et al., 2013). Other studies indicate the opposite (Verschueren et al., 2021; Molinaro
444 and Lizarazu, 2017; Di Liberto et al., 2018). In our opinion, the outcome may depend on the speech material. The
445 syllabic rate is often very close to 4 Hz, so envelope tracking to a slow speaker could be more dominant in the delta
446 band, while envelope tracking to a fast speaker could be more dominant in the theta band.

447 Envelope tracking responses can be analysed using a decoding or encoding model, or using a strategy that combines
448 a decoding and encoding approach in one model. In any case, the model requires an envelope feature extracted from
449 the stimulus waveform. In essence, the envelope reflects the modulation of the signal and can easily be obtained by
450 taking the absolute value of the Hilbert transform. Although this is a prevalent method, it is not the best choice as it
451 disregards auditory processing. Two important aspects of auditory processing need to be taken into account to better
452 approximate human envelope processing. First, the stimulus should be split into frequency bands before the actual
453 envelope extraction process to mimic how the basilar membrane in the cochlea divides a sound stimulus into different

454 auditory filters. Second, the compression and non-linear behaviour of the auditory system should be accounted for. To
455 incorporate these factors in the envelope extraction process, complex computational models of the auditory periphery
456 can be used (Yang et al., 2015; Bruce et al., 2018). However, Biesmans et al. (2017) evaluated various extraction
457 methods in an auditory attention decoding (AAD) paradigm and proposed a simplified approach. They found that a
458 combination of a gammatone filterbank, which simulates the auditory filters on the basilar membrane, followed by a
459 power law to account for compression and non-linearity in the auditory system, performed equally well as the more
460 complex and computationally expensive auditory models. An example envelope feature obtained using this technique
461 is provided in panel B of figure 2.

462 The results of a typical encoding modelling analysis using ridge regression for envelope tracking are visualised in
463 section 2 of figure 4. These results were obtained by highpass filtering the EEG responses above 0.5 Hz and refer-
464 encing to the average EEG response (32 participants; Vanthornhout et al. (more details regarding the preprocessing
465 are described in 2019)). Panel C presents the mean TRF, averaged over subjects and a channel selection (indicated
466 in pink on panel D). The TRFs of the individual subjects are visualised with a thin line to indicate the variance. The
467 TRF displays three distinct peaks. The P1 peak (50 ms), the N1 peak (93 ms) and the P2 peak (170 ms). This typical
468 P1-N1-P2 complex is also found in AEP studies with impulse-like stimuli and can thus be used to infer the neural
469 source of the peaks. The P1 peak originates in Heschl's Gyrus, and the N1 peak originates in the Superior Temporal
470 Gyrus (O'Sullivan et al., 2019b; Steinschneider et al., 2011). The origin of the P2 peak is less clear but is probably
471 in the (higher) auditory cortex (Godey et al., 2001). The topographic map shows negative weights for the temporal
472 channels and positive weights for the central channels. This distribution is an indication of a dipole located near the
473 auditory cortex. Without analyses in source space, the exact location is difficult to pinpoint.

474 Over the past decade, envelope tracking has been used to study, among others, how cortical speech processing is
475 affected by individual factors like age and hearing status. Decruy et al. (2019) and Brodbeck et al. (2018b) found
476 stronger envelope tracking for older participants compared to younger participants, even though older adults typically
477 have more difficulty understanding speech. Similarly, Decruy et al. (2020b) and Fuglsang et al. (2020) found increased
478 envelope tracking for hearing-impaired listeners compared to age-matched normal-hearing listeners. The enhanced
479 tracking in older listeners or listeners with a hearing impairment may be explained by a compensatory central gain
480 mechanism (Parthasarathy et al., 2019; De Villers-Sidani et al., 2010; Chambers et al., 2016), recruitment of additional
481 cortical resources (Brodbeck et al., 2018b; Gillis et al., 2021a) and increased listening effort and attention (Decruy
482 et al., 2020a; Vanthornhout et al., 2019; Lesenfants and Francart, 2020). This shows that it is also important to conduct
483 subject-specific analyses and not only at group-level measures. With an innovative artefact removal technique, Somers
484 et al. (2019) succeeded to analyse envelope tracking for cochlear implant listeners as well. For both hearing-impaired
485 listeners (with simulated amplification) (Decruy et al., 2020b) and cochlear implant listeners (Verschuere et al., 2019)
486 the tracking strength was significantly correlated to behaviourally-measured speech intelligibility, indicating a similar
487 relation with speech intelligibility as observed for normal hearing listeners (Vanthornhout et al., 2018).

488 One challenge with envelope tracking is that its functional interpretation is unclear. The main complicating factor is
489 that the envelope is highly correlated with linguistic cues, like the onsets of words and syllables. As such, the envelope
490 represents multiple unique features that all may contribute to the observed neural tracking response and are hard to
491 disentangle. In addition, the interpretation of envelope tracking is complicated because it is modulated by top-down
492 effects, such as attention and audio-visual integration (O’Sullivan et al., 2019a). A final challenge is that the exact
493 relation between envelope tracking and speech intelligibility remains a point of discussion (Ding and Simon, 2014;
494 Brodbeck and Simon, 2020). Multiple studies have shown that the level of envelope tracking reflects experimental
495 changes in speech intelligibility (Vanthornhout et al., 2018; Lesenfants et al., 2019; Verschueren et al., 2021), even in
496 the case of degraded speech with an intact envelope (Ding et al., 2014). However, it is unlikely that envelope tracking
497 is a direct reflection of successful speech intelligibility as neural tracking responses have been observed for non-speech
498 signals (Zuk et al., 2021) and foreign languages (Etard and Reichenbach, 2019). As such, envelope tracking is likely
499 necessary but not sufficient for speech intelligibility. Linguistic features can be used to gain further insight into how
500 the brain processes the meaning of speech, i.e. speech intelligibility.

501 **6. Neural tracking of linguistic features**

502 Recent studies focus on linguistic speech features in pursuit of an accurate neural marker of speech intelligibility.
503 While the f_0 and speech envelope are derived from the acoustic waveform of the speech, linguistic features are
504 derived from the content of the speech. Proper encoding of these features in the brain requires accurate linguistic and
505 not mere acoustic processing.

506 Linguistic features can be divided into two categories. Features in the first category denote lexical onsets. They rep-
507 resent (aspects of) a sequence of small building blocks that make up spoken language, e.g., sequences of phonemes,
508 phonetic features, words, or specific word categories like content and function words (Di Liberto et al., 2015; Lesen-
509 fans et al., 2019). These features are sparse arrays consisting of zeros with a fixed, non-zero entry (\sim impulse) at
510 the onset of each lexical building block (see features in light green on Panel C of figure 2). The neural responses to
511 lexical onset features are not straightforward to interpret. As phonemes, syllables, and words coincide with acoustic
512 cues, the associated neural response is neither purely lexical nor acoustic.

513 Features in the second category reflect higher-level linguistic aspects of the speech, e.g., how familiar, predictable or
514 surprising a word or phoneme is in its context. These features can be applied on three levels, which require different
515 amounts of linguistic context: (1) at the level of a phoneme (e.g., phoneme surprisal or cohort entropy (Di Liberto
516 et al., 2019; Brodbeck et al., 2018a)), (2) at the level of a word (e.g., word frequency or word surprisal (Weissbart
517 et al., 2019; Koskinen et al., 2020)), and (3) at a semantic contextual level (e.g., semantic dissimilarity (Broderick
518 et al., 2018)). These features are sparse arrays, similar to lexical onset features. However, in this case, the impulse
519 amplitude at each onset is not fixed but modulated by the linguistic information of the specific phoneme or word (see
520 features in dark green on Panel C of figure 2).

521 The fact that linguistic features are sparse arrays consisting of mostly zeroes with some non-zero entries (~ impulses),
522 makes them different from the continuous f0 and envelope features and poses challenges for response analysis. In
523 decoding modelling, the reconstructed feature is compared to the actual feature. However, sparse features are chal-
524 lenging to reconstruct from a continuous signal (the EEG), as we can only reconstruct a continuous signal from the
525 continuous input using a linear model. This problem does not occur for encoding modelling, where the non-sparse
526 predicted EEG is compared to the actual EEG. Therefore the encoding model is a more common choice for analysis
527 with linguistic features.

528 Panels E-H of figure 4 present a visualisation of the results of a typical encoding modelling analysis for linguistic
529 tracking with phoneme surprisal and word surprisal features. These results were obtained by filtering the EEG responses
530 between 0.5 and 25 Hz and referencing to the average EEG response (32 participants; Gillis et al. (more details
531 regarding the preprocessing are described in 2021b)). The TRFs at both phoneme (panel E) and word level (panel
532 G) show a negative response, situated centrally in the topography (panel F and H), around respectively 250 and 350
533 ms. The earlier response peak for phonemes compared to words is consistent with the hierarchy of the language
534 processing of these linguistic building blocks, i.e., the phonemes making up a word are processed before the word's
535 surprisal can be estimated. Moreover, the response to word surprisal resembles the N400 response, which is classically
536 observed in ERP paradigms (Lau et al., 2008). These congruent topographic responses indicate that this small and
537 specific language response can also be observed when listening to natural running speech rather than stand-alone
538 sentences.

539 Measuring neural tracking of linguistic features is an exciting avenue to test psycho-linguistic theories of speech
540 understanding. It is accepted that listeners use linguistic context to continuously adapt expectations of upcoming
541 concepts, words and phonemes, However, how these expectations are integrated with what is actually being perceived
542 is unclear. Brodbeck et al. (2021a) showed that the neural prediction of an upcoming phoneme or word relies on
543 contextual processing in a parallel manner, combining both bottom-up and top-down processing. Additional evidence
544 of the presence of top-down processing comes from Heilbron et al. (2020) who observed that higher-level predictions
545 influence the predictions at lower levels (i.e., word prediction affects the predictions at the phoneme level).

546 Additionally, linguistic features allow investigating to what extent speech is understood given the language proficiency.
547 Di Liberto et al. (2021) investigated neural tracking in Mandarin speakers with different levels of English proficiency.
548 Interestingly, the magnitude of central negative response to semantic dissimilarity around 400 ms increased with
549 proficiency.

550 Another exciting research path is the disentanglement of acoustic and linguistic neural processing. Verschueren et al.
551 (2022) disentangled acoustic and linguistic neural processing by changing the speech rate, which kept the linguistic
552 content the same while varying the acoustic properties and the intelligibility of the speech. As the speech rate in-
553 creased, the neural tracking of acoustic properties increased. This means that better time-locking was observed, even

554 though the speech became harder to understand. In contrast, neural tracking of linguistic properties decreased with
555 increasing speech rate. This indicates that linguistic tracking provides a more accurate objective measure of speech
556 intelligibility.

557 The two studies mentioned above indicate that neural tracking of linguistic features encodes aspects of neural language
558 processing. These findings open doors to study language development in young children or to objectively determine
559 speech understanding.

560 Linguistic speech features can also provide insight into age-related speech intelligibility deficits. We are aware of
561 two studies investigating speech intelligibility deficits in older adults. Although Mesik et al. (2021) did not report
562 differences, Broderick et al. (2021) reported that older adults rely less on semantic features than younger adults.
563 Furthermore, they showed that older adults who relied more on this semantic mechanism showed higher verbal fluency
564 than older adults with weaker semantic tracking.

565 Linguistic tracking is an up-and-coming research technique but has a few difficulties. Firstly, the inter-feature correla-
566 tion should be taken into account. The linguistic features coincide with the boundaries of phonemes and words. These
567 boundaries are often associated with high acoustic power; therefore, it is necessary to carefully control for acoustic
568 properties of the speech when evaluating linguistic features. If not, the speech tracking analysis might be biased to
569 find spurious significant linguistic features due to its correlations with acoustic features (Daube et al., 2019). We
570 proposed some approaches to deal with this inter-feature correlation in the section 3.1.

571 Secondly, the analyses are often based on encoding modelling due to sparse features. Prediction accuracies, i.e.
572 correlations, obtained with encoding models are typically small in magnitude: only around 3 to 7% of the variance
573 in the EEG signal can be explained by neural responses time-locked to the presented stimulus. Moreover, most
574 of this variance is explained by acoustic characteristics of the speech, as these lower-level acoustic features evoke
575 responses over large parts of the auditory system. In contrast, linguistic tracking targets the neural response from
576 a precisely localised neural process related to intelligibility. Therefore, the associated magnitudes of these neural
577 processes measured at the scalp level are much smaller. As the prediction accuracies of the encoding model are small
578 in magnitude, finding a significant improvement of the linguistic feature over and beyond acoustic features requires
579 enough observations (e.g. an improvement of ~1% corresponds to an increase in prediction accuracy of 3.4×10^{-4}
580 approach as described above (Gillis et al., 2021b)).

581 Thirdly, estimating the TRF to these linguistic features requires a lot of data. In Figure 1, a 15-minute story is used to
582 estimate the model, which is substantially shorter than most studies which evaluate linguistic tracking of speech in 45
583 to 60 minutes. This difference in the amount of data explains why the TRF pattern is noisier. The more data is used
584 to estimate the model, the better the brain response to linguistic tracking is characterised, and the more prominent
585 the peaks are. A method to overcome this is to estimate a subject-independent model whereby the data of different
586 subjects is combined to ensure enough data.

587 Another difficulty is that current studies investigating the effect of linguistic neural tracking only evaluate population
588 differences, such as comparing young to old or comparing a baseline to a complete model. As the effect is small
589 and the TRFs show a high between- and intra-subject variability, it is challenging to extract an objective measure at
590 a subject-specific level. Future research should focus on making the TRF patterns and prediction accuracies more
591 reliable and robust at a subject-specific level.

592 **7. Caveats**

593 A first caveat of the studies investigating neural tracking is that most studies focus on models assuming a linear rela-
594 tionship between speech and EEG responses. However, whether or not this assumption is valid remains unanswered.
595 Additionally, there is no ground truth of which speech features the brain tracks. Therefore, investigating which speech
596 features are tracked by the brain remains an explorative search which might lead to suboptimal results.

597 A second caveat is the comparison of neural tracking of speech with behavioural measures of speech understanding.
598 Behavioural measures can be obtained by, for example, sentence recall tests. However, when using continuous speech,
599 these tests are not suitable because longer text fragments cannot be recalled. As a result, the speaker of the behavioural
600 measures is often different from the speaker of the stimuli used for the neural measurements. This can be problematic
601 as speaker characteristics can affect the neural tracking of speech. As sentence recall tests are unsuitable, the current
602 field of research lacks a good evaluation of speech understanding for continuous text. Currently used metrics are
603 content questions or subjectively rated speech understanding, i.e., the participant's answer to the question 'how much
604 did you understand?'. To overcome the intersubject variability of these subjective ratings, the self-assessed Békesy
605 procedure can be applied to rescale the subjective ratings towards a more objective alternative (Decruy et al., 2018).
606 However, these metrics remain sub-optimal as content questions rely heavily on attention effects and subjectively
607 rated speech understanding introduces a bias.

608 Another caveat is that neural tracking is sometimes confused with neural entrainment. Both concepts underlie different
609 assumptions. Neural tracking assumes that the neural responses time-lock to different sound features. Speech features
610 evoke a cascade of different responses, such as responses that are time-locked to the acoustical, lexical, and linguistic
611 features. In contrast, neural entrainment assumes that the neural responses phase-lock to the stimulus. This theory
612 assumes that oscillations are present in the brain without stimulation. When a rhythmic sound, such as speech, is
613 heard, these oscillations reset their phases and become synchronised with the dominant phase of the speech signal
614 (Pelle and Davis, 2012). Phase-locking can be measured, for example, using inter-trial correlations (Ding et al.,
615 2014). As a result, neural entrainment involves more assumptions than neural tracking. Therefore, it is possible that
616 speech understanding can occur when there is no neural entrainment, for example, when the speech signal is very
617 arrhythmic (Pelle et al., 2013). On the other hand, having neural tracking does not guarantee that the speech signal
618 was understood correctly.

619 **8. Clinical applications of measurements of neural tracking responses**

620 To provide care to people with hearing problems, it is useful to review the merits and limitations of all (objective)
621 audiological measures and investigate how the measures may be combined to form a complete assessment of the
622 auditory system.

623 The current gold standard methods, i.e. tone and speech audiometry, have proven their worth. However, they are
624 challenging in crucial populations like young children or people with a cognitive impairment like dementia. Objective
625 measures for sound perception like the ABR and the ASSR have been introduced in the clinical toolset to remedy this.
626 However, there is no clinically available objective measure of speech intelligibility. Since speech intelligibility is the
627 basis for human communication, this is a significant gap to fill. Various populations may benefit from such a measure,
628 including young children, stroke patients (especially those with aphasia) and people with dementia.

629 Measurement of neural tracking is a versatile tool as the amount of tracking to the different speech features and thus
630 in different parts of the auditory system can be measured. Based on a single twenty-minute long EEG recording,
631 a wide range of speech processing abilities may be assessed simultaneously (incl. f0 tracking, envelope tracking,
632 phonetic processing, phonemic processing, syntactical processing and even linguistic and emotional processing). This
633 versatility may lead to an objective assessment of both auditory and language abilities. Moreover, measuring neural
634 tracking is easily automated, paving the way to improved automated screening, diagnostics, and automatic fitting of
635 auditory prostheses, or even auditory prostheses that continuously adapt themselves to the listener based on their brain
636 activity (Geirnaert et al., 2021).

637 Future studies preparing for clinical implementation may need to shift focus from group-level analyses towards
638 subject-specific analyses. Going towards a subject-specific analysis is key to allow its feasibility as a clinical marker.
639 Although the magnitude of neural tracking might be intrinsically different between different subjects, relative differ-
640 ences between different conditions can be used as a diagnostic marker. Additionally, future studies may focus on
641 which combination of stimulus features provides the most information and how these can be optimally analysed. As
642 the features are highly correlated, special care needs to be taken to investigate the effect of each feature (Gillis et al.,
643 2021b). Subsequent research efforts are also required to decide on the best speech stimuli (required to work well
644 for all types of tracking) and the best EEG measurement set-up, including the number of EEG electrodes and their
645 position (Montoya-Martínez et al., 2021) to reduce recording time. Furthermore, it is important to set best practices of
646 how the measures of neural tracking can be used (Crosse et al., 2021) and to have insight into how the preprocessing
647 influences the results (de Cheveigné and Nelken, 2019). It is also essential to validate the measures in a comprehen-
648 sive sample of the population, including participants of all ages and with various audiological and non-audiological
649 pathologies (for a review, see Palana et al. (2021)). Finally, the neural tracking results need to be transformed into an
650 easy-to-interpret set of scores and visualisations to allow for intuitive use by clinicians.

651 **9. Acknowledgements**

652 The authors would like to thank Bernd Accou, Wendy Verheijen and their students for collecting the dataset used for
653 the examples in this article. The research received funding from the European Research Council under the European
654 Unions Horizon 2020 research and innovation programme (grant agreement No. 637424, ERC starting grant to Tom
655 Francart). Jana Van Canneyt and Marlies Gillis are both supported by a PhD grant for Strategic Basic research by
656 the Research Foundation Flanders (FWO): project number 1S83618N and project number 1SA0620N, respectively.
657 Jonas Vanthornhout is funded by a postdoctoral grant from FWO, project number 1290821N.

658 **10. Declaration of interest**

659 The authors declare that author Tom Francart is involved in translational research which may lead to the commer-
660 cialisation of a product related to the presented research. Besides this, there are no conflicts of interest, financial, or
661 otherwise.

662 **References**

- 663 Accou, B., Jalilpour Monesi, M., Montoya, J., Van hamme, H., and Francart, T. (2021). Modeling the relationship between acoustic stimulus and
664 EEG with a dilated convolutional neural network. In *2020 28th European Signal Processing Conference (EUSIPCO)*, pages 1175–1179.
- 665 Aiken, S. J. and Picton, T. W. (2008). Human cortical responses to the speech envelope. *Ear and Hearing*, 29(2):139–157.
- 666 Aljarboa, G. S., Bell, S. L., and Simpson, D. M. (2022). Detecting cortical responses to continuous running speech using eeg data from only one
667 channel. *International Journal of Audiology*, pages 1–10.
- 668 Biesmans, W., Das, N., Francart, T., and Bertrand, A. (2017). Auditory-Inspired Speech Envelope Extraction Methods for Improved EEG-Based
669 Auditory Attention Detection in a Cocktail Party Scenario. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 25(5):402–
670 412.
- 671 Bollens, L., Francart, T., and Hamme, H. V. (2022). Learning subject-invariant representations from speech-evoked eeg using variational autoen-
672 coders. In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1256–1260.
- 673 Brennan, J. R. and Hale, J. T. (2019). Hierarchical structure guides rapid linguistic predictions during naturalistic listening. *PloS one*,
674 14(1):e0207741.
- 675 Brodbeck, C., Bhattasali, S., Heredia, A. C., Resnik, P., Simon, J. Z., and Lau, E. (2021a). Parallel processing in speech perception: Local and
676 global representations of linguistic context. *bioRxiv*, page 2021.07.03.450698.
- 677 Brodbeck, C., Das, P., Kulasingham, J. P., Bhattasali, S., Gaston, P., Resnik, P., and Simon, J. Z. (2021b). Eelbrain: A python toolkit for time-
678 continuous analysis with temporal response functions. *BioRxiv*.
- 679 Brodbeck, C., Hong, L. E., and Simon, J. Z. (2018a). Rapid Transformation from Auditory to Linguistic Representations of Continuous Speech.
680 *Current Biology*, 28(24):3976–3983.e5.
- 681 Brodbeck, C., Presacco, A., Anderson, S., and Simon, J. Z. (2018b). Over-representation of speech in older adults originates from early response
682 in higher order auditory cortex. *Acta Acustica united with Acustica*, 104(5):774–777.
- 683 Brodbeck, C., Presacco, A., and Simon, J. Z. (2018c). Neural Source Dynamics of Brain Responses to Continuous Stimuli: {{Speech}} Processing
684 from Acoustics to Comprehension. *NeuroImage*, 172:162–174.
- 685 Brodbeck, C. and Simon, J. Z. (2020). Continuous speech processing. *Current Opinion in Psychology*, 18:25–31.
- 686 Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J., and Lalor, E. C. (2018). Electrophysiological Correlates of Semantic
687 Dissimilarity Reflect the Comprehension of Natural, Narrative Speech. *Current Biology*, 28(5):803–809.e3.

- 688 Broderick, M. P., Di Liberto, G. M., Anderson, A. J., Rofes, A., and Lalor, E. C. (2021). Dissociable electrophysiological measures of natural
689 language processing reveal differences in speech comprehension strategy in healthy ageing. *Scientific Reports*, 11(1):1–12.
- 690 Bruce, I. C., Erfani, Y., and Zilany, M. S. (2018). A phenomenological model of the synapse between the inner hair cell and auditory nerve:
691 Implications of limited neurotransmitter release sites. *Hearing Research*, 360:40–54.
- 692 Brugge, J. F., Nourski, K. V., Oya, H., Reale, R. A., Kawasaki, H., Steinschneider, M., and Howard, M. A. (2009). Coding of Repetitive Transients
693 by Auditory Cortex on Heschl’s Gyrus. *Journal of Neurophysiology*, 102(4):2358–2374.
- 694 Chambers, A. R., Resnik, J., Yuan, Y., Whitton, J. P., Edge, A. S., Liberman, M. C., and Polley, D. B. (2016). Central Gain Restores Auditory
695 Processing following Near-Complete Cochlear Denervation. *Neuron*, 89(4):867–879.
- 696 Chandrasekaran, C., Trubanova, A., Stillitano, S., Caplier, A., and Ghazanfar, A. A. (2009). The natural statistics of audiovisual speech. *PLoS*
697 *Computational Biology*, 5(7).
- 698 Coffey, E. B., Musacchia, G., and Zatorre, R. J. (2017). Cortical Correlates of the Auditory Frequency-Following and Onset Responses: EEG and
699 fMRI Evidence. *The Journal of Neuroscience*, 37(4):830–838.
- 700 Coffey, E. B. J., Herholz, S. C., Chepesiuk, A. M. P., Baillet, S., and Zatorre, R. J. (2016). Cortical contributions to the auditory frequency-following
701 response revealed by MEG. *Nature Communications*, 7:11070.
- 702 Crosse, M. J., Di Liberto, G. M., Bednar, A., and Lalor, E. C. (2016a). The multivariate temporal response function (mTRF) toolbox: A MATLAB
703 toolbox for relating neural signals to continuous stimuli. *Frontiers in Human Neuroscience*, 10(NOV2016).
- 704 Crosse, M. J., Liberto, D., M, G., Bednar, A., and Lalor, E. C. (2016b). The {{Multivariate Temporal Response Function}} ({{mTRF}}) {{Toolbox}}:
705 {{A MATLAB Toolbox}} for {{Relating Neural Signals}} to {{Continuous Stimuli}}. *Front. Hum. Neurosci.*, 10.
- 706 Crosse, M. J., Zuk, N. J., Di Liberto, G. M., Nidiffer, A. R., Molholm, S., and Lalor, E. C. (2021). Linear modeling of neurophysiological responses
707 to speech and other continuous stimuli: methodological considerations for applied research. *Frontiers in Neuroscience*, 15.
- 708 Das, N., Bertrand, A., and Francart, T. (2018). EEG-based auditory attention detection: boundary conditions for background noise and speaker
709 positions. *Journal of Neural Engineering*, 15(6):066017.
- 710 Das, N., Vanthornhout, J., Francart, T., and Bertrand, A. (2019). Stimulus-Aware Spatial Filtering for Single-Trial Neural Response and Temporal
711 Response Function Estimation in High-Density {{EEG}} with Applications in Auditory Research. *NeuroImage*, page 116211.
- 712 Daube, C., Ince, R. A., and Gross, J. (2019). Simple Acoustic Features Can Explain Phoneme-Based Predictions of Cortical Responses to Speech.
713 *Current Biology*, 29(12):1924–1937.e9.
- 714 David, S. V., Mesgarani, N., and Shamma, S. A. (2007). Estimating Sparse Spectro-Temporal Receptive Fields with Natural Stimuli. *Network:*
715 *Computation in Neural Systems*, 18(3):191–212.
- 716 de Cheveigné, A., Di Liberto, G. M., Arzounian, D., Wong, D. D., Hjortkjaer, J., Fuglsang, S., and Parra, L. C. (2019). Multiway canonical
717 correlation analysis of brain data. *NeuroImage*, 186(November 2018):728–740.
- 718 de Cheveigné, A. and Nelken, I. (2019). Filters: When, Why, and How (Not) to Use Them. *Neuron*, 102(2):280–293.
- 719 De Cheveigné, A., Slaney, M., Fuglsang, S. A., and Hjortkjaer, J. (2021). Auditory stimulus-response modeling with a match-mismatch task.
720 *Journal of Neural Engineering*, 18(4):046040.
- 721 de Cheveigné, A., Wong, D. D., Di Liberto, G. M., Hjortkjaer, J., Slaney, M., and Lalor, E. (2018). Decoding the auditory brain with canonical
722 component analysis. *NeuroImage*, 172:206–216.
- 723 De Villers-Sidani, E., Alzghoul, L., Zhou, X., Simpson, K. L., Lin, R. C., and Merzenich, M. M. (2010). Recovery of functional and structural
724 age-related changes in the rat primary auditory cortex with operant training. *Proceedings of the National Academy of Sciences of the United*
725 *States of America*, 107(31):13900–13905.
- 726 Deckers, L., Das, N., Ansari, A. H., Bertrand, A., and Francart, T. (2018). EEG-based detection of the attended speaker and the locus of auditory
727 attention with convolutional neural networks. *bioRxiv*, page 475673.
- 728 Decruy, L., Das, N., Verschueren, E., and Francart, T. (2018). The self-assessed békesy procedure: validation of a method to measure intelligibility
729 of connected discourse. *Trends in hearing*, 22:2331216518802702.
- 730 Decruy, L., Lesenfants, D., Vanthornhout, J., and Francart, T. (2020a). Top-down modulation of neural envelope tracking: The interplay with

- 731 behavioral, self-report and neural measures of listening effort. *European Journal of Neuroscience*, European J(October 2019):3375–3393.
- 732 Decruy, L., Vanthornhout, J., and Francart, T. (2019). Evidence for enhanced neural tracking of the speech envelope underlying age-related
733 speech-in-noise difficulties. *Journal of Neurophysiology*, 122(2):601–615.
- 734 Decruy, L., Vanthornhout, J., and Francart, T. (2020b). Hearing impairment is associated with enhanced neural tracking of the speech envelope.
735 *Hearing Research*, 393:107961.
- 736 Di Liberto, G. M., Crosse, M. J., and Lalor, E. C. (2018). Cortical Measures of Phoneme-Level Speech Encoding Correlate with the Perceived
737 Clarity of Natural Speech. *eneuro*, 5(2):ENEURO.0084–18.2018.
- 738 Di Liberto, G. M., Nie, J., Yeaton, J., Khalighinejad, B., Shamma, S. A., and Mesgarani, N. (2021). Neural representation of linguistic feature
739 hierarchy reflects second-language proficiency. *Neuroimage*, 227:117586.
- 740 Di Liberto, G. M., O’Sullivan, J. A., and Lalor, E. C. (2015). Low-frequency cortical entrainment to speech reflects phoneme-level processing.
741 *Current Biology*, 25(19):2457–2465.
- 742 Di Liberto, G. M., Wong, D., Melnik, G. A., and de Cheveigné, A. (2019). Low-frequency cortical responses to natural speech reflect probabilistic
743 phonotactics. *Neuroimage*, 196:237–247.
- 744 Ding, N., Chatterjee, M., and Simon, J. Z. (2014). Robust cortical entrainment to the speech envelope relies on the spectro-temporal fine structure.
745 *NeuroImage*, 88:41–46.
- 746 Ding, N. and Simon, J. Z. (2012a). Emergence of Neural Encoding of Auditory Objects While Listening to Competing Speakers. *PNAS*,
747 109(29):11854–11859.
- 748 Ding, N. and Simon, J. Z. (2012b). Neural coding of continuous speech in auditory cortex during monaural and dichotic listening. *Journal of*
749 *Neurophysiology*, 107(1):78–89.
- 750 Ding, N. and Simon, J. Z. (2013). Adaptive temporal encoding leads to a background-insensitive cortical representation of speech. *Journal of*
751 *Neuroscience*, 33(13):5728–5735.
- 752 Ding, N. and Simon, J. Z. (2014). Cortical entrainment to continuous speech: Functional roles and interpretations. *Frontiers in Human Neuro-*
753 *science*, 8(MAY):1–7.
- 754 Drullman, R., Festen, J. M., and Plomp, R. (1994a). Effect of Reducing Slow Temporal Modulations on Speech Reception. *The Journal of the*
755 *Acoustical Society of America*, 95(5):2670–2680.
- 756 Drullman, R., Festen, J. M., and Plomp, R. (1994b). Effect of Temporal Envelope Smearing on Speech Reception. *The Journal of the Acoustical*
757 *Society of America*, 95(2):1053–1064.
- 758 Etard, O., Kegler, M., Braiman, C., Forte, A. E., and Reichenbach, T. (2019). Decoding of selective attention to continuous speech from the human
759 auditory brainstem response. *NeuroImage*, 200(May):1–11.
- 760 Etard, O. and Reichenbach, T. (2019). Neural Speech Tracking in the Theta and in the Delta Frequency Band Differentially Encode Clarity and
761 Comprehension of Speech in Noise. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, 39(29):5750–5759.
- 762 Forte, A. E., Etard, O., and Reichenbach, T. (2017). The human auditory brainstem response to running speech reveals a subcortical mechanism
763 for selective attention. *eLife*, 6:1–13.
- 764 Fuglsang, S. A., Märcher-Rørsted, J., Dau, T., and Hjortkjær, J. (2020). Effects of sensorineural hearing loss on cortical synchronization to
765 competing speech during selective attention. *Journal of Neuroscience*, 40(12):2562–2572.
- 766 Geirnaert, S., Vandecappelle, S., Alickovic, E., de Cheveigne, A., Lalor, E., Meyer, B., Miran, S., Francart, T., and Bertrand, A. (2021).
767 Electroencephalography-based Auditory Attention Decoding : Toward Neuro-Steered Hearing Devices. *Ieee Signal Processing Magazine*.
768 *Special issue on Signal Processing for Neurorehabilitation and Assistive Technologies*, 38(4):89–102.
- 769 Gillis, M., Decruy, L., Vanthornhout, J., and Francart, T. (2021a). Hearing loss is associated with delayed neural responses to continuous speech.
770 *bioRxiv*, 2021.01.21.
- 771 Gillis, M., Vanthornhout, J., Simon, J. Z., Francart, T., and Brodbeck, C. (2021b). Neural markers of speech comprehension: measuring EEG
772 tracking of linguistic speech representations, controlling the speech acoustics. *The Journal of Neuroscience*, (October):JN–RM–0812–21.
- 773 Godey, B., Schwartz, D., de Graaf, J. B., Chauvel, P., and Liégeois-Chauvel, C. (2001). Neuromagnetic source localization of auditory evoked

- 774 fields and intracerebral evoked potentials: a comparison of data in the same patients. *Clinical Neurophysiology*, 112(10):1850–1859.
- 775 Gross, J., Hoogenboom, N., Thut, G., Schyns, P., Panzeri, S., Belin, P., and Garrod, S. (2014). Speech Rhythms and Multiplexed Oscillatory
776 Sensory Coding in the Human Brain. *PLOS Biology*, 11(12):1–14.
- 777 Hamilton, L. S. and Huth, A. G. (2018). The revolution will not be controlled: natural stimuli in speech neuroscience. *Language, Cognition and
778 Neuroscience*, 35(5):573–582.
- 779 Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J.-D. D., Blankertz, B., and Bießmann, F. (2014). On the interpretation of weight vectors
780 of linear models in multivariate neuroimaging. *NeuroImage*, 87:96–110.
- 781 Heilbron, M., Armeni, K., Schoffelen, J.-M., Hagoort, P., and De Lange, F. P. (2020). A hierarchy of linguistic predictions during natural language
782 comprehension. *bioRxiv*, page 2020.12.03.410399.
- 783 Iotzov, I. and Parra, L. C. (2019). EEG can predict speech intelligibility. *Journal of Neural Engineering*, 16(3):036008.
- 784 Joris, P. X., Schreiner, C. E., and Rees, A. (2004). Neural Processing of Amplitude-Modulated Sounds. *Physiological Reviews*, 84(2):541–577.
- 785 Kaufeld, G., Bosker, H. R., Alday, P. M., Meyer, A. S., and Martin, A. E. (2020). Structure and meaning organize neural oscillations into a
786 content-specific hierarchy. *bioRxiv*, 40(49):9467–9475.
- 787 Kei, J., Smyth, V., Murdoch, B., and McPherson, B. (1999). Measuring the Understanding of Connected Discourse: An Overview of Methodology
788 and Clinical Applications in Rehabilitative Audiology. *Asia Pacific Journal of Speech, Language and Hearing*, 4(1):13–37.
- 789 Keidser, G., Naylor, G., Brungart, D. S., Caduff, A., Campos, J., Carlile, S., Carpenter, M. G., Grimm, G., Hohmann, V., Holube, I., Launer, S.,
790 Lunner, T., Mehra, R., Rapport, F., Slaney, M., and Smeds, K. (2020). The Quest for Ecological Validity in Hearing Science: What It Is, Why
791 It Matters, and How to Advance It. *Ear and hearing*, 41:5S–19S.
- 792 Kong, Y.-Y., Mullangi, A., and Ding, N. (2014). Differential Modulation of Auditory Responses to Attended and Unattended Speech in Different
793 Listening Conditions. *Hearing Research*, 316:73–81.
- 794 Koskinen, M., Kurimo, M., Gross, J., Hyvärinen, A., and Hari, R. (2020). Brain activity reflects the predictability of word sequences in listened
795 continuous speech: Brain activity predicts word sequences. *NeuroImage*, 219(May).
- 796 Kulasingham, J. P., Brodbeck, C., Presacco, A., Kuchinsky, S. E., Anderson, S., and Simon, J. Z. (2020). High gamma cortical processing of
797 continuous speech in younger and older listeners. *NeuroImage*, 222(June):117291.
- 798 Kulasingham, J. P. and Simon, J. Z. (2022). Algorithms for estimating time-locked neural response components in cortical processing of continuous
799 speech. *IEEE Transactions on Biomedical Engineering*.
- 800 Lalor, E. C. and Foxe, J. J. (2010). Neural responses to uninterrupted natural speech can be extracted with precise temporal resolution. *European
801 Journal of Neuroscience*, 31(1):189–193.
- 802 Lalor, E. C., Power, A. J., Reilly, R. B., and Foxe, J. J. (2009). Resolving Precise Temporal Processing Properties of the Auditory System Using
803 Continuous Stimuli. *Journal of Neurophysiology*, 102(1):349–359.
- 804 Lau, E. F., Phillips, C., and Poeppel, D. (2008). A cortical network for semantics: (de)constructing the N400.
- 805 Lesenfans, D. and Francart, T. (2020). The interplay of top-down focal attention and the cortical tracking of speech. *Scientific Reports*, 10(1):1–11.
- 806 Lesenfans, D., Vanthornhout, J., Verschueren, E., Decruy, L., and Francart, T. (2019). Predicting individual speech intelligibility from the cortical
807 tracking of acoustic- and phonetic-level speech representations. *Hearing Research*, 380:1–9.
- 808 Machens, C. K., Wehr, M. S., and Zador, A. M. (2004). Linearity of cortical receptive fields measured with natural sounds. *Journal of Neuroscience*,
809 24(5):1089–1100.
- 810 Martin, B. A., Tremblay, K. L., and Korczak, P. (2008). Speech evoked potentials: From the laboratory to the clinic. *Ear and Hearing*, 29(3):285–
811 313.
- 812 Mesik, J., Ray, L., and Wojtczak, M. (2021). Effects of Age on Cortical Tracking of Word-Level Features of Continuous Competing Speech.
813 *Frontiers in Neuroscience*, 15(April):1–21.
- 814 Molinaro, N. and Lizarazu, M. (2017). Delta (but Not Theta)-band Cortical Entrainment Involves Speech-specific Processing. *European Journal of
815 Neuroscience*, 48(7).
- 816 Monesi, M. J., Accou, B., Francart, T., and Van Hamme, H. (2021). Extracting different levels of speech information from eeg using an lstm-based

817 model. *arXiv preprint arXiv:2106.09622*.

818 Montoya-Martínez, J., Vanthornhout, J., Bertrand, A., and Francart, T. (2021). Effect of number and placement of EEG electrodes on measurement
819 of neural tracking of speech. *PLoS ONE*, 16(2 February):1–18.

820 Nourski, K. V., Reale, R. A., Oya, H., Kawasaki, H., Kovach, C. K., Chen, H., Howard, M. A., and Brugge, J. F. (2009). Temporal Envelope of
821 Time-Compressed Speech Represented in the Human Auditory Cortex. *Journal of Neuroscience*, 29(49):15564–15574.

822 O’Sullivan, A. E., Crosse, M. J., Di Liberto, G. M., de Cheveigné, A., and Lalor, E. C. (2021). Neurophysiological indices of audiovisual speech
823 processing reveal a hierarchy of multisensory integration effects. *Journal of Neuroscience*, 41(23):4991–5003.

824 O’Sullivan, A. E., Lim, C. Y., and Lalor, E. C. (2019a). Look at me when I’m talking to you: Selective attention at a multisensory cocktail party
825 can be decoded using stimulus reconstruction and alpha power modulations. *European Journal of Neuroscience*, 50(8):3282–3295.

826 O’Sullivan, J., Herrero, J., Smith, E., Schevon, C., McKhann, G. M., Sheth, S. A., Mehta, A. D., and Mesgarani, N. (2019b). Hierarchical Encoding
827 of Attended Auditory Objects in Multi-talker Speech Perception. *Neuron (Cambridge, Mass.)*, 104(6):1195—1209.e3.

828 O’Sullivan, J. A., Power, A. J., Mesgarani, N., Rajaram, S., Foxe, J. J., Shinn-Cunningham, B. G., Slaney, M., Shamma, S. A., and Lalor, E. C.
829 (2015). Attentional Selection in a Cocktail Party Environment Can Be Decoded from Single-Trial EEG. *Cerebral Cortex*, 25(7):1697–1706.

830 Palana, J., Schwartz, S., and Tager-Flusberg, H. (2021). Evaluating the use of cortical entrainment to measure atypical speech processing: A
831 systematic review. *Neuroscience & Biobehavioral Reviews*.

832 Parthasarathy, A., Bartlett, E. L., and Kujawa, S. G. (2019). Age-related Changes in Neural Coding of Envelope Cues: Peripheral Declines and
833 Central Compensation. *Neuroscience*, 407:21–31.

834 Peelle, J. E. and Davis, M. H. (2012). Neural Oscillations Carry Speech Rhythm through to Comprehension. *Frontiers in Psychology*,
835 3(September):1–17.

836 Peelle, J. E., Gross, J., and Davis, M. H. (2013). Phase-locked responses to speech in human auditory cortex are enhanced during comprehension.
837 *Cerebral cortex*, 23(6):1378–1387.

838 Petersen, E. B., Wöstmann, M., Obleser, J., and Lunner, T. (2016). Neural Tracking of Attended versus Ignored Speech Is Differentially Affected
839 by Hearing Loss. *Journal of Neurophysiology*, 117(1):18–27.

840 Pichora-Fuller, M. K., Kramer, S. E., Eckert, M. A., Edwards, B., Hornsby, B. W., Humes, L. E., Lemke, U., Lunner, T., Matthen, M., Mackersie,
841 C. L., Naylor, G., Phillips, N. A., Richter, M., Rudner, M., Sommers, M. S., Tremblay, K. L., and Wingfield, A. (2016). Hearing Impairment
842 and Cognitive Energy: The Framework for Understanding Effortful Listening (FUEL). *Ear & Hearing*, 37(1):5S–27S.

843 Picton, T. W. (2010). *Human Auditory Evoked Potentials*. Plural Pub.

844 Rosen, S. (1992). Temporal Information in Speech: Acoustic, Auditory and Linguistic Aspects. *Phil. Trans. R. Soc. Lond. B*, 336(1278):367–373.

845 Saiz-Alfá, M. and Reichenbach, T. (2020). Computational modeling of the auditory brainstem response to continuous speech. *Journal of Neural
846 Engineering*, 17(3):036035.

847 Särelä, J., Valpola, H., and Jordan, M. (2005). Denoising source separation. *Journal of machine learning research*, 6(3).

848 Shannon, R. V., Zeng, F.-G., Kamath, V., Wygonski, J., Ekelid, M., Series, N., and Oct, N. (1995). Speech Recognition with Primarily Temporal
849 Cues. *Source: Science, New Series*, 270(5234):303–304.

850 Somers, B., Francart, T., and Bertrand, A. (2018). A generic eeg artifact removal algorithm based on the multi-channel wiener filter. *Journal of
851 neural engineering*, 15(3):036007.

852 Somers, B., Verschuere, E., and Francart, T. (2019). Neural tracking of the speech envelope in cochlear implant users. *Journal of Neural
853 Engineering*, 16(1).

854 Steinschneider, M., Liégeois-Chauvel, C., and Brugge, J. F. (2011). *Auditory Evoked Potentials and Their Utility in the Assessment of Complex
855 Sound Processing*, pages 535–559. Springer US, Boston, MA.

856 Theunissen, F. E., Sen, K., and Doupe, A. J. (2000). Spectral-temporal receptive fields of nonlinear auditory neurons obtained using natural sounds.
857 *Journal of Neuroscience*, 20(6):2315–2331.

858 Tichko, P. and Skoe, E. (2017). Frequency-dependent fine structure in the frequency-following response: The byproduct of multiple generators.
859 *Hearing Research*, 348:1–15.

- 860 Van Canneyt, J., Wouters, J., and Francart, T. (2021a). Cortical compensation for hearing loss, but not age, in neural tracking of the fundamental
861 frequency of the voice. *Journal of Neurophysiology*, 126(3):791–802.
- 862 Van Canneyt, J., Wouters, J., and Francart, T. (2021b). Enhanced neural tracking of the fundamental frequency of the voice. *IEEE Transactions on*
863 *Biomedical Engineering (Early Access)*, x:1–1.
- 864 Van Canneyt, J., Wouters, J., and Francart, T. (2021c). Neural tracking of the fundamental frequency of the voice: the effect of voice characteristics.
865 *European Journal of Neuroscience*, 00(January):1–14.
- 866 Vanthornhout, J., Decruy, L., and Francart, T. (2019). Effect of Task and Attention on Neural Tracking of Speech. *Frontiers in Neuroscience*, 13.
- 867 Vanthornhout, J., Decruy, L., Wouters, J., Simon, J. Z., and Francart, T. (2018). Speech Intelligibility Predicted from Neural Entrainment of the
868 Speech Envelope. *JARO - Journal of the Association for Research in Otolaryngology*, 19(2):181–191.
- 869 Verschueren, E., Gillis, M., Decruy, L., Vanthornhout, J., and Francart, T. (2022). Speech understanding oppositely affects acoustic and linguistic
870 neural tracking in a speech rate manipulation paradigm. *bioRxiv*.
- 871 Verschueren, E., Somers, B., and Francart, T. (2019). Neural envelope tracking as a measure of speech understanding in cochlear implant users.
872 *Hearing Research*, 373:23–31.
- 873 Verschueren, E., Vanthornhout, J., and Francart, T. (2021). The effect of stimulus intensity on neural envelope tracking. *Hearing Research*,
874 403:108175.
- 875 Wang, X., Lu, T., Bendor, D., and Bartlett, E. (2008). Neural coding of temporal information in auditory thalamus and cortex. *Neuroscience*,
876 154(1):294–303.
- 877 Weissbart, H., Kandylaki, K. D., and Reichenbach, T. (2019). Cortical Tracking of Surprisal during Continuous Speech Comprehension. *Journal*
878 *of Cognitive Neuroscience*, pages 1–12.
- 879 Yang, M., Sheth, S. A., Schevon, C. A., II, G. M. M., and Mesgarani, N. (2015). Speech Reconstruction from Human Auditory Cortex with Deep
880 Neural Networks. *Interspeech*, page 5.
- 881 Zan, P., Presacco, A., Anderson, S., and Simon, J. Z. (2020). Exaggerated cortical representation of speech in older listeners: mutual information
882 analysis. *Journal of Neurophysiology*, 124(4):1152–1164.
- 883 Zuk, N. J., Murphy, J. W., Reilly, R. B., and Lalor, E. C. (2021). Envelope reconstruction of speech and music highlights stronger tracking of
884 speech at low frequencies. *PLOS Computational Biology*, 17(9):e1009358.