

1 **Title:** Early visual cortex tracks speech envelope in the absence of visual input

2

3

4 **Authors:** Evgenia Bednaya^{a*}, Bojana Mirkovic^b, Martina Berto^a, Emiliano
5 Ricciardi^a, Alice Martinelli^a, Alessandra Federici^a, Stefan Debener^{b,c}, and Davide
6 Bottari^a

7

8

9 **Affiliations:**

10 ^a Molecular Mind Laboratory, IMT School for Advanced Studies Lucca, Lucca, Italy

11 ^b Neuropsychology Lab, Department of Psychology, University of Oldenburg, Oldenburg, Germany

12 ^c Cluster of Excellence Hearing4all, University of Oldenburg, Oldenburg, Germany

13

14

15 ***Correspondence:**

16 Evgenia Bednaya

17 evgenia.bednaya@imtlucca.it

18 **Abstract**

19 Neural entrainment to continuous speech is typically observed within the language network
20 and can be modulated by both low-level acoustic features and high-level meaningful
21 linguistic units (e.g., phonemes, phrases, and sentences). Recent evidence showed that
22 visual cortex may entrain to speech envelope, however its putative role in the hierarchy of
23 speech processing remains unknown. We tested blindfolded participants who listened to
24 semantically meaningful or meaningless stories, either in quiet or embedded in multi-talker
25 babble noise. Entrainment to speech was assessed with forward linear modeling of
26 participants' EEG activity. We investigated (1) low-level acoustic effects by contrasting
27 neural tracking of speech presented in quiet or noise and (2) high-level linguistic effects by
28 contrasting neural tracking to meaningful or meaningless stories. Results showed that
29 envelope tracking was enhanced and delayed for speech embedded in noise compared to
30 quiet. When semantic information was missing, entrainment to speech envelope was
31 fastened and reduced. Source modeling revealed that envelope tracking engaged wide
32 neural networks beyond the auditory cortex, including early visual cortex. Surprisingly, while
33 no clear influence of semantic content was found, the magnitude of visual cortex
34 entrainment was affected by low-level features. The decrease of sound SNR-level
35 dampened visual cortex tracking, suggesting an active suppressing mechanism in
36 challenging listening conditions. Altogether, these findings provide further evidence of a
37 functional role of early visual cortex in the entrainment to continuous speech.

38

39

40 **Keywords:** EEG, envelope tracking, hierarchical speech processing, TRF, visual cortex

41

42

43

44

45

46

47

48 **Introduction**

49 Neuronal populations developed the ability to synchronize their activity (through aligning
50 the phase) to temporal regularities of a continuous input (Lakatos et al., 2019; Obleser &
51 Kayser, 2019). This neural entrainment influences several aspects of processing,
52 including language. In this context, neural activity entrained to amplitude modulations over
53 time of continuous speech (that is, the envelope) has been consistently reported (Ding &
54 Simon, 2014). The exact functional meaning of the entrainment to the speech envelope is
55 still unclear. Several studies showed that intelligible speech is not mandatory for neural
56 tracking (Howard & Poeppel, 2010; Luo & Poeppel, 2007). However, during
57 comprehension, phase-locked responses to speech in the auditory cortex are enhanced
58 (Gross et al., 2013; Peelle et al., 2013). Moreover, entrainment to an attended speaker's
59 speech envelope in noisy environments appears to play a role in solving the so-called
60 "cocktail party" (Cherry, 1953) problem (Ding, Chatterjee, & Simon, 2014; Riecke,
61 Formisano, Sorger, Başkent, & Gaudrain, 2018). Based on this evidence, entrainment to
62 speech envelope may be involved in promoting the perception of linguistic information
63 (Poeppel & Assaneo, 2020) and facilitating speech comprehension (Ahissar et al., 2001;
64 Luo & Poeppel, 2007), especially in challenging acoustic environments (e.g., Kerlin,
65 Shahin, & Miller, 2010; Zion Golumbic et al., 2013). Importantly, neural entrainment to
66 temporal dynamics of speech is modulated by low-level acoustic features (Ding et al.,
67 2014) and high-level meaningful linguistic units, such as phonetic information, phrases,
68 and sentences (Di Liberto, O'Sullivan, & Lalor, 2015).

69 Neural entrainment does not only occur for the auditory input of speech (A. E. O'Sullivan,
70 Crosse, Liberto, Cheveigné, & Lalor, 2021; Plass, Brang, Suzuki, & Grabowecy, 2020).
71 Recent magnetoencephalography (MEG) studies revealed that the early visual areas
72 entrain even to silent lip movements (Bourguignon, Baart, Kapnoula, & Molinaro, 2018,
73 2020; Hauswald, Lithari, Collignon, Leonardelli, & Weisz, 2018). This neural tracking is
74 modulated by audiovisual congruences and boosts speech comprehension in noisy
75 conditions (Park, Kayser, Thut, & Gross, 2016). The contribution of visual cortices in
76 language processing is not limited to visual or audiovisual representations of spoken
77 language. There is scattered evidence that the early visual cortex is also active during
78 purely auditory stimulation (Brang et al., 2022; Petro, Paton, & Muckli, 2017; Vetter, Smith,
79 & Muckli, 2014) and while listening to spoken language (e.g., Martinelli et al., 2020;
80 Seydell-Greenwald, Wang, Newport, Bi, & Striem-Amit, 2021; Wolmetz, Poeppel, & Rapp,
81 2011). Importantly, such activations cannot be explained by semantic-based imagery
82 alone but rather seem to reflect genuine responses to language input; in fact, the visual
83 cortex also responds to abstract concepts with low imaginability rates (Seydell-Greenwald
84 et al., 2021). Overall, this evidence highlights a putative role of the visual cortex in mapping

85 temporal modulations of incoming sounds, especially in the absence of competing retinal
86 input (Martinelli et al., 2020; Vetter et al., 2014). However, the exact role of the visual
87 cortex in the hierarchy of speech processing remains unclear.

88 Here, we investigated the neural tracking of speech envelope when visual input is absent.
89 Using electroencephalography (EEG), we recorded neural responses of blindfolded
90 individuals while they were listening to stories presented in isolation (Quiet) or combined
91 with multi-talker babble noise at different signal-to-noise ratios (SNR; Noise). Stories
92 comprised either meaningful (speech) or meaningless (jabberwocky) narration. We used
93 a temporal response function (TRF) to model neural tracking of broadband speech
94 envelope (in 2-8 Hz range; as in: Hausfeld, Riecke, Valente, & Formisano, 2018; Mirkovic,
95 Debener, Jaeger, & De Vos, 2015; J. A. O'Sullivan et al., 2015). TRF approach allows
96 linear mapping between neurophysiological responses and continuous speech stimuli
97 (Crosse, Di Liberto, Bednar, & Lalor, 2016; Crosse et al., 2021) and has been used to
98 measure entrainment to speech in both clear and challenging listening conditions (e.g.,
99 Decruy, Vanthornhout, & Francart, 2019; Di Liberto et al., 2015; Ding et al., 2014; Ding,
100 Melloni, Zhang, Tian, & Poeppel, 2016; Ding & Simon, 2014; Legendre, Andrillon, Koroma,
101 & Kouider, 2019; J. A. O'Sullivan et al., 2015).

102 To disambiguate the effects of lower-level acoustic and higher-level linguistic processing
103 using continuous naturalistic stimuli, we built a hierarchical model. We specifically
104 assessed the effects of (i) low-level acoustic features by contrasting TRFs resulting from
105 listening to stories presented in quiet vs. in noise, and (ii) high-level linguistic information
106 by contrasting TRFs resulting from listening to meaningful (speech) vs. meaningless
107 (jabberwocky) stories, both embedded in noise. Finally, we tested how low-level and high-
108 level information effects are distributed at the source level, with a focus on whether and
109 how speech envelope information is mapped in the visual cortex in the absence of
110 competing visual information.

111

112

113 **Materials and Methods**

114 **Participants**

115 Nineteen native speakers of the Italian language took part in the study (N = 19; age:
116 median = 28; min = 22; max = 32; females = 12; all right-handed). We excluded one
117 participant because of an error in the presentation script during EEG acquisition and three
118 more participants due to their inability to complete the experiment, resulting in a final
119 sample of fifteen participants (N = 15; age: median = 28; min = 22; max = 30; females=
120 10). All participants self-reported the absence of any hearing problems and neurological
121 disorders. The experimental protocol was approved by the local ethics committee and

122 conducted following the Declaration of Helsinki. All participants were informed in advance
123 that they would be blindfolded during the experiment, signed written informed consent
124 prior to the study, and received monetary compensation for their participation.

125

126 **Stimuli**

127 We used two types of target stories: (i) meaningful (speech) and (ii) meaningless
128 (jabberwocky) narration. Meaningful stories were extracted from the fiction book for teens
129 *Polissena del Porcello* by (Pitzorno, 1993). Meaningless stories were extracted from the
130 books containing nonsense, metasemantic (jabberwocky) poems and texts: *Gnòsi delle*
131 *fànfole* by (Maraini, 2019) and *Esercizi di Stile* by (Queneau, 1947/1983). Note that
132 syntactic information is preserved in jabberwocky stories, whereas semantic information
133 is absent or significantly reduced.

134 Target stories were narrated by a trained Italian actress. We registered stories in a
135 soundproof booth, using a video camera with an external condenser microphone
136 (Olympus ME51S) at sampling frequency of 48000 kHz. To create stimuli for our EEG
137 experiment, we extracted the audio material from the recorded files and edited them in
138 Audacity® software (version 2.3.0, <https://www.audacityteam.org/>) and with a custom
139 code using Signal Processing toolbox incorporated in MATLAB (version R2018b, Natick,
140 Massachusetts: The MathWorks Inc.). Specifically, we: (i) inspected raw audio files for
141 pronunciation errors and long breaths, consequently removing them, (ii) downsampled
142 audio to 44100 Hz, set to 16-bit and converted from Stereo to Mono, (iii) truncated long
143 pauses and silent periods exceeding 0.5 s to 0.5 s, (iv) trimmed resulting files to the same
144 length (~ 15 min), (v) identified the noise floor of the frequencies comprising the noise via
145 “Get Noise Profile” feature and subsequently removed low-amplitude background noise
146 using the Noise Reduction built-in feature based on an algorithm using Fourier analysis,
147 (vi) normalized resulting files to the same common root-mean-square (RMS) value to
148 ensure no variation of loudness across stories. Natural variations of loudness within each
149 story were preserved.

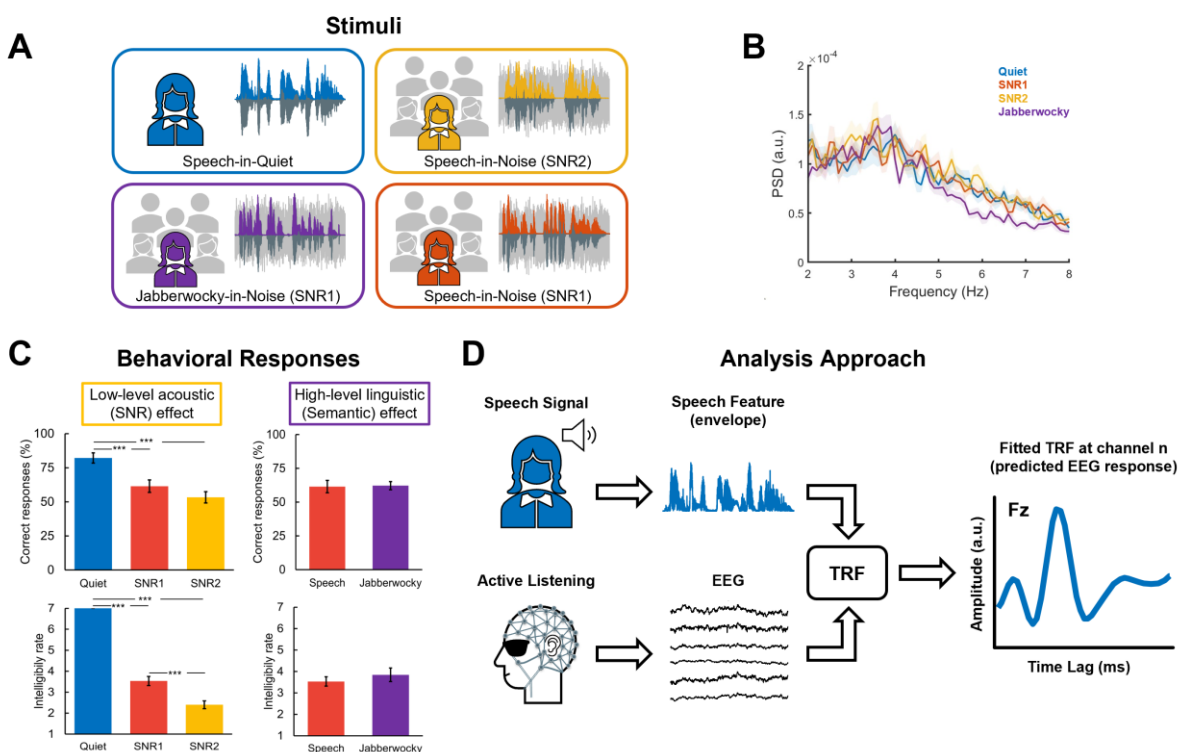
150 We combined the target stories with a five-talker babble to construct stimuli in which the
151 target story was embedded in the noise. Here, we used the babble noise, which is a non-
152 stationary noise that works well both as an energetic and informational masker, efficiently
153 reducing intelligibility and speech quality (Brungart, 2001; X. Wang & Xu, 2021). The
154 babble noise was a mixture of five different voices (2 females, 3 males, all native Italian
155 speakers). Every speaker was recorded in the soundproof booth, reading several non-
156 related extracts from the fiction book *La Strada* by (McCarthy, 2006/2014). These
157 individual recordings were registered and edited with the similar routine described above
158 for the target stimuli. Then, individual recordings were superimposed, resulting in multi-

159 talker babble. Finally, the initial 500 ms of the multi-talker babble got discarded to eliminate
 160 a part that did not contain all five talkers.

161 The first 5 s of the resulting multi-talker babble were set to zero/"muted," followed by 5 s
 162 of fade-in to make it easier for the participants to identify and track the target stories in the
 163 multi-talker babble noise. Then, with custom MATLAB scripts, we normalized the target
 164 stories and the babble to a common RMS value to make sure there would be no story or
 165 any of its segments standing out from the noise, and then superimposed the stories and
 166 the babble at two SNR levels (SNR1 = +3.52 dB, – for both meaningful and meaningless
 167 stories, and SNR2 = +1.74 dB, – for meaningful story only; see supplementary material
 168 for more details). As the last step, we normalized all the resulting audio files for all
 169 conditions once again to a common RMS value to achieve equal loudness across the
 170 stimuli and consequently verified each file's spectrogram in Audacity.

171 Altogether, we constructed stimuli to generate four experimental conditions: 1) *Speech-*
 172 *in-Quiet*, 2) *Speech-in-Noise at SNR1*, 3) *Speech-in-Noise at SNR2*, and 4) *Jabberwocky-*
 173 *in-Noise at SNR1* (Figure 1A). Each experimental condition contained a particular story
 174 divided into three parts of ~ 5 min, therefore the total duration of continuous speech stimuli
 175 per condition was ~ 15 min.

176 To test the effect of low-level acoustic (SNR) information, we compared neural tracking in
 177 *Speech-in-Quiet* condition and *Speech-in-Noise at SNR2* condition. To test the effect of
 178 high-level linguistic (semantic) information, we compared neural tracking in *Speech-in-*
 179 *Noise at SNR1* condition and *Jabberwocky-in-Noise at SNR1* condition.



180

181 **Figure 1: Stimuli, Behavioral Responses, and Analysis Approach. (A)** Stimuli
182 consisted of continuous (i) meaningful (Speech) and (ii) meaningless (Jabberwocky)
183 stories presented either in quiet (Quiet) or as embedded in the multi-talker babble noise
184 at a different signal-to-noise ratio (SNR1; SNR2). The babble noise was a mixture of five
185 voices (2 females, 3 males) reading extracts from a book. The acoustic envelopes were
186 extracted for further analysis through the Hilbert transform and filtering in the range
187 between 2 and 8 Hz. **(B)** Power spectra density estimates of normalized acoustic
188 envelopes were obtained using Welch's method with a 10 s Hamming window and half-
189 overlap. Bold lines indicate average across trials; shaded areas indicate standard error of
190 the mean. **(C)** Behavioral Responses represented by correct responses (Top) and
191 intelligibility rates (Bottom). Barplots display mean \pm SE across participants. Asterisks
192 indicate statistically significant differences ($***p < 0.001$). **(D)** Neural tracking of the speech
193 envelope was estimated using the forward encoding approach – Temporal Response
194 Function (TRF). Ridge regression-based linear models (TRFs) were fitted to participants'
195 neural data, obtained during active listening, to predict EEG response for of a given EEG
196 channel from speech envelope.

197

198 **Task and Experimental Procedure**

199 Participants performed four blocks, each consisted of one experimental condition. During
200 the first block, they always listened to the story without background noise (i.e., *Speech-*
201 *in-Quiet* condition). This was done to help the participants habituate both to the (target)
202 narrator's voice and the experimental design since this condition was the easiest to attend.
203 The order of the remaining three blocks was randomized across participants. Each of the
204 four blocks consisted of a story that lasted ~ 15 min divided into three parts ~ 5 min (see
205 supplementary material for further details). Participants listened to each part of the story
206 only once, without repetition, therefore avoiding the possibility of predicting the content of
207 the story. To maintain the continuity of the storyline within each block, each part within
208 each story followed the previous part chronologically.

209 We instructed participants to attentively listen to the target story (narrated by the female
210 voice and guided by the first 5 s of the audio) while ignoring babble noise in the
211 background. To ensure that the participants were actively attending to the stimuli, at the
212 end of each part, they answered three specific Yes/No questions about the part of the
213 story that they just listened to (for example, "*Il cane di Lucrezia è un San Bernardo?* [Is
214 Lucrezia's dog a Saint Bernard?]; see supplementary material for the full list of
215 questions). If they were not sure about the correct answer between the two, they had to
216 choose the answer that seemed to them the most probable. To answer, participants
217 pressed corresponding buttons on the response panel with their index and middle fingers.

218 At the end of each part, we asked participants to self-report intelligibility rates of the target
219 story on a Likert scale (where 1 – absolutely non-intelligible, 7 - very intelligible) and let
220 them have a short break lasting ~ 2 min. We also ensured that none of the participants
221 was familiar with or recently exposed to the target stories. Moreover, we informally
222 assessed a participant's comfort, alertness, and motivation to continue the experiment
223 during short and long breaks. We removed the blindfolding mask during the breaks
224 between each blocks (every 15 minutes) for the participants' comfort and in order to avoid
225 inducing short term cross-modal plasticity effects resulting from the prolonged visual
226 deprivation (Landry, Shiller, & Champoux, 2013; Lazzouni, Voss, & Lepore, 2012; Merabet
227 et al., 2008).

228 The experiment was controlled with E-Prime® software (version 3.0, Schneider et al.,
229 2002). All instructions and speech stimuli were presented through a single front-facing
230 loudspeaker (Bose Companion® series III multimedia speaker system, USA) placed in
231 front of the participants at approximately 1 m distance from their heads. Stimuli were
232 delivered at ~ 60 dB sound pressure level (SPL), measured at the participant's ear, and
233 reported by all the participants as comfortable volume.

234 To accurately measure the actual onset time of our stimuli, we administered a timing-test
235 using Audio/Visual (AV) Device (Electrical Geodesics, Inc.) compatible with E-Prime
236 software and NetStation system. The measured average delay in time was constant and
237 about + 5 ms regarding the stimulus onset.

238

239 **EEG Recording**

240 Before starting the experiment, each participant received a brief instruction and had a
241 short (~ 1 min) "training" session on how to control over muscular artifacts through
242 monitoring their EEG signal displayed on the computer screen. Then, we applied the
243 blindfolding mask to the participant, and they were reminded to keep their eyes open
244 during the EEG recordings, though blinking was permitted whenever they wanted.
245 Moreover, we recorded resting-state EEG data for about 2 minutes at the beginning of
246 each experiment while the participant kept their eyes open. Obtained resting-state data
247 served as calibration data to attenuate EEG artifacts during the preprocessing step.

248 During the tasks, the participants were seated comfortably in a chair in a dark,
249 soundproofed booth (BOXY, B-Beng s.r.l., Italy). The EEG recordings were acquired at a
250 sampling rate of 500 Hz using NetStation5 software together with a Net Amps 400 EGI
251 amplifier connected to 64 electrodes HydroCel Geodesic Sensor Net (Electrical
252 Geodesics, Inc.), all signals referenced to vertex (additional channel E65/Cz). For data
253 visualization purposes only, the data were band-pass filtered online using the digital filter
254 from 1.0 to 100 Hz, and online digital anti-alias filter aligning EEG recordings with real-

255 time events was kept on. Electrode impedances were kept below 50 k Ω and were checked
256 between the blocks (when the blindfolding mask was reapplied).

257 Participants were encouraged to take a break after each block and get enough rest before
258 continuing. They also were reminded about the importance of staying attentive, keeping
259 eyes open while blindfolded, and avoiding excessive movements during the EEG
260 recordings.

261

262 **EEG Preprocessing**

263 We preprocessed continuous EEG raw data offline using custom MATLAB (version
264 R2018b, Mathworks Inc., Natick, MA) scripts together with EEGLAB toolbox (version
265 14.1.2b, Delorme & Makeig, 2004) for MATLAB.

266 First, the EEG data were submitted to cleaning with Artifact Subspace Reconstruction
267 (ASR) - an automated artifact attenuation algorithm (clean_rawdata plug-in, version 2.1)
268 for EEGLAB toolbox. We applied the default flatline criterion of 5 s, together with default
269 transition band parameters ([0.25 0.75]). ASR algorithm was chosen due to its objective
270 and reproducible evaluation of artifactual components in EEG data. ASR is based on
271 Principal Component Analysis (PCA) sliding window and effectively attenuates high-
272 variance signal components in the EEG data (including eye blinks, eye movements, and
273 motion artifacts). Specifically, first, the algorithm automatically identifies the most artifact-
274 free part of the data (here, the resting-state data) to use it as the calibration data to
275 compute the statistics. Next, a 500 ms PCA sliding window with 50% overlap is applied
276 across all the channels to identify "bad" principal components. Then, the algorithm
277 identifies the subspaces in which the signal exceeds 5 standard deviations away from the
278 calibration data as corrupted and rejects them. Finally, it reconstructs the high variance
279 subspaces using a mixing matrix calculated based on the calibration data.

280 The artifact attenuated EEG data were preprocessed as follows: (i) re-referenced from
281 E65/Cz electrode to a common average reference, (ii) band-pass filtered from 0.1 to 40
282 Hz (low-pass: FIR filter, filter order: 100, window type: Hann; high-pass: FIR filter, filter
283 order: 500, window type: Hann), (iii) downsampled to 250 Hz, (iv) epoched according to
284 the onset of acoustic stimuli (related to each part of the story), adjusting to measured +5
285 ms onset delay in time and discarding the first 5 s of target-speech alone and 5 s of fade-
286 in for the babble noise, (v) band-pass filtered between 2 and 8 Hz (filter type and
287 parameters the same as described above), (vi) downsampled to 64 Hz, (vii) EEG data
288 corresponding to each of the three ~ 5 min parts of the story were concatenated, (viii) and
289 segmented into 1 min long trials, resulting in 12 trials per block per subject (N = 12). The
290 preprocessed EEG data for each trial were z-scored to optimize cross-validation
291 procedure during encoding (Crosse et al., 2016).

292

293 **Extraction of Acoustic Envelope**

294 First, audio files containing relevant parts of the target stories were concatenated and
295 segmented into corresponding 1 min long trials, resulting in 12 trials per speech envelope
296 per subject (N = 12) (Figure 1B). Next, the acoustic envelope per each trial was obtained
297 taking the absolute value of the Hilbert transform of the original target stories (i.e., without
298 babble noise) followed by a low-pass filtering using a 3rd-order Butterworth filter with a
299 cut-off frequency of 8 Hz (filtfilt function in MATLAB) and downsampling the resulting
300 signal to 64 Hz, so to be matched with the EEG data (e.g., Mirkovic et al., 2015; J. A.
301 O'Sullivan et al., 2015). Finally, the resultant extracted envelopes were normalized by
302 dividing by maximum value.

303

304 **Estimation of TRF**

305 We modeled where and how the neural response to the speech envelope of the target
306 stories is encoded in the brain, using a linear prediction approach known as temporal
307 response function (TRF) (Figure 1D). The TRF approach, incorporated in mTRF toolbox
308 (Crosse et al., 2016), allows to predict previously unseen EEG response from the stimulus
309 and has been used to model the neural tracking of acoustic and linguistic properties of
310 naturalistic continuous speech (Drennan & Lalor, 2019; Obleser & Kayser, 2019).

311 The TRF is a mathematical function that is based on the ridge regression and could be
312 described as follows:

$$313 \quad r(t, n) = \sum_{\tau} w(\tau, n) s(t - \tau) + \varepsilon(t, n),$$

314 where $t = 0, 1, \dots, T$ is time, $r(t, n)$ is the EEG response from an individual channel, $s(t)$
315 is the stimulus feature(s) (e.g., speech envelope), τ is the range of time-lags between s
316 and r , $w(\tau, n)$ are the regression weights over time-lags, and $\varepsilon(t)$ is a residual response
317 at each channel not explained by the TRF model (Crosse et al., 2016). Specifically, TRF
318 can be viewed as a filter that describes the linear relationship between a continuous
319 speech stimulus and a continuous neural response for a specified range of time-lags
320 related to stimulus occurrence (Crosse et al., 2016).

321 The important assumptions about the TRF include the fact that it reflects the same neural
322 generators as cortical auditory evoked potentials (CAEPs) resulting in their comparable
323 topographies and that it can be used to measure neural tracking of speech envelope (Lalor
324 & Foxe, 2010; Lalor, Power, Reilly, & Foxe, 2009). We fitted separate models (TRFs) to
325 predict response in each EEG channel, using time-lags from -100 to 600 ms related to
326 stimulus onset, typically used to capture CAEP components. Here we estimated the TRF

327 using the envelope estimated between 2 and 8 Hz as previously performed (Legendre et
328 al., 2019; Mirkovic et al., 2015; J. A. O'Sullivan et al., 2015).

329 The TRF models were trained using a leave-one-out cross-validation procedure, keeping
330 all but one trial for training the model to predict EEG response from the stimuli and using
331 a left-out trial for testing. Thus, a prediction model was obtained for every single trial, and
332 then the final averaging across trials, within participants and conditions was performed,
333 resulting in a grand average TRF model.

334

335 **Regularization Parameter Estimation**

336 Regression models are exposed to overfitting the training data, that is, fitting the random
337 noise rather than true relationships between variables and failing to generalize to unseen
338 data. The problem of overfitting needs to be accounted for before making any
339 interpretations from the resulting model since it could be misleading. Ridge regularization
340 prevents the model from overfitting by penalizing the model weights, forcing them to be
341 smaller, towards 0, so the model could become better generalized.

342 To control for model overfitting, we empirically identified the optimal regularization
343 parameter (λ) of TRF models through leave-one-out cross-validation procedure, using a
344 grid of ridge values ($\lambda = \{10^{-6}, 10^{-5}, \dots, 1, 10, \dots, 10^5, 10^6\}$), for time-lags from -100
345 to 600 ms. The regularization parameter λ was determined based on the mean squared
346 error (MSE) value between the actual and predicted EEG responses. The optimal
347 regularization parameter was the one yielding the lowest MSE on the testing data (here,
348 identified as $\lambda=10^3$) and kept constant across channels, participants, and conditions
349 allowing to generalize across them at the group level.

350

351 **Spatiotemporal Characteristics**

352 Forward model weights are directly physiologically interpretable (Haufe et al., 2014) and
353 allow us to get an insight about which channels contribute most to neural tracking of the
354 speech envelope. The resulting topographical plots with TRF weights obtained per each
355 individual time-lag window can be interpreted similarly to CAEPs in terms of both
356 amplitude and direction (Lalor, Pearlmutter, Reilly, McDarby, & Foxe, 2006; Lalor et al.,
357 2009). We investigated spatiotemporal characteristics of forward model weights by fitting
358 the TRFs at different individual time-lags between the EEG response and the speech
359 envelopes, using a sliding time-lag window of 45 with 30 ms overlap in a time-lag range
360 from -115 to 620 ms. Finally, the estimate of forward model weights allowed us to directly
361 transfer the data into source space avoiding further transformations (Haufe et al., 2014).

362

363 **Chance-level Estimation by Permutation Testing (Control)**

364 To assess the ability of TRF models to predict neural responses (i.e., neural tracking) and
365 verify that neural tracking was well above chance, we computed "null distributed" TRF
366 model (Combrisson & Jerbi, 2015). We used a permutation-based approach with
367 *mTRFpermute* function, incorporated in mTRF-toolbox (Crosse et al., 2016, 2021).
368 Specifically, this approach cross-validates models, iteratively (1000 iterations) fitting TRFs
369 on randomly mismatched pairings of speech envelopes/EEG responses and evaluating
370 the models on matched data. This procedure was done separately for each trial,
371 participant, and condition, and then grand averaged to get the average "null" TRF model,
372 which served as a baseline ("control").

373

374 **Source Estimation**

375 Forward modeling allowed us to investigate the TRFs and better understand how the
376 information about the envelope of continuous stimuli is encoded in the brain. Specifically,
377 we tested how low-level (SNR) acoustic and high-level linguistic (Semantic) effects are
378 distributed at sensor and source levels. Furthermore, we investigated whether and how
379 the visual cortex is activated for neural tracking of the speech envelope in blindfolded
380 individuals when competing retinal input is absent.

381 We performed source localization in Brainstorm software (Tadel, Baillet, Mosher,
382 Pantazis, & Leahy, 2011) together with custom MATLAB scripts and the pipeline for EEG
383 source estimation introduced by Stropahl and colleagues (2018; see also Bottari et al.,
384 2020) that we adapted to the TRF data. Specifically, source localization was performed
385 using dynamical Statistical Parametric Mapping (dSPM, Dale et al., 2000). A Boundary
386 Element Model (BEM) was computed for each participant using default parameters to
387 calculate the forward solution and constrain source locations to the cortical surface. We
388 used a standard electrode layout together with a standard anatomy template (ICBM152)
389 for all participants. The model resulted in a single dipole oriented perpendicularly to the
390 cortical surface for each vertex since dipole orientations were constrained to the cortical
391 surface. We did not perform individual noise modeling since TRF has no clear nor true
392 baseline period. Instead, we used an identity matrix as a noise covariance matrix, with the
393 assumption of equal unit variance of noise on every sensor.

394 We created visual regions of interest (ROIs) based on predefined scouts from the
395 Destrieux atlas (Destrieux, Fischl, Dale, & Halgren, 2010) implemented in FreeSurfer
396 (Fischl, 2012) and available in Brainstorm. Visual ROIs were selected for the left and right
397 hemispheres and included primary (V1; Calcarine sulcus) and secondary (V2, Lingual
398 gyrus) visual cortex, defined as the '*S_calcarine*' and the '*G_oc-tem_med-Lingual*' scouts
399 in the atlas, correspondingly. These visual ROIs were selected based on recently reported
400 evidence of their involvement in speech processing not only in blind but also sighted

401 individuals, albeit to a lower extent (Martinelli et al., 2020; Petro et al., 2017; Seydell-
402 Greenwald et al., 2021; Van Ackeren, Barbero, Mattioni, Bottini, & Collignon, 2018; Vetter
403 et al., 2020, 2014). Upon the ROIs creation, their time-series were extracted and
404 submitted to the analysis.

405

406 **Statistical Analysis**

407 Participants' behavioral responses concerning comprehension of the story were computed
408 as the average correct responses (in %) across all three parts of the story, resulting in
409 nine scores per participant for each condition. Intelligibility rates from each participant
410 were computed similarly, by averaging across all three parts of the story. Statistical
411 analysis of behavioral responses to assess low-level acoustic (SNR) effect was conducted
412 using one-way repeated measure ANOVA. Post-hoc comparisons were made with two-
413 tailed paired t-tests. Statistical analysis of behavioral responses to assess high-level
414 linguistic (semantic) effect was performed using two-tailed paired t-tests.

415 As a sanity check, we first performed comparisons between the TRFs of each condition
416 with the "null" TRF through paired t-tests, with the significance threshold set at $p < 0.05$
417 (one-tailed) and corrected for multiple comparisons with the false-discovery rate (FDR) at
418 0.05 (Benjamini & Hochberg, 1995), at two electrodes selected a priori on the midline
419 frontocentral (Fz) and the occipital (Oz) scalp locations, over a range of post-stimulus
420 time-lags between 0 and 600 ms

421 To access differences in the spatiotemporal profile of averaged TRFs between conditions,
422 we performed non-parametric cluster-based permutation tests (Maris & Oostenveld, 2007)
423 in FieldTrip toolbox (Oostenveld, Fries, Maris, & Schoffelen, 2011). A cluster was defined
424 along *electrodes x time-lags* dimensions, with extension criteria set to at least two
425 neighboring electrodes. The t-statistic for adjusted *electrode x time-lag* pairs exceeding a
426 preset critical threshold of 5% (cluster alpha = 0.05) was summed, and the adjusted pairs
427 formed the clusters. Then, two-tailed tests were performed at the whole brain level (across
428 all electrodes and time-lags from 0 to 600 ms), using the Monte-Carlo method with 1000
429 permutations. The maximum of the summed t-statistic in the observed data was compared
430 with a random partition formed by permuting the experimental condition labels, resulting
431 in a critical p-value for each cluster. In case the cluster-based p-value was less than 0.025
432 (corresponding to a critical alpha level of 0.05 for two-tailed testing, accounting for both
433 positive and negative clusters), we rejected our null hypothesis that there were no
434 differences between TRFs for two conditions.

435 Finally, cluster-based statistics on sources at the whole-brain level were performed in
436 Brainstorm, across all electrodes and time-lags from 0 to 600 ms, using Monte-Carlo
437 method with 1000 permutations, alpha = 0.05, two-tailed (meaning alpha = 0.025 per each

438 tail), cluster alpha = 0.05, and neighboring criteria for electrodes set for 2. Analysis of
439 visual ROIs time-series between conditions was performed using paired t-tests, with the
440 significance threshold set at $p < 0.05$ (one-tailed) and correcting for multiple comparisons
441 with the FDR-method at 0.05.

442

443

444 **Results**

445 **Behavioral responses**

446 To ensure that the participants successfully understood the content of the target stories,
447 they were asked to answer three Yes/No questions at the end of each segment (5
448 minutes). Moreover, participants were asked to self-rate the intelligibility of each part of
449 the target story from 1 (absolutely non-intelligible) to 7 (very intelligible).

450

451 *Low-level acoustic (SNR) effect*

452 As expected, both comprehension scores and intelligibility rates gradually decreased with
453 SNR (Figure 1C). Comprehension scores, converted to percentage of correct responses,
454 decreased as a function of noise (*Speech-in-Quiet* mean \pm SE: $82.22 \pm 3.72\%$; *Speech-*
455 *in-Noise at SNR1* mean \pm SE: $61.48 \pm 4.58\%$; *Speech-in-Noise at SNR2* mean \pm SE:
456 $53.33 \pm 4.09\%$). A repeated measures ANOVA with a correction confirmed that listening
457 condition significantly affected participants' comprehension ($F(2, 28) = 16.14$, $p = 0.00002$,
458 Huynh-Feldt corrected). Post-hoc comparisons showed that correct responses for
459 *Speech-in-Quiet* were significantly higher than for *Speech-in-Noise at SNR1* ($t(14) = 4.40$,
460 $p = 0.0006$) and *Speech-in-Noise at SNR2* ($t(14) = 5.46$, $p = 0.0001$), but no significant
461 difference emerged between *Speech-in-Noise at SNR1* and *Speech-in-Noise at SNR2*
462 ($t(14) = 1.43$, $p = 0.17$).

463 Intelligibility rates were in line with comprehension scores, dramatically dropping from
464 *Speech-in-Quiet* (rated 7 by all participants, and thus reaching a ceiling which prevented
465 comparisons with other conditions; see Liu & Wang, 2021; Šimkovic & Träuble, 2019) to
466 *Speech-in-Noise at SNR1* (mean \pm SE: 3.53 ± 0.22) and further significantly dropping at
467 *Speech-in-Noise at SNR2* (mean \pm SE: 2.40 ± 0.19 ; *Speech-in-Noise at SNR1* vs.
468 *Speech-in-Noise at SNR2*: $t(14) = 5.90$, $p < 0.0001$).

469

470 *High-level linguistic (Semantic) effect*

471 We found no difference in correct responses and intelligibility rates between *Speech-in-*
472 *Noise at SNR1* and *Jabberwocky-in-Noise at SNR1* (all p-values > 0.05 ; correct
473 responses, mean \pm SE: *Speech-in-Noise at SNR1*: $61.48 \pm 4.58\%$; *Jabberwocky-in-Noise*

474 at SNR1: $62.22 \pm 3.03\%$; intelligibility rates, mean \pm SE: *Speech-in-Noise at SNR1*: 3.53
475 ± 0.22 ; *Jabberwocky-in-Noise at SNR1*: 3.84 ± 0.32 . Results indicated that participants
476 were able to equally attend target stories embedded in noise (SNR1), regardless of
477 semantic information.

478

479 **Neural tracking**

480 *Low-level acoustic (SNR) effect at the sensor level*

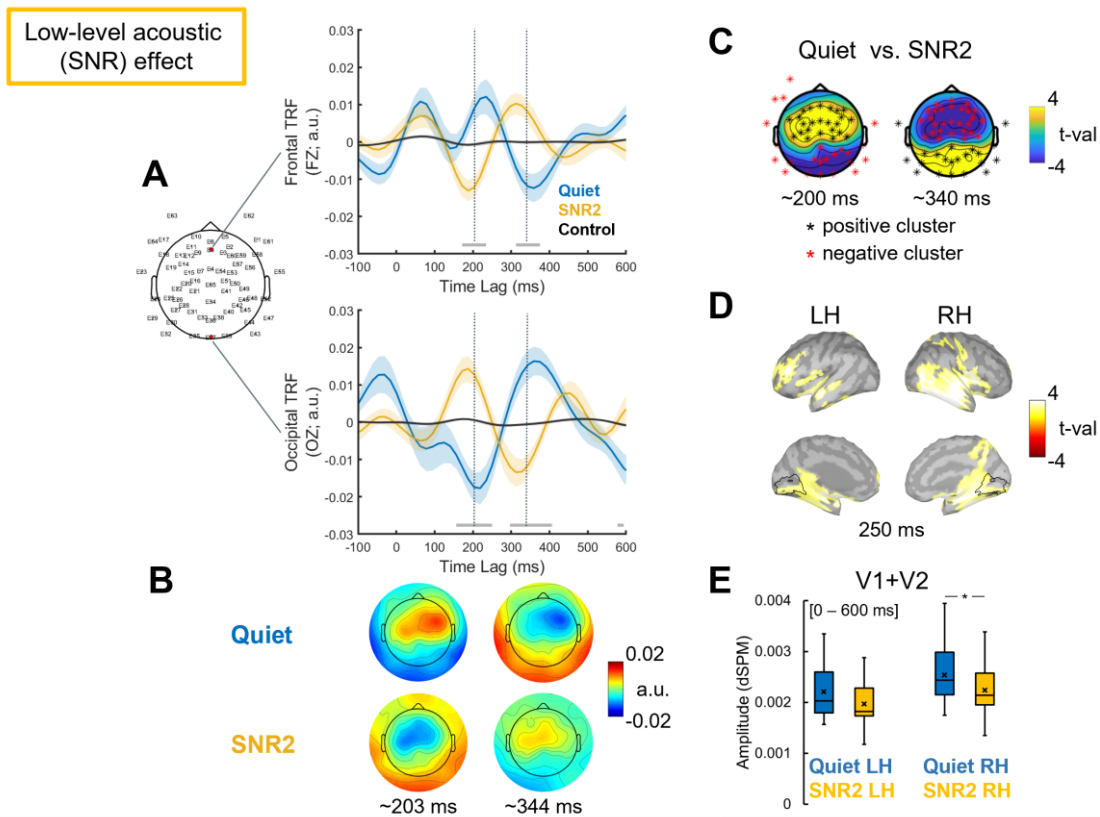
481 First, we examined the temporal profile of SNR effect at preselected representative
482 electrodes: frontal (Fz) and occipital (Oz; Figure 2A). The TRFs for *Speech-in-Quiet* and
483 *Speech-in-Noise at SNR2* were significantly different from the "null" TRF ($p < 0.05$, one-
484 tailed, FDR-corrected), suggesting that TRFs indeed reflected neural tracking of the
485 speech envelope (Supplementary Figure S1).

486 To access the effect of SNR on neural tracking of the speech envelope, we compared the
487 TRFs of *Speech-in-Quiet* and *Speech-in-Noise at SNR2* (the most challenging) conditions
488 (Figure 2). The Cluster-based permutation test revealed significant differences between
489 the TRFs for the two conditions ($p < 0.025$; cluster-corrected). A positive ($p = 0.002$,
490 corrected) and a negative ($p = 0.002$, corrected) clusters were identified at time-lags
491 interval 150 – 250 ms. Other pair of positive ($p = 0.001$, corrected) and negative clusters
492 ($p = 0.002$, corrected) were also found at time-lags interval 290 – 410 ms (Figure 2C).
493 Both effects extended over fronto-central and parieto-occipital electrodes. Results showed
494 that TRF to *Speech-in-Noise at SNR2* was delayed and increased in magnitude compared
495 to *Speech-in-Quite condition* (Figure 2A, B and C).

496

497 *Low-level acoustic (SNR) effect at the source level*

498 The cluster-based permutation test, performed at the whole brain level, contrasting TRFs
499 for *Speech-in-Quiet* vs. *Speech-in-Noise at SNR2*, revealed that SNR effect was localized
500 in both hemispheres (Figure 2D): a significant cluster was found in the left hemisphere (p
501 $= 0.008$, corrected), lasting from ~ 0 to ~ 484 ms, and another one in the right hemisphere
502 ($p = 0.028$, corrected), lasting from ~ 141 to ~ 312 ms (see Supplementary Figure S2).
503 The effect was observed mostly over bilateral temporal cortex, and also included parts of
504 the bilateral parietal cortex, insular cortex, visual cortex, and left prefrontal cortex.



505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

Figure 2: Low-level acoustic (SNR) effect. (A) Grand averaged temporal response functions (TRFs) for *Speech-in-Quiet* (Quiet, blue), *Speech-in-Noise at SNR2* (SNR2, yellow), and "null" TRF (Control, black). TRFs displayed over time-lags at frontal Fz and occipital Oz electrodes, marked with red on the electrode layout. Shaded areas represent the standard error of the mean (SE) across participants. Grey horizontal bars above the x-axis indicate time-lags at which TRFs for *Speech-in-Quiet* and *Speech-in-Noise at SNR2* differed significantly at these representative electrodes (series of paired two-tailed t-tests, $p < 0.05$, FDR-corrected). Grey dotted vertical lines indicate time-lags with the maximal difference between TRFs for *Speech-in-Quiet* and *Speech-in-Noise at SNR2*. (B) Topographic representations of TRFs, displayed at time-lags marked by grey dotted vertical lines on A. (C) The results of the cluster-based permutation test contrasting TRFs for *Speech-in-Quiet* vs. *Speech-in-Noise at SNR2*, displayed around time-lags marked by grey dotted vertical lines on A. Significant ($p < 0.05$, corrected for two tails, $p < 0.025$ for each tail) positive and negative clusters comprised the electrodes marked in black and in red asterisks, respectively. (D) Differences at the source level, contrasting TRFs for *Speech-in-Quiet* vs. *Speech-in-Noise at SNR2* at the whole-brain level ($p < 0.05$, corrected for two tails). Lateral and medial views of the left (LH) and right (RH) hemispheres, displayed at the time-lag corresponding to the peak in the temporal profile (i.e., 250 ms). Bright yellow (positive t-values) indicates greater activation for Quiet over SNR2. Black contours indicate the ROIs borders (V1 and V2) in both hemispheres based

526 on the Destrieux cortical atlas. **(E)** Activations obtained at the source space in visual ROIs.
527 Boxplots display source activation for each condition, *averaged* over the ROIs (V1 + V2)
528 and across the 0 to 600 ms time window, in the left (LH) and right (RH) hemispheres,
529 respectively. The line through the boxplot indicates the median, × marker indicates the
530 mean, lines indicate pairwise statistical comparisons (* $p < 0.05$, one-tailed).

531

532 *Visual cortex ROIs*

533 To test whether and how the visual cortex is taking part in neural tracking of speech and
534 speech comprehension in blindfolded individuals, we performed source analysis on TRFs,
535 using predefined ROIs in the visual cortex comprising V1 and V2.

536 The contrast *Speech-in-Quiet vs. Speech-in-Noise at SNR2* survived cluster-correction
537 for multiple-comparisons in the left ($p = 0.008$, corrected) and right hemispheres ($p =$
538 0.028 , corrected; Supplementary Figure S3). Extracted time-series from V1 and V2
539 showed a similar pattern, with the magnitude of source activation for TRF in *Speech-in-*
540 *Quiet* being larger than for TRF in *Speech-in-Noise at SNR2* at multiple time points (see
541 Supplementary Figure S3 reporting uncorrected results). Averaged activation across time
542 points in combined ROIs (V1 + V2) was significantly larger for TRF in *Speech-in-Quiet*
543 than for TRF in *Speech-in-Noise at SNR2* in the right hemisphere ($p = 0.04$, one-tailed),
544 but not in the left hemisphere ($p = 0.08$, one-tailed) (Figure 2E). These results suggest the
545 dampening of visual cortex activity in case of challenging auditory inputs.

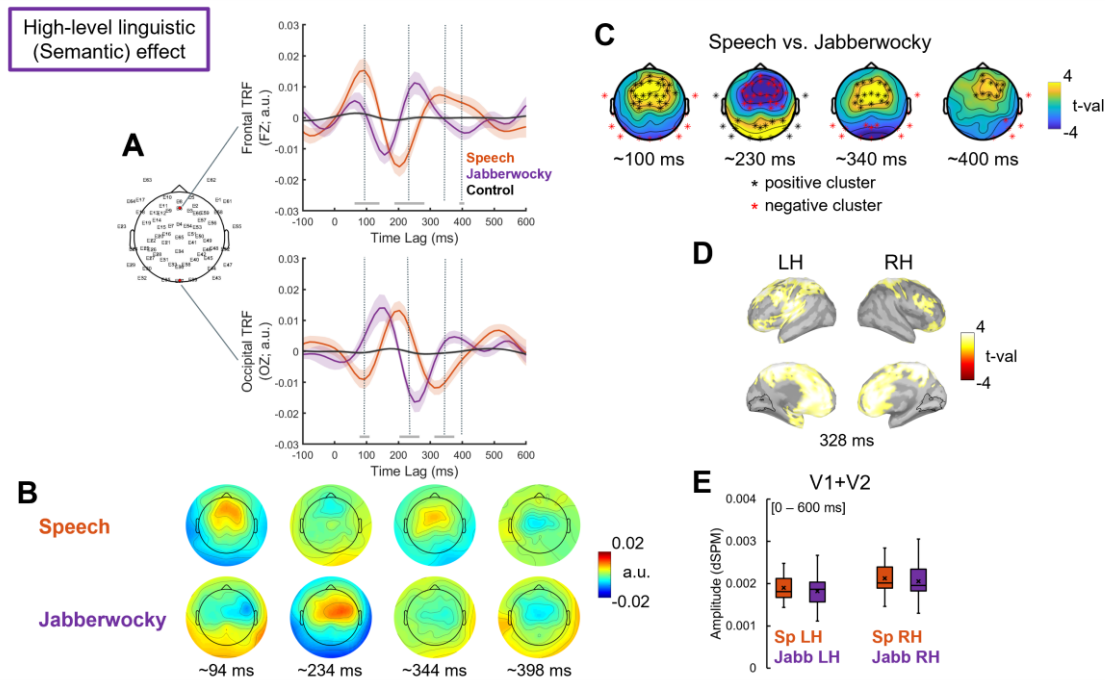
546

547 *High-level linguistic (Semantic) effect at the sensor level*

548 At the two electrodes of interest (The TRFs for *Speech-in-Noise at SNR1* and
549 *Jabberwocky-in-Noise at SNR1* significantly differed from the "null" TRF ($p < 0.05$, one-
550 tailed, FDR-corrected), suggesting that the estimated TRFs indeed reflected neural
551 tracking of the speech envelope (Supplementary Figure S1).

552 To access the effect of semantic information on neural tracking, we compared the TRFs
553 of *Speech-in-Noise at SNR1* and *Jabberwocky-in-Noise at SNR1* conditions (Figure 3).

554 Cluster-based permutation test on TRFs revealed statistically significant differences
555 between two conditions ($p < 0.025$; corrected). Three pairs of positive and negative
556 clusters were identified at time-lags intervals of 70 – 165 ms (positive: $p = 0.001$,
557 corrected; negative: $p = 0.01$, corrected), 200 – 290 ms (positive: $p = 0.001$, corrected;
558 negative: $p = 0.001$, corrected), and 310 – 430 ms (positive: $p = 0.003$, corrected;
559 negative: $p = 0.01$, corrected), comprising fronto-central electrodes and parieto-occipital
560 electrodes (Figure 3C). Results revealed that the TRFs of *Speech-in-Noise at SNR1* was
561 higher and delayed compared to the TRF of *Jabberwocky-in-Noise at SNR1* (see Figure
562 3A, B and C).



563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

Figure 3: High-level (Semantic) effect. (A) Grand averaged temporal response functions (TRFs) for *Speech-in-Noise at SNR1* (Speech, red), for *Jabberwocky-in-Noise at SNR1* (Jabberwocky, purple), and "null" TRF (Control, black). TRFs displayed over time-lags at frontal Fz and occipital Oz electrodes, marked with red on the electrode layout. Shaded areas represent the standard error of the mean (SE) across participants. Grey horizontal bars above the x-axis indicate time-lags at which TRFs for *Speech-in-Noise at SNR1* and *Jabberwocky-in-Noise at SNR1* differed significantly (running paired two-tailed t-tests, $p < 0.05$, FDR-corrected). Grey dotted vertical lines indicate time-lags with the maximal difference between TRFs for *Speech-in-Noise at SNR1* and *Jabberwocky-in-Noise at SNR1*. (B) Topographic representations of TRFs, displayed at time-lags marked by grey dotted vertical lines on A. (C) The results of the cluster-based permutation test contrasting TRFs for *Speech-in-Noise at SNR1* and *Jabberwocky-in-Noise at SNR1*, displayed around time-lags marked by grey dotted vertical lines on A. Significant ($p < 0.05$, cluster-corrected for two tails, meaning $p < 0.025$ each tail) positive and negative clusters comprised the electrodes marked in black and in red asterisks, respectively. (D) Differences at the source level, contrasting TRFs for *Speech-in-Noise at SNR1* and *Jabberwocky-in-Noise at SNR1* at the whole brain level ($p < 0.05$, cluster-corrected for two tails). Lateral and medial views of the left (LH) and right (RH) hemispheres, displayed at the time-lag corresponding to the peaks in the temporal profile. Bright yellow (positive t-values)/dark red (negative t-values) colors indicate greater activation for Speech/Jabberwocky, respectively. Black contours indicate the ROIs borders (union of V1 and V2) in both hemispheres based on the Destrieux cortical atlas. (E) Activations obtained at the source space in visual ROIs. Boxplots display source activation for each

587 condition, *averaged* over the ROIs (V1 + V2) and across the 0 to 600 ms time window in
588 the left (LH) and right (RH) hemispheres, respectively. The line through the boxplot
589 indicates the median, × marker indicates the mean.

590

591 *High-level linguistic (Semantic) effect at the source level*

592 Cluster-based permutation test, contrasting TRFs for *Speech-in-Noise at SNR1* and
593 *Jabberwocky-in-Noise at SNR1* at the whole-brain level, revealed two clusters in both
594 hemispheres: one in the left hemisphere ($p = 0.002$, corrected), extending over all time
595 points, and one in the right hemisphere ($p = 0.006$, corrected), lasting from ~0 to ~531 ms
596 (Supplementary Figure S2), with maximum activation ~ 330 ms (Figure 3D). The effect
597 extended primarily over the left auditory cortex and a large portion of the bilateral fronto-
598 parietal network at earlier time points and extended to the anterior temporal lobe (ATL) at
599 later time points (Supplementary Figure S2).

600

601 *Visual cortex ROIs*

602 In the visual ROI the Semantic effect did not survive cluster-correction for multiple-
603 comparisons ($p > 0.05$, corrected for two tails), and extracted time-series from ROIs did
604 not differ between source TRFs for *Speech-in-Noise at SNR1* and *Jabberwocky-in-Noise*
605 *at SNR1* ($p > 0.05$) (Supplementary Figure S3 and Figure 2E).

606

607

608 **Discussion**

609 We used a hierarchical model to investigate entrainment to continuous speech envelope
610 in blindfolded individuals, assessing (1) the effects of low-level acoustic and high-level
611 linguistic information on neural tracking and (2) testing how these effects are distributed
612 at the source level, with the focus on the visual cortex. To address the role of low-level
613 acoustic, we compared the entrainment to target stories presented in quiet or multi-talker
614 babble noise. Results revealed that TRF was delayed and higher in magnitude at latencies
615 between 100 and 300 ms when SNR decreased. This finding suggests that neural tracking
616 requires greater resources in case of concurrent masking noise. Next, we also addressed
617 the role of high (semantic) level of speech processing on neural entrainment by comparing
618 TRFs to meaningful and meaningless stories. Results indicated delayed and higher TRFs
619 when semantic information is present. Source modeling suggested that entrainment to
620 continuous speech in noise engaged a spread activation beyond the auditory cortex,
621 including linguistic and attentional networks. Finally, in the absence of retinal input, we
622 found evidence that the visual cortex entrained to the speech envelope. However, the
623 magnitude of such entrainment was degraded with concurrent background noise,

624 suggesting a suppressing mechanism helping to focus auditory attention in challenging
625 listening conditions.

626

627 **Effects of low-level acoustic (SNR) processing on neural tracking of speech** 628 **envelope**

629 We demonstrated that speech envelope tracking in noise, compared to quiet, was
630 characterized by larger amplitude and delayed latency of the TRF responses and by the
631 reversed polarity of the TRFs topography distributions over fronto-central parieto-occipital
632 electrodes (Figure 2A, B).

633 The TRF time-courses were consistent with previous studies reporting amplitudes and
634 latencies being affected by concurrent noise (Brodbeck, Jiao, Hong, & Simon, 2020; Ding
635 & Simon, 2013; Fiedler, Wöstmann, Herbst, & Obleser, 2019; Gustafson, Billings,
636 Hornsby, & Key, 2019; Zendel, West, Belleville, & Peretz, 2019) as well as enhanced N1
637 and N2 amplitudes in noise compared to quiet (Papesh, Billings, & Baltzell, 2015).

638 Increased frontal negativity around 100 ms (N1) is associated with attention-dependent
639 processes in response to auditory changes (Hansen & Hillyard, 1980; Näätänen, 1982).

640 The enhanced envelope tracking observed here for the N1-like response to speech in
641 noise compared to quiet may reflect the use of more resources for the encoding of acoustic
642 variations at earlier stages of speech processing when intelligibility gets degraded by
643 noise (Alain, Quan, McDonald, & Van Roon, 2009; Näätänen & Picton, 1987; Parbery-
644 Clark, Marmel, Bair, & Kraus, 2011).

645 Additional differences were observed around the second negative peak, corresponding to
646 the N2 component. The TRF peak around this component was smaller and delayed for
647 speech in noise compared to speech in quiet. Delayed N2 response is associated with
648 attentive speech processing in challenging acoustic conditions (Balkenhol, Wallhäusser-
649 Franke, Rotter, & Servais, 2020; Billings, Tremblay, Stecker, & Tolin, 2009; Finke,
650 Büchner, Ruigendijk, Meyer, & Sandmann, 2016). Again, differences in this time range
651 (between 100 and 300 ms after stimulus onset) possibly reflect changes in the degree of
652 attention required to encode incoming stimuli effectively. Particularly, delayed TRF peak
653 response may reflect participants' effort in keeping track of meaningful information over
654 time in the degraded signal. Compensatory mechanisms may be involved in segregating
655 speech from noise. Previous evidence reported stronger envelope tracking of attended
656 speech with increased background noise in hearing-impaired and elderly individuals
657 compared to hearing younger adults (Brodbeck, Presacco, Anderson, & Simon, 2018;
658 Decruy, Vanthornhout, & Francart, 2020; Presacco, Simon, & Anderson, 2016). Both
659 internal (hearing loss) and external (background noise) factors can produce acoustic

660 distortion, which may result in increased listening effort (Van Engen & Peelle, 2014) and
661 enhanced envelope tracking.

662 There is a debate whether envelope tracking is enhanced (Ding et al., 2014; Ding & Simon,
663 2013; Fuglsang, Dau, & Hjortkjær, 2017; Presacco et al., 2016) or reduced (Desai et al.,
664 2021; Ding & Simon, 2013; Kurthen et al., 2021; Vanthornhout, Decruy, Wouters, Simon,
665 & Francart, 2018; L. Wang, Wu, & Chen, 2020) with decreasing SNR. Our behavioral
666 results showed that comprehension scores and intelligibility rates were directly
667 proportional to SNR levels. Our results on TRFs also add to the findings that envelope
668 tracking increases with noise and when listening becomes more challenging.

669

670 **Effects of high-level linguistic (Semantic) processing on neural tracking of speech** 671 **envelope**

672 Topographical distributions of the TRFs suggest the involvement of distinct neural
673 generators when semantic content is present or absent (Figure 3B). Moreover, the
674 temporal dynamics of TRFs for meaningful story was characterized by a more prominent
675 P1 peak and generally delayed P1-N1-P2-N2 complex, as compared to meaningless story
676 (Figure 3A).

677 At a relatively early processing stage (around 100 ms), we observed stronger neural
678 tracking of the speech envelope for meaningful story than for meaningless story over
679 fronto-central electrodes (Figure 3A, B). This finding could seem surprising since auditory
680 P1 is often associated with pre-attentive processes such as onset detection and sensory
681 gating (Huotilainen et al., 1998; Miller, Graham, & Schafer, 2021; Thoma et al., 2003;
682 Waldo et al., 1992). Predictive models of speech processing provide a plausible
683 explanation for this result. Semantic content generates expectations about upcoming
684 stimuli and limits the degree of uncertainty about what was heard (Poepffel, Idsardi, & van
685 Wassenhove, 2008), affecting early auditory encoding (Broderick, Anderson, & Lalor,
686 2019) and neural tracking of the speech envelope (Di Liberto et al., 2018; Kaufeld et al.,
687 2020). Meaningful information may provide regularities in meaningful story, making it more
688 predictable than meaningless story.

689 Moreover, it is possible that envelope tracking of meaningless story may not be affected
690 by the background noise as much as meaningful story due to the difference in the degree
691 of informational masking. It is possible that meaningless story could "pop-out" from the
692 background multi-talker babble noise due to lower informational masking compared to
693 meaningful story. Under the linguistic similarity hypothesis (Van Engen & Bradlow, 2007),
694 informational masking is more efficient when background babble noise has more linguistic
695 similarity with the target speech stream (e.g., same spoken language, known accent)
696 compared to a different or unknown language, accent and semantically anomalous

697 speech (Brouwer, Van Engen, Calandruccio, & Bradlow, 2012; Brungart, 2001;
698 Calandruccio, Van Engen, Dhar, & Bradlow, 2010; Cooke, Garcia Lecumberri, & Barker,
699 2008; Garcia Lecumberri & Cooke, 2006; Van Engen, 2010; Van Engen & Bradlow, 2007).
700 Therefore, it could have been easier for participants to segregate from the background
701 noise meaningless story than meaningful story.

702

703 **Two distributed networks are engaged in envelope tracking of continuous speech**

704 Source analysis of TRFs highlighted temporal and fronto-parietal regions traditionally
705 involved in speech and language comprehension (Hertrich, Dietrich, & Ackermann, 2020).
706 Key regions for low-level acoustic effect tested here involved the bilateral temporal cortex,
707 parts of the parietal, insular, and visual cortices bilaterally, and the left prefrontal cortex
708 (Figure 2D). Naturalistic speech stimuli are complex and resemble everyday listening
709 conditions, thus leading to extended activations and involvement of higher-order cortical
710 regions (Alexandrou, Saarinen, Kujala, & Salmelin, 2020; Hamilton & Huth, 2020). For
711 example, narrative speech involves widely distributed bilateral neural activity that tracks
712 hierarchically organized speech representations at multiple cortical sites and temporal
713 windows (de Heer, Huth, Griffiths, Gallant, & Theunissen, 2017; Di Liberto et al., 2015;
714 Huth, de Heer, Griffiths, Theunissen, & Gallant, 2016; Lerner, Honey, Silbert, & Hasson,
715 2011; Poeppel, 2003; Puschmann, Regev, Baillet, & Zatorre, 2021). Neuroimaging studies
716 reported distributed cortical activations beyond the auditory cortex (comprising higher-
717 order associative brain structures and attentional networks) during effortful listening (see
718 Alain, Du, Bernstein, Barten, & Banai, 2018 for a meta-analysis).

719 Higher-level linguistic processing was assessed by contrasting meaningful and
720 meaningless stories (Speech vs. Jabberwocky) and resulted in higher activation for
721 meaningful story, mainly involving the left auditory cortex, a large portion of bilateral fronto-
722 parietal network, and the left anterior temporal lobe later in time (Figure 3D). Overall
723 source modeling results of TRFs indicate that low-level acoustic effects mainly involved a
724 bilateral temporo-parietal network, while higher-level linguistic effects primarily involved a
725 left dominant fronto-temporal network. These results support the notion that successful
726 speech comprehension requires multiple extended networks beyond the temporal lobe to
727 process the acoustic signal at multiple and parallel hierarchical levels (Davis & Johnsrude,
728 2003, 2007; de Heer et al., 2017; Hickok & Poeppel, 2007; Peelle, 2012; Peelle,
729 Johnsrude, & Davis, 2010)

730

731 **Early visual cortex's entrainment to speech envelope in blindfolded individuals is** 732 **reduced by background noise**

733 We performed source analysis on the TRFs from preselected visual ROIs (V1 and V2) to
734 assess whether the visual cortex contributes to neural envelope tracking in blindfolded
735 individuals. While source estimates of EEG activity should be taken with caution, results
736 suggested early visual cortex's involvement in envelope tracking, especially for low-level
737 acoustic speech processing (Figure 3D, E).

738 A recent fMRI study showed that the visual cortex of blindfolded individuals displayed
739 some degree of synchrony to audio tracks from movies and narratives, suggesting that
740 auditory information can reach the visual cortices (Loiotile, Cusack, & Bedny, 2019).
741 Overall, numerous fMRI findings supported the notion that the visual cortex is functionally
742 engaged in processing non-visual stimuli in sighted individuals (Facchini & Aglioti, 2003;
743 Merabet et al., 2008; Poirier et al., 2006; Qin & Yu, 2013; Ricciardi et al., 2011; Sathian,
744 2005; Seydell-Greenwald et al., 2021; Vetter et al., 2014; Zangaladze, Epstein, Grafton,
745 & Sathian, 1999).

746 Interestingly, we observed a decrease in total signal magnitude for speech in noise
747 compared to speech in quiet. This difference emerged in particular for the right visual
748 cortex (although a trend also existed in the left hemisphere; Figure 2E). Hemispheric
749 asymmetry is not surprising, as previous evidence already showed the right hemisphere
750 dominance for several aspects of natural speech processing, especially for tracking of
751 slow temporal modulations within the delta-theta range (Alexandrou, Saarinen, Mäkelä,
752 Kujala, & Salmelin, 2017; Poeppel, 2003). More importantly, this finding aligns with the
753 evidence that the early visual cortex is sensitive to acoustic SNR effects (Bishop & Miller,
754 2009).

755 These results seem to suggest that the activity of the visual cortex could be modulated
756 during continuous speech tracking. However, its activity gets suppressed if the attentional
757 network becomes more engaged in tracking relevant auditory information in a challenging
758 listening environment. Human neuroimaging studies reported cross-modal deactivation of
759 the visual cortex by auditory stimuli during active listening or passive stimulation (with the
760 instructions to concentrate on the stimuli) and suggested that such suppression can be
761 top-down modulated by attention as task demands increase (e.g., Hairston et al., 2008;
762 Johnson & Zatorre, 2006; Laurienti et al., 2002). Several other studies found suppression
763 effects of sound on visual perception. Such cross-modal suppression has been suggested
764 to reduce the magnitude of the percept of a weaker or less relevant modality input
765 considered as a perceptual noise (Hidaka & Ide, 2015).

766 Overall, our results align with recent evidence reporting that the visual cortex can
767 contribute to auditory information processing in sighted individuals (Brang et al., 2022;
768 Martinelli et al., 2020; Seydell-Greenwald et al., 2021; Vetter et al., 2014). Here, we
769 observed that the visual cortex is more engaged in processing when speech signal is

770 intelligible and clear (i.e., presented in quiet). Differences in mapping speech envelope in
771 the visual cortex for low-level acoustic representations exist and might reflect cross-modal
772 visual cortex suppression. Such suppression could be top-down modulated and attributed
773 to auditory attention (Cate et al., 2009), which plays an essential role in segregating
774 relevant speech information in challenging listening conditions and when congruent visual
775 input is unavailable.

776 It could be argued that mental imagery mechanisms may drive the visual cortex's
777 response to speech. Previous studies observed an overlap in neural representations in
778 the occipital areas between perception and visual imagery, stemming from common top-
779 down influences (see Dijkstra, Bosch, & Gerven, 2019 for a review). However, V1 has
780 been shown to encode auditory information regardless of imageability (Martinelli et al.,
781 2020; Seydell-Greenwald et al., 2021; Vetter et al., 2020, 2014). Thus, the role of the early
782 visual cortex in auditory processing may not be merely ascribed to an imagery effect. If
783 that was the case, when contrasting *Speech-in-Noise* and *Jabberwocky-in-Noise*, we
784 could have observed higher visual cortex's responses in meaningful condition compared
785 to meaningless one, since only the former contained visually imaginable information.
786 However, no significant difference in the visual cortex's entrainment to speech envelope
787 was found between these conditions.

788

789

790 **Limitations and future research perspectives**

791 It is important to acknowledge the challenges of EEG-based source modeling, as spatial
792 resolution of EEG is generally known to be relatively poor, making it difficult to identify
793 exact brain sources that generate the neuronal activity measured on the scalp. EEG-
794 based source modeling majorly suffers from an ill-posed inverse problem and can also
795 result in misleading activity patterns due to, for instance, low SNR, unrealistic head
796 models, invalid constraints, and so on. More accurate EEG source localization requires
797 digitized electrode positions and individual anatomical scans of participants, which can
798 diminish source estimation uncertainty (Shirazi and Huang, 2019; Michel and Brunet,
799 2019; Zorzos et al., 2021) but were not available in our study. Therefore, EEG source
800 estimates should be interpreted with caution. However, it is worth noting that we used a
801 validated pipeline for source modeling estimation (Stropahl et al., 2018; Bottari et al.,
802 2020). Moreover, the same source modeling was performed across different conditions;
803 thus, similar errors should be attributed to activations for each condition. While the exact
804 location of the activity cannot be ensured with the present data, our results suggested that
805 the activity of posterior cortices was modulated only by low-level and not high-level speech
806 processing.

807 A further limitation pertains the input data we used for the encoding. We modeled neural
808 tracking of the speech signal based on a single feature: the speech envelope comprising
809 specific bandwidth frequencies (2-8 Hz). The envelope represents slow-variate temporal
810 modulations of the speech signal. It contains multiple acoustic and linguistic cues
811 important for continuous speech segmentation into smaller units, and therefore it has been
812 hypothesized to be crucial for speech comprehension (Luo & Poeppel, 2007; Shannon,
813 Zeng, Kamath, Wygonski, & Ekelid, 1995; Zoefel, 2018). However, it has also been argued
814 that focusing on the envelope alone might not get the complete picture of the neural
815 mechanism underlying speech comprehension (Obleser, Herrmann, & Henry, 2012).
816 Recent studies reported that the inclusion of multiple speech features, such as
817 spectrogram, phonemes, and phonetic features in the model sometimes result in a better
818 model performance represented by a more robust neural tracking response (e.g.,
819 Brodbeck, Hong, & Simon, 2018; Di Liberto et al., 2015, 2018; Lesenfants, Vanthornhout,
820 Verschueren, Decruy, & Francart, 2019). Future research may include multiple speech
821 features to build a multivariate model to assess neural speech tracking in the brain and
822 how the visual cortex maps speech information when visual input is absent.

823

824

825 **Conclusion**

826 Overall, our results indicate low-level acoustic and high-level linguistic processes affecting
827 envelope tracking of continuous speech. Envelope tracking may play a role in supporting
828 active listening in challenging conditions and is enhanced when SNR decreases, and
829 when segregation of target speech from the background noise becomes more difficult (i.e.,
830 due to linguistic similarity). Tracking speech signal embedded in noise requires spread
831 networks of activation, including linguistic and attentional regions beyond the auditory
832 cortex. In the absence of retinal input, the visual cortex might entrain to the speech
833 envelope, however, the functional role of such activity remains to be ascertained. The
834 magnitude of entrainment is degraded by concurrent noise, suggesting a suppressing
835 mechanism aimed at focusing resources within the auditory attention network in case of
836 challenging listening conditions. Conversely, no clear impact of semantic content was
837 found in the visual cortex, suggesting that the magnitude of such entrainment is generally
838 affected by low-level speech features.

839

840

841 **Acknowledgements**

842 The authors thank Chiara Battaglini who helped with stimuli recordings and preparation.
843 We also thank Chiara Maccioni for lending her voice for stimuli recordings. Davide Bottari
844 is funded by PRIN 2017 research grant. Prot. 20177894ZH.

845

846 **Author contributions**

847 Conceptualization, D.B., E.B., B.M., S.D.; Methodology, D.B., E.B., A.M.; Formal Analysis,
848 E.B., B.M., D.B.; Investigation, E.B., A.M.; Data Curation – E.B.; Writing – Original Draft,
849 E.B., M.B., A.F., D.B.; Writing – Review and Editing, E.B., D.B., M.B., A.F., B.M., S.D.,
850 E.R., A.M.; Visualization, E.B., D.B.; Resources & Funding, D.B.

851

852 **Declaration of interests**

853 The authors declare no competing interests.

854

855

856 **References**

- 857 Ahissar, E., Nagarajan, S., Ahissar, M., Protopapas, A., Mahncke, H., & Merzenich, M. M.
858 (2001). Speech comprehension is correlated with temporal response patterns recorded
859 from auditory cortex. *Proceedings of the National Academy of Sciences*, 98(23), 13367–
860 13372. <https://doi.org/10.1073/pnas.201400998>
- 861 Alain, C., Du, Y., Bernstein, L. J., Barten, T., & Banai, K. (2018). Listening under difficult
862 conditions: An activation likelihood estimation meta-analysis. *Human Brain Mapping*,
863 39(7), 2695–2709. <https://doi.org/10.1002/hbm.24031>
- 864 Alain, C., Quan, J., McDonald, K., & Van Roon, P. (2009). Noise-induced increase in
865 human auditory evoked neuromagnetic fields. *European Journal of Neuroscience*, 30(1),
866 132–142. <https://doi.org/10.1111/j.1460-9568.2009.06792.x>
- 867 Alexandrou, A. M., Saarinen, T., Kujala, J., & Salmelin, R. (2020). Cortical entrainment:
868 What we can learn from studying naturalistic speech perception. *Language, Cognition and*
869 *Neuroscience*, 35(6), 681–693. <https://doi.org/10.1080/23273798.2018.1518534>
- 870 Alexandrou, A. M., Saarinen, T., Mäkelä, S., Kujala, J., & Salmelin, R. (2017). The right
871 hemisphere is highlighted in connected natural speech production and perception.
872 *NeuroImage*, 152, 628–638. <https://doi.org/10.1016/j.neuroimage.2017.03.006>
- 873 Balkenhol, T., Wallhäusser-Franke, E., Rotter, N., & Servais, J. J. (2020). Changes in
874 Speech-Related Brain Activity During Adaptation to Electro-Acoustic Hearing. *Frontiers in*
875 *Neurology*, 11. <https://www.frontiersin.org/article/10.3389/fneur.2020.00161>
- 876 Benjamini, Y., & Hochberg, Y. (1995). Controlling the False Discovery Rate: A Practical
877 and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series*
878 *B (Methodological)*, 57(1), 289–300.
- 879 Billings, C. J., Tremblay, K. L., Stecker, G. C., & Tolin, W. M. (2009). Human evoked
880 cortical activity to signal-to-noise ratio and absolute signal level. *Hearing Research*,
881 254(1), 15–24. <https://doi.org/10.1016/j.heares.2009.04.002>

- 882 Bishop, C. W., & Miller, L. M. (2009). A Multisensory Cortical Network for Understanding
883 Speech in Noise. *Journal of Cognitive Neuroscience*, 21(9), 1790–1804.
884 <https://doi.org/10.1162/jocn.2009.21118>
- 885 Bottari, D., Bednaya, E., Dormal, G., Villwock, A., Dzhelyova, M., Grin, K., ... Röder, B.
886 (2020). EEG frequency-tagging demonstrates increased left hemispheric involvement and
887 crossmodal plasticity for face processing in congenitally deaf signers. *NeuroImage*, 223,
888 117315. doi: [10.1016/j.neuroimage.2020.117315](https://doi.org/10.1016/j.neuroimage.2020.117315)
- 889 Bourguignon, M., Baart, M., Kapnoula, E. C., & Molinaro, N. (2018). *Hearing through lip-*
890 *reading: The brain synthesizes features of absent speech.* bioRxiv.
891 <https://doi.org/10.1101/395483>
- 892 Bourguignon, M., Baart, M., Kapnoula, E. C., & Molinaro, N. (2020). Lip-Reading Enables
893 the Brain to Synthesize Auditory Features of Unknown Silent Speech. *Journal of*
894 *Neuroscience*, 40(5), 1053–1065. <https://doi.org/10.1523/JNEUROSCI.1101-19.2019>
- 895 Brang, D., Plass, J., Sherman, A., Stacey, W. C., Wasade, V. S., Grabowecky, M., ...
896 Suzuki, S. (2022). Visual cortex responds to sound onset and offset during passive
897 listening. *Journal of Neurophysiology*, 127(6), 1547–1563.
898 <https://doi.org/10.1152/jn.00164.2021>
- 899 Brodbeck, C., Hong, L. E., & Simon, J. Z. (2018). Rapid Transformation from Auditory to
900 Linguistic Representations of Continuous Speech. *Current Biology*, 28(24), 3976–3983.e5.
901 <https://doi.org/10.1016/j.cub.2018.10.042>
- 902 Brodbeck, C., Jiao, A., Hong, L. E., & Simon, J. Z. (2020). Neural speech restoration at
903 the cocktail party: Auditory cortex recovers masked speech of both attended and ignored
904 speakers. *PLOS Biology*, 18(10), e3000883. <https://doi.org/10.1371/journal.pbio.3000883>
- 905 Brodbeck, C., Presacco, A., Anderson, S., & Simon, J. Z. (2018). Over-Representation of
906 Speech in Older Adults Originates from Early Response in Higher Order Auditory Cortex.
907 *Acta Acustica United with Acustica*, 104(5), 774–777. <https://doi.org/10.3813/AAA.919221>
- 908 Broderick, M. P., Anderson, A. J., & Lalor, E. C. (2019). Semantic Context Enhances the
909 Early Auditory Encoding of Natural Speech. *Journal of Neuroscience*, 39(38), 7564–7575.
910 <https://doi.org/10.1523/JNEUROSCI.0584-19.2019>
- 911 Brouwer, S., Van Engen, K. J., Calandruccio, L., & Bradlow, A. R. (2012). Linguistic
912 contributions to speech-on-speech masking for native and non-native listeners: Language
913 familiarity and semantic content. *The Journal of the Acoustical Society of America*, 131(2),
914 1449–1464. <https://doi.org/10.1121/1.3675943>
- 915 Brungart, D. S. (2001). Informational and energetic masking effects in the perception of
916 two simultaneous talkers. *The Journal of the Acoustical Society of America*, 109(3), 1101–
917 1109. <https://doi.org/10.1121/1.1345696>
- 918 Calandruccio, L., Van Engen, K., Dhar, S., & Bradlow, A. R. (2010). The Effectiveness of
919 Clear Speech as a Masker. *Journal of Speech, Language, and Hearing Research : JSLHR*,
920 53(6), 1458–1471. [https://doi.org/10.1044/1092-4388\(2010\)09-0210](https://doi.org/10.1044/1092-4388(2010)09-0210)
- 921 Cate, A. D., Herron, T. J., Yund, E. W., Stecker, G. C., Rinne, T., Kang, X., Petkov, C. I.,
922 Disbrow, E. A., & Woods, D. L. (2009). Auditory Attention Activates Peripheral Visual
923 Cortex. *PLOS ONE*, 4(2), e4645. <https://doi.org/10.1371/journal.pone.0004645>

- 924 Ceponienė, R., Alku, P., Westerfield, M., Torki, M., & Townsend, J. (2005). ERPs
925 differentiate syllable and nonphonetic sound processing in children and adults.
926 *Psychophysiology*, 42(4), 391–406. <https://doi.org/10.1111/j.1469-8986.2005.00305.x>
- 927 Cherry, E. C. (1953). Some Experiments on the Recognition of Speech, with One and with
928 Two Ears. *The Journal of the Acoustical Society of America*, 25(5), 975–979.
929 <https://doi.org/10.1121/1.1907229>
- 930 Combrisson, E., & Jerbi, K. (2015). Exceeding chance level by chance: The caveat of
931 theoretical chance levels in brain signal classification and statistical assessment of
932 decoding accuracy. *Journal of Neuroscience Methods*, 250, 126–136.
933 <https://doi.org/10.1016/j.jneumeth.2015.01.010>
- 934 Cooke, M., Garcia Lecumberri, M. L., & Barker, J. (2008). The foreign language cocktail
935 party problem: Energetic and informational masking effects in non-native speech
936 perception. *The Journal of the Acoustical Society of America*, 123(1), 414–427.
937 <https://doi.org/10.1121/1.2804952>
- 938 Crosse, M. J., Di Liberto, G. M., Bednar, A., & Lalor, E. C. (2016). The Multivariate
939 Temporal Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural
940 Signals to Continuous Stimuli. *Frontiers in Human Neuroscience*, 10.
941 <https://doi.org/10.3389/fnhum.2016.00604>
- 942 Crosse, M. J., Zuk, N. J., Di Liberto, G. M., Nidiffer, A. R., Molholm, S., & Lalor, E. C.
943 (2021). Linear Modeling of Neurophysiological Responses to Speech and Other
944 Continuous Stimuli: Methodological Considerations for Applied Research. *Frontiers in*
945 *Neuroscience*, 15, 1350. <https://doi.org/10.3389/fnins.2021.705621>
- 946 Dale, A. M., Liu, A. K., Fischl, B. R., Buckner, R. L., Belliveau, J. W., Lewine, J. D., &
947 Halgren, E. (2000). Dynamic Statistical Parametric Mapping: Combining fMRI and MEG
948 for High-Resolution Imaging of Cortical Activity. *Neuron*, 26(1), 55–67.
949 [https://doi.org/10.1016/S0896-6273\(00\)81138-1](https://doi.org/10.1016/S0896-6273(00)81138-1)
- 950 Davis, M. H., & Johnsrude, I. S. (2003). Hierarchical Processing in Spoken Language
951 Comprehension. *Journal of Neuroscience*, 23(8), 3423–3431.
952 <https://doi.org/10.1523/JNEUROSCI.23-08-03423.2003>
- 953 Davis, M. H., & Johnsrude, I. S. (2007). Hearing speech sounds: Top-down influences on
954 the interface between audition and speech perception. *Hearing Research*, 229(1–2), 132–
955 147. <https://doi.org/10.1016/j.heares.2007.01.014>
- 956 de Heer, W. A., Huth, A. G., Griffiths, T. L., Gallant, J. L., & Theunissen, F. E. (2017). The
957 Hierarchical Cortical Organization of Human Speech Processing. *The Journal of*
958 *Neuroscience*, 37(27), 6539–6557. <https://doi.org/10.1523/JNEUROSCI.3267-16.2017>
- 959 Decruy, L., Vanthornhout, J., & Francart, T. (2019). Evidence for enhanced neural tracking
960 of the speech envelope underlying age-related speech-in-noise difficulties. *Journal of*
961 *Neurophysiology*, 122(2), 601–615. <https://doi.org/10.1152/jn.00687.2018>
- 962 Decruy, L., Vanthornhout, J., & Francart, T. (2020). Hearing impairment is associated with
963 enhanced neural tracking of the speech envelope. *Hearing Research*, 393, 107961.
964 <https://doi.org/10.1016/j.heares.2020.107961>
- 965 Delorme, A., & Makeig, S. (2004). EEGLAB: An open source toolbox for analysis of single-
966 trial EEG dynamics including independent component analysis. *Journal of Neuroscience*
967 *Methods*, 134(1), 9–21. <https://doi.org/10.1016/j.jneumeth.2003.10.009>

- 968 Desai, M., Holder, J., Villarreal, C., Clark, N., Hoang, B., & Hamilton, L. S. (2021).
969 Generalizable EEG Encoding Models with Naturalistic Audiovisual Stimuli. *Journal of*
970 *Neuroscience*, 41(43), 8946–8962. <https://doi.org/10.1523/JNEUROSCI.2891-20.2021>
- 971 Destrieux, C., Fischl, B., Dale, A., & Halgren, E. (2010). Automatic parcellation of human
972 cortical gyri and sulci using standard anatomical nomenclature. *NeuroImage*, 53(1), 1–15.
973 <https://doi.org/10.1016/j.neuroimage.2010.06.010>
- 974 Di Liberto, G. M., O'Sullivan, J. A., & Lalor, E. C. (2015). Low-Frequency Cortical
975 Entrainment to Speech Reflects Phoneme-Level Processing. *Current Biology*, 25(19),
976 2457–2465. <https://doi.org/10.1016/j.cub.2015.08.030>
- 977 Di Liberto, G. M., Peter, V., Kalashnikova, M., Goswami, U., Burnham, D., & Lalor, E. C.
978 (2018). Atypical cortical entrainment to speech in the right hemisphere underpins
979 phonemic deficits in dyslexia. *NeuroImage*, 175, 70–79.
980 <https://doi.org/10.1016/j.neuroimage.2018.03.072>
- 981 Dijkstra, N., Bosch, S. E., & Gerven, M. A. J. van. (2019). Shared Neural Mechanisms of
982 Visual Perception and Imagery. *Trends in Cognitive Sciences*, 23(5), 423–434.
983 <https://doi.org/10.1016/j.tics.2019.02.004>
- 984 Ding, N., Chatterjee, M., & Simon, J. Z. (2014). Robust cortical entrainment to the speech
985 envelope relies on the spectro-temporal fine structure. *NeuroImage*, 88, 41–46.
986 <https://doi.org/10.1016/j.neuroimage.2013.10.054>
- 987 Ding, N., Melloni, L., Zhang, H., Tian, X., & Poeppel, D. (2016). Cortical tracking of
988 hierarchical linguistic structures in connected speech. *Nature Neuroscience*, 19(1), 158–
989 164. <https://doi.org/10.1038/nn.4186>
- 990 Ding, N., & Simon, J. Z. (2012). Emergence of neural encoding of auditory objects while
991 listening to competing speakers. *Proceedings of the National Academy of Sciences*,
992 109(29), 11854–11859. <https://doi.org/10.1073/pnas.1205381109>
- 993 Ding, N., & Simon, J. Z. (2013). Adaptive Temporal Encoding Leads to a Background-
994 Insensitive Cortical Representation of Speech. *Journal of Neuroscience*, 33(13), 5728–
995 5735. <https://doi.org/10.1523/JNEUROSCI.5297-12.2013>
- 996 Ding, N., & Simon, J. Z. (2014). Cortical entrainment to continuous speech: Functional
997 roles and interpretations. *Frontiers in Human Neuroscience*, 8.
998 <https://doi.org/10.3389/fnhum.2014.00311>
- 999 Drennan, D. P., & Lalor, E. C. (2019). Cortical Tracking of Complex Sound Envelopes:
1000 Modeling the Changes in Response with Intensity. *ENeuro*, 6(3).
1001 <https://doi.org/10.1523/ENEURO.0082-19.2019>
- 1002 Facchini, S., & Aglioti, S. M. (2003). Short term light deprivation increases tactile spatial
1003 acuity in humans. *Neurology*, 60(12), 1998–1999.
1004 <https://doi.org/10.1212/01.wnl.0000068026.15208.d0>
- 1005 Fiedler, L., Wöstmann, M., Herbst, S. K., & Obleser, J. (2019). Late cortical tracking of
1006 ignored speech facilitates neural selectivity in acoustically challenging conditions.
1007 *NeuroImage*, 186, 33–42. <https://doi.org/10.1016/j.neuroimage.2018.10.057>
- 1008 Finke, M., Büchner, A., Ruigendijk, E., Meyer, M., & Sandmann, P. (2016). On the
1009 relationship between auditory cognition and speech intelligibility in cochlear implant users:
1010 An ERP study. *Neuropsychologia*, 87, 169–181.
1011 <https://doi.org/10.1016/j.neuropsychologia.2016.05.019>

- 1012 Fischl, B. (2012). FreeSurfer. *NeuroImage*, 62(2), 774–781.
1013 <https://doi.org/10.1016/j.neuroimage.2012.01.021>
- 1014 Fuglsang, S. A., Dau, T., & Hjortkjær, J. (2017). Noise-robust cortical tracking of attended
1015 speech in real-world acoustic scenes. *NeuroImage*, 156, 435–444.
1016 <https://doi.org/10.1016/j.neuroimage.2017.04.026>
- 1017 Garcia Lecumberri, M. L., & Cooke, M. (2006). Effect of masker type on native and non-
1018 native consonant perception in noise. *The Journal of the Acoustical Society of America*,
1019 119(4), 2445–2454. <https://doi.org/10.1121/1.2180210>
- 1020 Gross, J., Hoogenboom, N., Thut, G., Schyns, P., Panzeri, S., Belin, P., & Garrod, S.
1021 (2013). Speech Rhythms and Multiplexed Oscillatory Sensory Coding in the Human Brain.
1022 *PLOS Biology*, 11(12), e1001752. <https://doi.org/10.1371/journal.pbio.1001752>
- 1023 Gustafson, S. J., Billings, C. J., Hornsby, B. W. Y., & Key, A. P. (2019). Effect of competing
1024 noise on cortical auditory evoked potentials elicited by speech sounds in 7- to 25-year-old
1025 listeners. *Hearing Research*, 373, 103–112. <https://doi.org/10.1016/j.heares.2019.01.004>
- 1026 Hairston, W. D., Hodges, D. A., Casanova, R., Hayasaka, S., Kraft, R., Maldjian, J. A., &
1027 Burdette, J. H. (2008). Closing the mind's eye: Deactivation of visual cortex related to
1028 auditory task difficulty. *NeuroReport*, 19(2), 151–154.
1029 <https://doi.org/10.1097/WNR.0b013e3282f42509>
- 1030 Hamilton, L. S., & Huth, A. G. (2020). The revolution will not be controlled: Natural stimuli
1031 in speech neuroscience. *Language, Cognition and Neuroscience*, 35(5), 573–582.
1032 <https://doi.org/10.1080/23273798.2018.1499946>
- 1033 Hansen, J. C., & Hillyard, S. A. (1980). Endogenous brain potentials associated with
1034 selective auditory attention. *Electroencephalography and Clinical Neurophysiology*, 49(3–
1035 4), 277–290. [https://doi.org/10.1016/0013-4694\(80\)90222-9](https://doi.org/10.1016/0013-4694(80)90222-9)
- 1036 Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J.-D., Blankertz, B., & Bießmann,
1037 F. (2014). On the interpretation of weight vectors of linear models in multivariate
1038 neuroimaging. *NeuroImage*, 87, 96–110.
1039 <https://doi.org/10.1016/j.neuroimage.2013.10.067>
- 1040 Hausfeld, L., Riecke, L., Valente, G., & Formisano, E. (2018). Cortical tracking of multiple
1041 streams outside the focus of attention in naturalistic auditory scenes. *NeuroImage*, 181,
1042 617–626. <https://doi.org/10.1016/j.neuroimage.2018.07.052>
- 1043 Hauswald, A., Lithari, C., Collignon, O., Leonardelli, E., & Weisz, N. (2018). A Visual
1044 Cortical Network for Deriving Phonological Information from Intelligible Lip Movements.
1045 *Current Biology*, 28(9), 1453-1459.e3. doi: [10.1016/j.cub.2018.03.044](https://doi.org/10.1016/j.cub.2018.03.044)
- 1046 Hertrich, I., Dietrich, S., & Ackermann, H. (2020). The Margins of the Language Network
1047 in the Brain. *Frontiers in Communication*, 5.
1048 <https://www.frontiersin.org/article/10.3389/fcomm.2020.519955>
- 1049 Hickok, G., & Poeppel, D. (2007). The cortical organization of speech processing. *Nature*
1050 *Reviews Neuroscience*, 8(5), 393–402. doi: [10.1038/nrn2113](https://doi.org/10.1038/nrn2113)
- 1051 Hidaka, S., & Ide, M. (2015). Sound can suppress visual perception. *Scientific Reports*,
1052 5(1), 10483. <https://doi.org/10.1038/srep10483>
- 1053 Howard, M. F., & Poeppel, D. (2010). Discrimination of Speech Stimuli Based on Neuronal
1054 Response Phase Patterns Depends on Acoustics But Not Comprehension. *Journal of*
1055 *Neurophysiology*, 104(5), 2500–2511. <https://doi.org/10.1152/jn.00251.2010>

- 1056 Huotilainen, M., Winkler, I., Alho, K., Escera, C., Virtanen, J., Ilmoniemi, R. J.,
1057 Jääskeläinen, I. P., Pekkonen, E., & Näätänen, R. (1998). Combined mapping of human
1058 auditory EEG and MEG responses. *Electroencephalography and Clinical*
1059 *Neurophysiology*, 108(4), 370–379. [https://doi.org/10.1016/s0168-5597\(98\)00017-3](https://doi.org/10.1016/s0168-5597(98)00017-3)
- 1060 Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E., & Gallant, J. L. (2016).
1061 Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*,
1062 532(7600), 453–458. <https://doi.org/10.1038/nature17637>
- 1063 Johnson, J. A., & Zatorre, R. J. (2006). Neural substrates for dividing and focusing
1064 attention between simultaneous auditory and visual events. *NeuroImage*, 31(4), 1673–
1065 1681. <https://doi.org/10.1016/j.neuroimage.2006.02.026>
- 1066 Kaufeld, G., Bosker, H. R., Oever, S. ten, Alday, P. M., Meyer, A. S., & Martin, A. E.
1067 (2020). Linguistic Structure and Meaning Organize Neural Oscillations into a Content-
1068 Specific Hierarchy. *Journal of Neuroscience*, 40(49), 9467–9475.
1069 <https://doi.org/10.1523/JNEUROSCI.0302-20.2020>
- 1070 Kerlin, J. R., Shahin, A. J., & Miller, L. M. (2010). Attentional Gain Control of Ongoing
1071 Cortical Speech Representations in a "Cocktail Party." *The Journal of Neuroscience*,
1072 30(2), 620–628. <https://doi.org/10.1523/JNEUROSCI.3631-09.2010>
- 1073 Kurthen, I., Galbier, J., Jagoda, L., Neuschwander, P., Giroud, N., & Meyer, M. (2021).
1074 Selective attention modulates neural envelope tracking of informationally masked speech
1075 in healthy older adults. *Human Brain Mapping*, 42(10), 3042–3057.
1076 <https://doi.org/10.1002/hbm.25415>
- 1077 Lakatos, P., Gross, J., & Thut, G. (2019). A New Unifying Account of the Roles of Neuronal
1078 Entrainment. *Current Biology*, 29(18), R890–R905.
1079 <https://doi.org/10.1016/j.cub.2019.07.075>
- 1080 Lator, E. C., & Foxe, J. J. (2010). Neural responses to uninterrupted natural speech can
1081 be extracted with precise temporal resolution. *European Journal of Neuroscience*, 31(1),
1082 189–193. <https://doi.org/10.1111/j.1460-9568.2009.07055.x>
- 1083 Lator, E. C., Pearlmuter, B. A., Reilly, R. B., McDarby, G., & Foxe, J. J. (2006). The
1084 VESPA: A method for the rapid estimation of a visual evoked potential. *NeuroImage*,
1085 32(4), 1549–1561. <https://doi.org/10.1016/j.neuroimage.2006.05.054>
- 1086 Lator, E. C., Power, A. J., Reilly, R. B., & Foxe, J. J. (2009). Resolving Precise Temporal
1087 Processing Properties of the Auditory System Using Continuous Stimuli. *Journal of*
1088 *Neurophysiology*, 102(1), 349–359. <https://doi.org/10.1152/jn.90896.2008>
- 1089 Landry, S. P., Shiller, D. M., & Champoux, F. (2013). Short-term visual deprivation
1090 improves the perception of harmonicity. *Journal of Experimental Psychology: Human*
1091 *Perception and Performance*, 39(6), 1503–1507. <https://doi.org/10.1037/a0034015>
- 1092 Laurienti, P. J., Burdette, J. H., Wallace, M. T., Yen, Y.-F., Field, A. S., & Stein, B. E.
1093 (2002). Deactivation of Sensory-Specific Cortex by Cross-Modal Stimuli. *Journal of*
1094 *Cognitive Neuroscience*, 14(3), 420–429. <https://doi.org/10.1162/089892902317361930>
- 1095 Lazzouni, L., Voss, P., & Lepore, F. (2012). Short-term crossmodal plasticity of the
1096 auditory steady-state response in blindfolded sighted individuals. *European Journal of*
1097 *Neuroscience*, 35(10), 1630–1636. <https://doi.org/10.1111/j.1460-9568.2012.08088.x>

- 1098 Legendre, G., Andrillon, T., Koroma, M., & Kouider, S. (2019). Sleepers track informative
1099 speech in a multitalker environment. *Nature Human Behaviour*, 3(3), 274.
1100 <https://doi.org/10.1038/s41562-018-0502-5>
- 1101 Lerner, Y., Honey, C. J., Silbert, L. J., & Hasson, U. (2011). Topographic Mapping of a
1102 Hierarchy of Temporal Receptive Windows Using a Narrated Story. *Journal of*
1103 *Neuroscience*, 31(8), 2906–2915. <https://doi.org/10.1523/JNEUROSCI.3684-10.2011>
- 1104 Lesenfants, D., Vanthornhout, J., Verschueren, E., Decruy, L., & Francart, T. (2019).
1105 Predicting individual speech intelligibility from the cortical tracking of acoustic- and
1106 phonetic-level speech representations. *Hearing Research*, 380, 1–9.
1107 <https://doi.org/10.1016/j.heares.2019.05.006>
- 1108 Liu, Q., & Wang, L. (2021). T-Test and ANOVA for data with ceiling and/or floor effects.
1109 *Behavior Research Methods*, 53(1), 264–277. [https://doi.org/10.3758/s13428-020-01407-](https://doi.org/10.3758/s13428-020-01407-2)
1110 [2](https://doi.org/10.3758/s13428-020-01407-2)
- 1111 Loiotile, R. E., Cusack, R., & Bedny, M. (2019). Naturalistic Audio-Movies and Narrative
1112 Synchronize "Visual" Cortices across Congenitally Blind But Not Sighted Individuals.
1113 *Journal of Neuroscience*, 39(45), 8940–8948. [https://doi.org/10.1523/JNEUROSCI.0298-](https://doi.org/10.1523/JNEUROSCI.0298-19.2019)
1114 [19.2019](https://doi.org/10.1523/JNEUROSCI.0298-19.2019)
- 1115 Luo, H., & Poeppel, D. (2007). Phase Patterns of Neuronal Responses Reliably
1116 Discriminate Speech in Human Auditory Cortex. *Neuron*, 54(6), 1001–1010.
1117 <https://doi.org/10.1016/j.neuron.2007.06.004>
- 1118 Maraini, F. (2019). Gnòsi delle fànfole. La nave di Teseo.
- 1119 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-
1120 data. *Journal of Neuroscience Methods*, 164(1), 177–190.
1121 <https://doi.org/10.1016/j.jneumeth.2007.03.024>
- 1122 Martinelli, A., Handjaras, G., Betta, M., Leo, A., Cecchetti, L., Pietrini, P., Ricciardi, E., &
1123 Bottari, D. (2020). *Auditory features modelling demonstrates sound envelope*
1124 *representation in striate cortex* bioRxiv. <https://doi.org/10.1101/2020.04.15.043174>
- 1125 McCarthy, C. (2014). *The Road* (M. Testa, Trans.). Einaudi. (Original work published
1126 2006).
- 1127 Merabet, L. B., Hamilton, R., Schlaug, G., Swisher, J. D., Kiriakopoulos, E. T., Pitskel, N.
1128 B., Kauffman, T., & Pascual-Leone, A. (2008). Rapid and Reversible Recruitment of Early
1129 Visual Cortex for Touch. *PLOS ONE*, 3(8), e3046.
1130 <https://doi.org/10.1371/journal.pone.0003046>
- 1131 Michel, C. M., & Brunet, D. (2019). EEG Source Imaging: A Practical Review of the
1132 Analysis Steps. *Frontiers in Neurology*, 10. <https://doi.org/10.3389/fneur.2019.00325>
- 1133 Miller, S. E., Graham, J., & Schafer, E. (2021). Auditory Sensory Gating of Speech and
1134 Nonspeech Stimuli. *Journal of Speech, Language, and Hearing Research*, 64(4), 1404–
1135 1412. https://doi.org/10.1044/2020_JSLHR-20-00535
- 1136 Mirkovic, B., Debener, S., Jaeger, M., & De Vos, M. (2015). Decoding the attended speech
1137 stream with multi-channel EEG: Implications for online, daily-life applications. *Journal of*
1138 *Neural Engineering*, 12(4), 046007. <https://doi.org/10.1088/1741-2560/12/4/046007>
- 1139 Müller, J. A., Wendt, D., Kollmeier, B., Debener, S., & Brand, T. (2019). Effect of Speech
1140 Rate on Neural Tracking of Speech. *Frontiers in Psychology*, 10.
1141 <https://doi.org/10.3389/fpsyg.2019.00449>

- 1142 Näätänen, R. (1982). Processing negativity: An evoked-potential reflection of selective
1143 attention. *Psychological Bulletin*, 92(3), 605–640. [https://doi.org/10.1037/0033-](https://doi.org/10.1037/0033-2909.92.3.605)
1144 [2909.92.3.605](https://doi.org/10.1037/0033-2909.92.3.605)
- 1145 Näätänen, R., & Picton, T. (1987). The N1 wave of the human electric and magnetic
1146 response to sound: A review and an analysis of the component structure.
1147 *Psychophysiology*, 24(4), 375–425. <https://doi.org/10.1111/j.1469-8986.1987.tb00311.x>
- 1148 Obleser, J., Herrmann, B., & Henry, M. (2012). Neural Oscillations in Speech: Don't be
1149 Enslaved by the Envelope. *Frontiers in Human Neuroscience*, 6.
1150 <https://doi.org/10.3389/fnhum.2012.00250>
- 1151 Obleser, J., & Kayser, C. (2019). Neural Entrainment and Attentional Selection in the
1152 Listening Brain. *Trends in Cognitive Sciences*. <https://doi.org/10.1016/j.tics.2019.08.004>
- 1153 Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J.-M. (2011). FieldTrip: Open source
1154 software for advanced analysis of MEG, EEG, and invasive electrophysiological data.
1155 *Computational Intelligence and Neuroscience*, 2011, 156869.
1156 <https://doi.org/10.1155/2011/156869>
- 1157 O'Sullivan, A. E., Crosse, M. J., Liberto, G. M. D., Cheveigné, A. de, & Lalor, E. C. (2021).
1158 Neurophysiological Indices of Audiovisual Speech Processing Reveal a Hierarchy of
1159 Multisensory Integration Effects. *Journal of Neuroscience*, 41(23), 4991–5003.
1160 <https://doi.org/10.1523/JNEUROSCI.0906-20.2021>
- 1161 O'Sullivan, J. A., Power, A. J., Mesgarani, N., Rajaram, S., Foxe, J. J., Shinn-
1162 Cunningham, B. G., Slaney, M., Shamma, S. A., & Lalor, E. C. (2015). Attentional
1163 Selection in a Cocktail Party Environment Can Be Decoded from Single-Trial EEG.
1164 *Cerebral Cortex*, 25(7), 1697–1706. <https://doi.org/10.1093/cercor/bht355>
- 1165 Papesh, M. A., Billings, C. J., & Baltzell, L. S. (2015). Background noise can enhance
1166 cortical auditory evoked potentials under certain conditions. *Clinical Neurophysiology: Official Journal of the International Federation of Clinical Neurophysiology*, 126(7), 1319–
1167 1330. <https://doi.org/10.1016/j.clinph.2014.10.017>
- 1169 Parbery-Clark, A., Marmel, F., Bair, J., & Kraus, N. (2011). What subcortical–cortical
1170 relationships tell us about processing speech in noise. *European Journal of Neuroscience*,
1171 33(3), 549–557. <https://doi.org/10.1111/j.1460-9568.2010.07546.x>
- 1172 Park, H., Kayser, C., Thut, G., & Gross, J. (2016). Lip movements entrain the observers'
1173 low-frequency brain oscillations to facilitate speech intelligibility. *eLife*, 5.
1174 <https://doi.org/10.7554/eLife.14521>
- 1175 Peelle, J. (2012). The hemispheric lateralization of speech processing depends on what
1176 "speech" is: A hierarchical perspective. *Frontiers in Human Neuroscience*, 6, 309.
1177 <https://doi.org/10.3389/fnhum.2012.00309>
- 1178 Peelle, J. E., Gross, J., & Davis, M. H. (2013). Phase-Locked Responses to Speech in
1179 Human Auditory Cortex are Enhanced During Comprehension. *Cerebral Cortex*, 23(6),
1180 1378–1387. <https://doi.org/10.1093/cercor/bhs118>
- 1181 Peelle, J. E., Johnsrude, I., & Davis, M. H. (2010). Hierarchical processing for speech in
1182 human auditory cortex and beyond. *Frontiers in Human Neuroscience*, 4.
1183 <https://doi.org/10.3389/fnhum.2010.00051>

- 1184 Petro, L. S., Paton, A. T., & Muckli, L. (2017). Contextual modulation of primary visual
1185 cortex by auditory signals. *Philosophical Transactions of the Royal Society B: Biological*
1186 *Sciences*, 372(1714), 20160104. <https://doi.org/10.1098/rstb.2016.0104>
- 1187 Pitzorno, B. (1993). Polissena del Porcello [Polissena and her Pig]. Mondadori.
- 1188 Plass, J., Brang, D., Suzuki, S., & Grabowecky, M. (2020). Vision perceptually restores
1189 auditory spectral dynamics in speech. *Proceedings of the National Academy of Sciences*,
1190 117(29), 16920–16927. <https://doi.org/10.1073/pnas.2002887117>
- 1191 Poeppel, D. (2003). The analysis of speech in different temporal integration windows:
1192 Cerebral lateralization as 'asymmetric sampling in time.' *Speech Communication*, 41(1),
1193 245–255. [https://doi.org/10.1016/S0167-6393\(02\)00107-3](https://doi.org/10.1016/S0167-6393(02)00107-3)
- 1194 Poeppel, D., & Assaneo, M. F. (2020). Speech rhythms and their neural foundations.
1195 *Nature Reviews Neuroscience*, 21(6), 322–334. [https://doi.org/10.1038/s41583-020-](https://doi.org/10.1038/s41583-020-0304-4)
1196 [0304-4](https://doi.org/10.1038/s41583-020-0304-4)
- 1197 Poeppel, D., Idsardi, W. J., & van Wassenhove, V. (2008). Speech perception at the
1198 interface of neurobiology and linguistics. *Philosophical Transactions of the Royal Society*
1199 *of London. Series B, Biological Sciences*, 363(1493), 1071–1086.
1200 <https://doi.org/10.1098/rstb.2007.2160>
- 1201 Poirier, C., Collignon, O., Scheiber, C., Renier, L., Vanlierde, A., Tranduy, D., Veraart, C.,
1202 & De Volder, A. G. (2006). Auditory motion perception activates visual motion areas in
1203 early blind subjects. *NeuroImage*, 31(1), 279–285.
1204 <https://doi.org/10.1016/j.neuroimage.2005.11.036>
- 1205 Presacco, A., Simon, J. Z., & Anderson, S. (2016). Evidence of degraded representation
1206 of speech in noise, in the aging midbrain and cortex. *Journal of Neurophysiology*, 116(5),
1207 2346–2355. <https://doi.org/10.1152/jn.00372.2016>
- 1208 Puschmann, S., Regev, M., Baillet, S., & Zatorre, R. J. (2021). MEG Intersubject Phase
1209 Locking of Stimulus-Driven Activity during Naturalistic Speech Listening Correlates with
1210 Musical Training. *Journal of Neuroscience*, 41(12), 2713–2722.
1211 <https://doi.org/10.1523/JNEUROSCI.0932-20.2020>
- 1212 Qin, W., & Yu, C. (2013). Neural Pathways Conveying Novisual Information to the Visual
1213 Cortex. *Neural Plasticity*, 2013, 864920. <https://doi.org/10.1155/2013/864920>
- 1214 Queneau, R. (1983). Exercices de style (U. Eco, Trans.). Einaudi. (Original work published
1215 1947).
- 1216 Ricciardi, E., Basso, D., Sani, L., Bonino, D., Vecchi, T., Pietrini, P., & Miniussi, C. (2011).
1217 Functional inhibition of the human middle temporal cortex affects non-visual motion
1218 perception: A repetitive transcranial magnetic stimulation study during tactile speed
1219 discrimination. *Experimental Biology and Medicine*, 236(2), 138–144.
1220 <https://doi.org/10.1258/ebm.2010.010230>
- 1221 Riecke, L., Formisano, E., Sorger, B., Baskent, D., & Gaudrain, E. (2018). Neural
1222 Entrainment to Speech Modulates Speech Intelligibility. *Current Biology*, 28(2), 161-
1223 169.e5. <https://doi.org/10.1016/j.cub.2017.11.033>
- 1224 Sathian, K. (2005). Visual cortical activity during tactile perception in the sighted and the
1225 visually deprived. *Developmental Psychobiology*, 46(3), 279–286.
1226 <https://doi.org/10.1002/dev.20056>

- 1227 Schneider, W., Eschman, A., & Zuccolotto, A. (2002). E-Prime Reference Guide.
1228 Pittsburg, PA: Psychology Software Tools.
- 1229 Seydell-Greenwald, A., Wang, X., Newport, E., Bi, Y., & Striem-Amit, E. (2021). Primary
1230 visual cortex is activated by spoken language comprehension. *Journal of Vision*, 21(9),
1231 2256. <https://doi.org/10.1167/jov.21.9.2256>
- 1232 Shannon, R. V., Zeng, F.-G., Kamath, V., Wygonski, J., & Ekelid, M. (1995). Speech
1233 Recognition with Primarily Temporal Cues. *Science*.
1234 <https://doi.org/10.1126/science.270.5234.303>
- 1235 Shirazi, S. Y., & Huang, H. J. (2019). More Reliable EEG Electrode Digitizing Methods
1236 Can Reduce Source Estimation Uncertainty, but Current Methods Already Accurately
1237 Identify Brodmann Areas. *Frontiers in Neuroscience*, 13, 1159.
1238 <https://doi.org/10.3389/fnins.2019.01159>
- 1239 Šimkovic, M., & Träuble, B. (2019). Robustness of statistical methods when measure is
1240 affected by ceiling and/or floor effect. *PLOS ONE*, 14(8), e0220889.
1241 <https://doi.org/10.1371/journal.pone.0220889>
- 1242 Stropahl, M., Bauer, A.-K. R., Debener, S., & Bleichner, M. G. (2018). Source-Modeling
1243 Auditory Processes of EEG Data Using EEGLAB and Brainstorm. *Frontiers in*
1244 *Neuroscience*, 12, 309. <https://doi.org/10.3389/fnins.2018.00309>
- 1245 Tadel, F., Baillet, S., Mosher, J. C., Pantazis, D., & Leahy, R. M. (2011). Brainstorm: A
1246 User-Friendly Application for MEG/EEG Analysis. *Computational Intelligence and*
1247 *Neuroscience*, 2011, e879716. <https://doi.org/10.1155/2011/879716>
- 1248 Thoma, R. J., Hanlon, F. M., Moses, S. N., Edgar, J. C., Huang, M., Weisend, M. P., Irwin,
1249 J., Sherwood, A., Paulson, K., Bustillo, J., Adler, L. E., Miller, G. A., & Cañive, J. M. (2003).
1250 Lateralization of Auditory Sensory Gating and Neuropsychological Dysfunction in
1251 Schizophrenia. *American Journal of Psychiatry*, 160(9), 1595–1605.
1252 <https://doi.org/10.1176/appi.ajp.160.9.1595>
- 1253 Van Ackeren, M. J., Barbero, F. M., Mattioni, S., Bottini, R., & Collignon, O. (2018).
1254 Neuronal populations in the occipital cortex of the blind synchronize to the temporal
1255 dynamics of speech. *ELife*, 7, e31640. <https://doi.org/10.7554/eLife.31640>
- 1256 Van Engen, K. J. (2010). Similarity and familiarity: Second language sentence recognition
1257 in first- and second-language multi-talker babble. *Speech Communication*, 52(11–12),
1258 943–953. <https://doi.org/10.1016/j.specom.2010.05.002>
- 1259 Van Engen, K. J., & Bradlow, A. R. (2007). Sentence recognition in native- and foreign-
1260 language multi-talker background noise. *The Journal of the Acoustical Society of America*,
1261 121(1), 519–526. <https://doi.org/10.1121/1.2400666>
- 1262 Van Engen, K. J., & Peelle, J. E. (2014). Listening effort and accented speech. *Frontiers*
1263 *in Human Neuroscience*, 8. <https://www.frontiersin.org/article/10.3389/fnhum.2014.00577>
- 1264 Vanthornhout, J., Decruy, L., Wouters, J., Simon, J. Z., & Francart, T. (2018). Speech
1265 Intelligibility Predicted from Neural Entrainment of the Speech Envelope. *Journal of the*
1266 *Association for Research in Otolaryngology*, 19(2), 181–191.
1267 <https://doi.org/10.1007/s10162-018-0654-z>
- 1268 Vetter, P., Bola, Ł., Reich, L., Bennett, M., Muckli, L., & Amedi, A. (2020). Decoding
1269 Natural Sounds in Early "Visual" Cortex of Congenitally Blind Individuals. *Current Biology*,
1270 30(15), 3039–3044.e2. <https://doi.org/10.1016/j.cub.2020.05.071>

- 1271 Vetter, P., Smith, F. W., & Muckli, L. (2014). Decoding sound and imagery content in early
1272 visual cortex. *Current Biology: CB*, 24(11), 1256–1262.
1273 <https://doi.org/10.1016/j.cub.2014.04.020>
- 1274 Waldo, M., Gerhardt, G., Baker, N., Drebing, C., Adler, L., & Freedman, R. (1992).
1275 Auditory sensory gating and catecholamine metabolism in schizophrenic and normal
1276 subjects. *Psychiatry Research*, 44(1), 21–32. [https://doi.org/10.1016/0165-
1277 1781\(92\)90066-c](https://doi.org/10.1016/0165-1781(92)90066-c)
- 1278 Wang, L., Wu, E. X., & Chen, F. (2020). Robust EEG-Based Decoding of Auditory
1279 Attention With High-RMS-Level Speech Segments in Noisy Conditions. *Frontiers in
1280 Human Neuroscience*, 14. <https://doi.org/10.3389/fnhum.2020.557534>
- 1281 Wang, X., & Xu, L. (2021). Speech perception in noise: Masking and unmasking. *Journal
1282 of Otology*, 16(2), 109–119. <https://doi.org/10.1016/j.joto.2020.12.001>
- 1283 Wolmetz, M., Poeppel, D., & Rapp, B. (2011). What Does the Right Hemisphere Know
1284 about Phoneme Categories? *Journal of Cognitive Neuroscience*, 23(3), 552–569.
1285 <https://doi.org/10.1162/jocn.2010.21495>
- 1286 Zangaladze, A., Epstein, C. M., Grafton, S. T., & Sathian, K. (1999). Involvement of visual
1287 cortex in tactile discrimination of orientation. *Nature*, 401(6753), 587–590.
1288 <https://doi.org/10.1038/44139>
- 1289 Zendel, B. R., West, G. L., Belleville, S., & Peretz, I. (2019). Musical training improves the
1290 ability to understand speech-in-noise in older adults. *Neurobiology of Aging*, 81, 102–115.
1291 <https://doi.org/10.1016/j.neurobiolaging.2019.05.015>
- 1292 Zion Golumbic, E. M., Ding, N., Bickel, S., Lakatos, P., Schevon, C. A., McKhann, G. M.,
1293 Goodman, R. R., Emerson, R., Mehta, A. D., Simon, J. Z., Poeppel, D., & Schroeder, C.
1294 E. (2013). Mechanisms Underlying Selective Neuronal Tracking of Attended Speech at a
1295 'Cocktail Party.' *Neuron*, 77(5), 980–991. <https://doi.org/10.1016/j.neuron.2012.12.037>
- 1296 Zoefel, B. (2018). Speech Entrainment: Rhythmic Predictions Carried by Neural
1297 Oscillations. *Current Biology*, 28(18), R1102–R1104.
1298 <https://doi.org/10.1016/j.cub.2018.07.048>
- 1299 Zorzos, I., Kakkos, I., Ventouras, E. M., & Matsopoulos, G. K. (2021). Advances in
1300 Electrical Source Imaging: A Review of the Current Approaches, Applications and
1301 Challenges. *Signals*, 2(3), 378–391. <https://doi.org/10.3390/signals2030024>