

Distinct roles of SNR, speech Intelligibility, and attentional effort on neural speech tracking in noise

Xiaomin He^{a,b} Vinay S Raghavan^{a,b}, Nima Mesgarani^{a,b*}

^a Department of Electrical Engineering, Columbia University, New York, NY, USA

^b Mortimer B. Zuckerman Mind Brain Behavior Institute, Columbia University, New York, NY, USA

* *Corresponding author*

Correspondence: nima@ee.columbia.edu; xh2369@columbia.edu

Conflicts of Interest: None declared

Funding sources: This work was supported by the National Institutes of Health (NIH) - National Institute on Deafness and Other Communication Disorders (NIDCD) (DC014279) and a grant from Marie-Josée and Henry R. Kravis.

1 **Abstract**

2 Robust neural encoding of speech in noise is influenced by several factors, including signal-to-
3 noise ratio (SNR), speech intelligibility (SI), and attentional effort (AE). Yet, the interaction and
4 distinct role of these factors remain unclear. In this study, fourteen native English speakers
5 performed selective speech listening tasks at various SNR levels while EEG responses were
6 recorded. Attentional performance was assessed using a repeated word detection task, and
7 attentional effort was inferred from subjects' gaze velocity. Results indicate that both SNR and SI
8 enhance neural tracking of target speech, with distinct effects influenced by the previously
9 overlooked role of attentional effort. Specifically, at high levels of SI, increasing SNR leads to
10 reduced attentional effort, which in turn decreases neural speech tracking. Our findings highlight
11 the importance of differentiating the roles of SNR, SI, and AE in neural speech processing and
12 advance our understanding of how noisy speech is processed in the auditory pathway.

13 **Keywords**

14 Neural speech tracking, SNR, speech intelligibility, attentional effort

15 **1. Introduction**

16 The neural encoding of speech in noise is an essential process that enables speech
17 comprehension in complex auditory scenes. Various objective and subjective factors influence
18 how the auditory cortex processes noisy speech. Objective factors include the signal-to-noise
19 ratio (SNR), representing the physical properties of the acoustic signal and its masking by
20 background noise. Speech intelligibility (SI), on the other hand, is a subjective measure that
21 reflects the listener's ability to recognize spoken words and depends not only on SNR but also on
22 the listener's auditory processing capabilities (Nilsson et al., 1994; Sharma et al., 2013).
23 Attentional performance (AP) is another subjective factor that pertains to the listener's ability to

24 selectively concentrate on one speech stream among many and filter out unwanted sounds in
25 complex auditory scenes. Another related yet distinct factor is attentional effort (AE) (Sarter et al.,
26 2006), which involves the cognitive resources expended to focus on the attended talker while
27 ignoring distractions and is influenced by listener's engagement, fatigue, and overall task difficulty
28 (Bruya and Tang, 2018; Sarter et al., 2006; Strauss and Francis, 2017). While these factors are
29 interconnected, they are mechanistically distinct. SNR is an external, quantifiable measure,
30 whereas intelligibility and attention are subjective experiences that vary across individuals, even
31 in identical acoustic settings. This differentiation underscores the complexity of auditory
32 processing and the gaps in our understanding of how these elements collectively influence neural
33 speech encoding.

34

35 Past research has extensively studied the neural encoding of speech in noise, emphasizing the
36 role of SNR and speech intelligibility. Studies demonstrated that increasing SNR generally
37 enhances intelligibility and neural speech encoding (Das et al., 2018; Decruy et al., 2020a; Ding
38 and Simon, 2013; Lesenfants et al., 2019; Vanthornhout et al., 2018). Others have used varying
39 degrees of visual congruency to modulate intelligibility and examined its impact on neural speech
40 encoding (Crosse et al., 2015; lotzov and Parra, 2019). Further work has identified different
41 response components that differentially reflect SNR or intelligibility, such as frequency bands
42 (Etard and Reichenbach, 2019; Vanthornhout et al., 2018), temporal components (Decruy et al.,
43 2020a; Ding and Simon, 2013; Yasmin et al., 2023), and response latency (Yasmin et al., 2023).
44 However, these findings often imply a monotonic relationship between SNR, intelligibility, and
45 neural encoding, which oversimplifies the dynamic interaction among these features (Krueger et
46 al., 2017). For instance, increasing noise levels under specific conditions can enhance neural
47 tracking (Das et al., 2018; Lesenfants et al., 2019), and higher intelligibility does not always
48 correlate with increased neural encoding (Etard and Reichenbach, 2019). Moreover, past
49 research typically used standard speech-in-noise tasks to measure intelligibility, often separating

50 this assessment from the task used to evaluate neural responses. Such intelligibility tasks typically
51 involve asking subjects to repeat short sentences heard in noisy environments (Feng and Chen,
52 2022; Nilsson et al., 1994; Sharma et al., 2013), a method that may only partially capture the
53 complexities of real-world listening due to its limited engagement with challenging factors such as
54 attention span and effort. It has been shown that attentional performance significantly modulates
55 neural speech encoding (Ding and Simon, 2013; Mesgarani and Chang, 2012; O'Sullivan et al.,
56 2015), where SNR can considerably change target speech intelligibility (Brungart et al., 2001) and
57 the attentional effort required to maintain focus on a talker (Cui and Herrmann, 2023; Dimitrijevic
58 et al., 2019; Zekveld et al., 2006). These findings suggest a complex interplay between SNR,
59 intelligibility, attention, and their impact on neural speech encoding (Devocht et al., 2017; Krueger
60 et al., 2017; Yasmin et al., 2023), highlighting a critical gap in our holistic understanding of how
61 these factors individually and collectively shape neural encoding of speech in noise.

62
63 Our study aims to address the need for a comprehensive analysis integrating these dimensions
64 (SNR, SI, and AE) to fully elucidate their combined impact on neural encoding. We examined
65 neural responses to speech in noise through a multifaceted approach incorporating a high-
66 resolution range of SNR values. We used a repeated word detection task (Kirchner, 1958; Laffere
67 et al., 2020; Marinato and Baldauf, 2019), designed to continually monitor subjects' behavior in a
68 manner that integrates attentional performance with the assessment of speech intelligibility,
69 allowing us to capture the variability of attentional engagement in natural listening conditions.
70 Additionally, we estimated attentional effort by analyzing gaze velocity (Ala et al., 2020; Ciccarelli
71 et al., 2019; Gopher, 1973) to understand their collective impact on EEG signals. Our findings
72 advance our holistic understanding of noisy speech processing in the auditory cortex and have
73 practical implications for designing auditory technologies to improve speech perception under
74 challenging listening conditions.

75 **2. Materials and Methods**

76 **2.1. Participants**

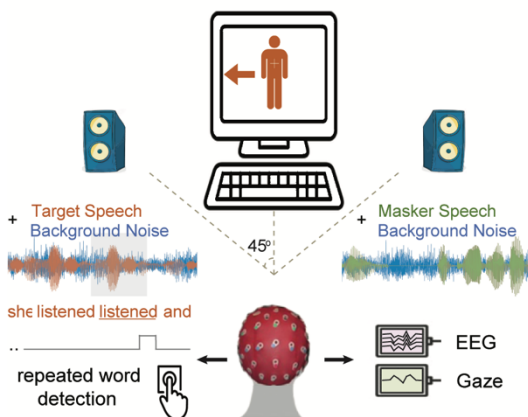
77 Fourteen native American English speakers (7 males; mean \pm standard deviation (SD) age, 24.86
78 \pm 4.4 years) with self-reported normal hearing participated in the experiment. The study followed
79 the protocol approved by the Institutional Review Board of Columbia University (Protocol Number:
80 AAAR7230). Participants were paid for their time as well as a bonus based on their task
81 performance (1-back detection hit rate).

82 **2.2. Experiment Procedures**

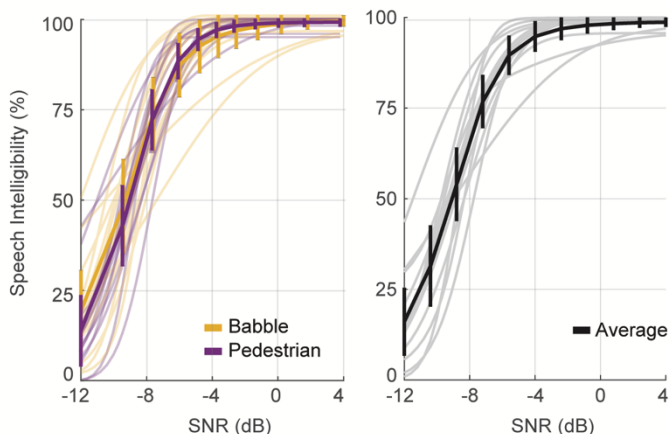
83 **2.2.1. Experiment 1: Measuring Intelligibility by Connected Speech Test**

84 Speech intelligibility (SI) was measured with the Connected Speech Test (CST) (Cox et al., 1987)
85 in experiment 1. Subjects listened to a series of connected short sentences from daily familiar
86 topics with one sentence at a time. The sentences were normalized to 65dB and were covered
87 with noises at different SNRs ranging from -12 dB to 4 dB. Subjects were asked to verbally repeat
88 the words they heard. Stimuli were synthesized by Google Text-To-Speech API (WaveNet) (Oord
89 et al., 2016) with four different voices (2 males and 2 females), and played by two loudspeakers
90 placed at \pm 45 degrees. Experimenters recorded subjects' responding accuracy and regressed
91 for individual SI curves afterward.

A. Task schematic



B. Psychometric curve: Speech Intelligibility vs. SNR



92

93 **Fig 1. Illustration of the task schematic and psychometric curves.** (A) The general task
94 schematic. Subjects were instructed to focus on the target talker while ignoring the masker talker
95 and the background noise. EEG, gaze activities, and button press were recorded while subjects
96 performed the tasks. The subjects press a buzzer when they hear a repeated word in the target
97 stream. (B) Speech intelligibility (SI) is measured using a connected speech task to derive
98 psychometric curves for pedestrian and babble noises. No significant difference appears between
99 the noise types.

100 **2.2.2 Experiment 2: Multi-talker Speech-in-Noise Perception Test**

101 American English podcast stories were synthesized by Google Text-to-Speech API with the same
102 setting as experiment 1. 160 trials of context-continuous stories (average length ~35s) were also
103 played by two loudspeakers placed to ± 45 degrees of subjects (Fig 1A). During each trial,
104 subjects were presented with two speech streams covered by naturalistic background noises
105 (babble or street noise). The target speech was normalized to 65 dB, and its SNR ranged from -
106 12 dB to 4 dB.

107 Subjects were instructed to focus on the speaker whose gender and direction were specified by
108 the icon on the monitor. Three repeated words were inserted in both speech streams. For
109 simplicity, we selected semantically important keywords as repeated words, excluding articles,
110 prepositions, and conjunctions. During the experiment, subjects needed to press the buzzer
111 whenever they captured a repeated word from the target speaker. After each round of 16 trials,
112 experimenters calculated the buzzer responses to the repeated words (1-back detection hit rate)

113 and reported to subjects as feedback. Experimenters also asked subjects to summarize the
114 stories they heard briefly. However, only 1-back detection hit rate was evaluated as a basis for
115 compensation, to avoid unnecessary memory load for subjects.

116 Subjects were also instructed to keep their gaze on the monitor and minimize head movement
117 during each trial. To facilitate tracking of target speech in extremely adverse trials with low SNRs,
118 a 3-second window was used at the start of each trial where masker speech and noise gradually
119 increased to the pre-set SNR. These windows were removed in later analyses.

120 **2.3. Data Acquisition and Preprocessing**

121 In Experiment 1, the word recalling accuracy for each SNR bin was manually recorded for later
122 regressing the *SNR-SI* psychometric curve ([Fig 1B](#); Details also in **2.4.2 Speech Perceptual**
123 **Attributes**).

124 In Experiment 2, buzzer responses to the repeated words, 64-channel EEG, and eye-tracking
125 data were recorded for each trial. Among them, buzzer responses and EEG were recorded by
126 g.HIAMP (g.tec, Australia). Eye tracking data were calibrated and acquired from Tobii Pro Nano
127 (Tobii, Sweden). All data was streamed from Simulink (Mathworks, MA, USA) at 1200 Hz with a
128 60Hz notch. Afterward, EEG data were downsampled to 100Hz with an anti-aliasing filter.
129 Channels with unusual standard deviations were automatically detected and replaced using
130 spherical interpolation of the remaining channels (Delorme and Makeig, 2004; Kang et al., 2015;
131 Perrin et al., 1989).

132 Speech envelopes for both target and masker speakers were firstly extracted by a nonlinear,
133 iterative (NLI) method (Horwitz-Martin et al., 2016) and secondly downsampled to 100Hz to match
134 with the EEG recordings. Finally, each envelope was z-scored to zero mean and unit variance.

135 Blink detection and gaze tracking were completed and preprocessed automatically by Tobii Pro
136 SDK (Tobii, Sweden) with a sampling frequency of 60Hz. The gaze coordinates were normalized
137 to (0,0) and (1,1) within the screen.

138 **2.4. Measurement of objective and subjective features**

139 **2.4.1. Speech objective attribute: Signal-to-Noise ratio (SNR) of target speech.**

140 In both experiments, the volume of target speeches was normalized to 65 dB. The masker speech
141 and bi-channel noises were at equalized volume to form the SNRs distribution from -12 dB to 4
142 dB. The SNRs were computed in the following formula:

$$143 \quad SNR_{target} = 10 \log_{10} \frac{P(target)}{P(masker) + 2P(noise)} \quad (1)$$

144 **P**: power of stimuli

145 The bi-channel noises used in this study for SNR adjustments are:

- 146 • The babble noise was 10-speaker babble derived from the AzBio test (Spahr et al., 2012).
- 147 • The street noise was the pedestrian area recording from CHiME3 (Barker et al., 2015) but
148 with any salient interference removed (e.g. car horn, high-pitch car brake sound, intelligible
149 pedestrians' talking, etc.).

150 Noise audios were truncated as long as the formal trial (~35s per trial).

151 **2.4.2. Speech perceptual attribute: Speech Intelligibility (SI)**

152 Speech Intelligibility (SI) was measured by the Connected Speech Test (Cox et al., 1987). In
153 Experiment 1, experimenters manually filed the subjects' word recalling accuracy for each SNR
154 bin in the range of -12 dB to 4 dB. Then, the psychometric curve between SNRs and word recall

155 accuracy was fitted by *psignifit* toolbox, which implements the maximum-likelihood method
156 described by (Wichmann and Hill, 2001a, 2001b) , and a customized logistic function:

$$SI = lw + \frac{up - lw}{1 + \exp^{-gr*(SNR-ths)}} \quad (2)$$

157 **lw**: lower bound (defined to approach 0); **up**: upper bound (defined to approach 1); **gr**:
158 growth rate; **SNR**: signal-to-noise ratio, range from -12 to 4 dB; **ths**: threshold when SI =
159 50%, defined in the range of SNR.

160 Among the two approaches, the one producing a regressed curve with higher R^2 and lower $RMSE$
161 was selected. From the selected curve, the corresponding SI for each SNR in Experiment 2 was
162 read.

163 2.4.3. Attention Measures

- 164 • Attentional performance (AP): Single-trial 1-back detection hit rate (HR)

165 As mentioned above, single-trial 1-back detection hit rate (HR), as a measure of subjects'
166 performance in terms of attention focus for each trial, was computed by the buzzer-hitting
167 performance (Kirchner, 1958; Laffere et al., 2020; Marinato and Baldauf, 2019). There
168 were 3 words inserted for each trial. Therefore, the range of HR was $[0, \frac{1}{3}, \frac{2}{3}, 1]$. Intuitively,
169 high HR in a trial indicates a better attentional performance.

- 170 • Attentional effort (AE): Gaze Velocity (GV)

171 Concentration periods are associated with suppressed irrelevant physiological activities,
172 evidenced by reduced ocular movements (saccade and micro-saccade rate) and blink
173 rates, as well as prolonged fixation (Abeles et al., 2020; Braga et al., 2016; Contadini-

174 Wright et al., 2023; Cui and Herrmann, 2023). Oculomotor activity, with its anatomical
175 overlap with the attention-related network (Corbetta et al., 1998), is, therefore, a valuable
176 metric for evaluating attentional effort with superiority in stability across age (Bruenech,
177 2008), and the relationship has been justified by (Ala et al., 2020; Ciccarelli et al., 2019;
178 Gopher, 1973). In this paper, we use averaged gaze velocity (GV) to quantify attentional
179 effort for each trial, as it reflects the overall intensity of oculomotor activity, including
180 saccade and micro-saccade. A higher GV, indicating more frequent oculomotor activity,
181 suggests reduced attentional effort (Ala et al., 2020; Ciccarelli et al., 2019; Gopher, 1973).

182 The gaze coordinates were recorded and normalized between (0,0) and (1,1) by the Tobii
183 Pro Nano screen-based eye tracker. To calculate actual gaze angular velocity, we first
184 restored relative coordinates to screen size, then computed and averaged the absolute
185 value of the derivative of gaze coordinates over time within each trial. Finally, using the
186 (approximately) 0.6m distance from the subject's seat to the screen, we calculated gaze
187 velocity (GV) using trigonometric functions (Diaz et al., 2013).

188 **2.4. Attention Decoding**

189 The classical approach for auditory attention decoding (AAD) is to model the linear projection
190 between neural electrophysiological recordings and speech features (O'Sullivan et al., 2015),
191 such as speech envelope. Once the model is trained, AAD correlates—the correlations between
192 the speech features and their reconstruction from neural recordings—are compared for target and
193 masker speech to decode auditory attention. Over and above the conventional forward or
194 backward regularized linear model that applies transformation solely on one side of this projection
195 (either neural signal or speech features), the Canonical Correlation Analysis (CCA) approach
196 transforms both neural recordings and speech features for significantly better correlations scores
197 (Dähne et al., 2015; de Cheveigné et al., 2018).

198 We adopt the CCA algorithms for target speech decoding. Speech envelope, which is the slow
199 modulation of speech and proved to be feasible for neural speech tracking with EEG (Horton et
200 al., 2014; Mesgarani et al., 2009; O’Sullivan et al., 2015), was extracted from both target and
201 masker speech. Envelopes for clean speech and the multi-talker EEG recordings were first
202 downsampled to 100Hz for a sampling rate match. Second, stimuli and neural recordings were
203 windowed for overlapping receptive fields. Time-lagged matrices were produced for envelopes
204 and EEG recordings. For EEG, the receptive fields were 400 ms and for speech envelope, the
205 receptive fields were 200ms. Third, for each subject, subject-wise CCA-based linear models for
206 both speeches were trained in a leave-one-out cross-validation setting.

207 The stimuli-response mapping is quantified by the trained model and evaluated by Pearson’s
208 correlation between transformed stimuli and neural responses. As the target and masker stimuli
209 are not identically encoded in the brain (Ding and Simon, 2012a), we computed this correlation
210 for both the target and masker speech. The correlation for the target is referred to as r_T , and for
211 the masker speech is r_M . We also defined r_D as their difference ($r_D = r_T - r_M$) to quantitatively
212 represent the different intensities of neural entrainment caused by selective attention. $r_D > 0$
213 indicates a successful attention-decoded trial.

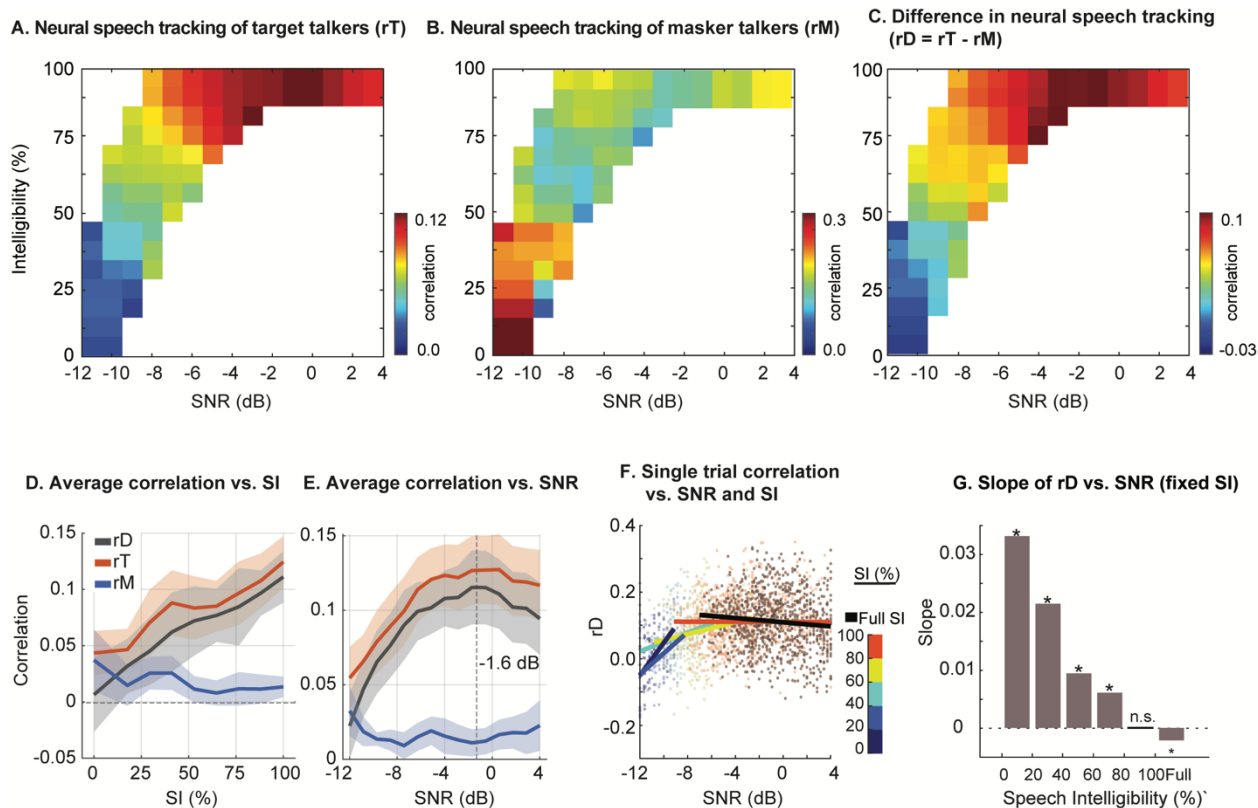
214 Moreover, to investigate the neural modulation pattern under varying speech conditions, we
215 estimated temporal response functions (Ding and Simon, 2012b; Lalor et al., 2009) of target
216 speech using regularized linear regression. This approach minimizes the mean-square error
217 between the actual neural recordings and the predicted values. The training and prediction
218 processes for each subject were also conducted in a leave-one-out cross-validation fashion using
219 the mTRF toolbox (Crosse et al., 2016).

220 **3. Result**

221 Fourteen participants were instructed to perform selective listening tasks, focusing on a target
222 speaker's speech (attended stream) while ignoring a masker speaker (non-target, unattended)
223 and background noises. We recorded 64-channel EEG, gaze velocity, and buzzer press
224 responses to capture the participants' neural and behavioral responses in real-time ([Fig 1A](#)). For
225 each participant, speech intelligibility (SI) was measured using a connected speech test (Cox et
226 al., 1987) prior to the actual experiment. As the difference in psychometric curves between
227 different types of noise was negligible ([Fig 1B](#), left), we adopted an average psychometric curve
228 for each subject to streamline the analysis ([Fig 1B](#), right).

229 **3.1. Distinct Impacts of SNR and SI on Neural Speech Tracking: SI Enhances Tracking,** 230 **While High SNR Reduces It**

231 To measure the strength of neural tracking for target and masker speech, we trained CCA-based
232 linear models to quantify neural speech tracking for each talker. We computed Pearson's
233 correlation between the transformed speech envelope and neural recordings to measure the
234 strength of neural speech tracking for the target (r_T) and masker (r_M) speech. The difference in
235 correlation between target and masker speech ($r_D = r_T - r_M$) was used as a single measure to
236 reflect how well participants followed the target speech while suppressing the masker speech.



237

238 **Fig 2. The relationship between neural speech tracking for target and masker speech with**
 239 **variables SNRs and SI. (A) rT: the correlation of the target neural speech tracking across**
 240 **different SNR and SI values; (B) rM: the correlation of masker neural speech tracking across**
 241 **different SNR and SI values; (C) rD: the difference between neural speech tracking of target and**
 242 **masker speech streams, rD = rT - rM. (D) Average neural speech tracking across SI. (E) Averaged**
 243 **neural speech tracking across SNR. (F) rD across SNR for different groups of SI, with a linear**
 244 **line fitted to the data sample distribution of rD. Scatters and the fitted line are color-coded by SI.**
 245 **(G) Slope of rD vs. SNR with SI fixed. Significant slopes with 95% confidence intervals not**
 246 **containing 0 are marked with asterisks.**
 247

248 To accurately assess the impact of SNR and SI on target and masker neural speech tracking, it
 249 is crucial to distinguish between these two highly correlated factors. We addressed this by
 250 analyzing their effect on neural speech tracking as a function of both SI and SNR. Our analysis
 251 revealed distinct patterns in how SNR and SI affect neural speech tracking. Fig 2A-2C show the
 252 averaged neural tracking correlations across subjects for target (rT) and masker speech (rM) and
 253 their difference (rD) for different SI and SNR values. We found that while increasing SNR and SI
 254 generally increase rT (enhanced neural tracking of the target speech) and decrease rM
 255 (suppressed neural tracking of the masker speech), this relationship shifts when SI is sufficiently

256 high. Specifically, under high SI conditions (i.e., $SI > 80\%$), increasing SNR reduces rT and
257 increases rM , indicating decreased neural tracking of the target speech while increasing the
258 tracking of the masker speaker. The average plots across SI and SNR in [Fig 2D and 2E](#) further
259 illustrate these effects, showing that SI has a nearly monotonic relationship with rT , rM , and rD ,
260 while SNR's impact on these values reverses beyond approximately -1.6 dB. This indicates that
261 in easier listening conditions, when the target talker is highly intelligible, increasing the SNRs of
262 the masker talker can paradoxically reduce its neural speech tracking. In [Fig 2F](#), we quantized
263 SI into 6 bins from 0 to full intelligibility, each denoted by a different color. Trials categorized as
264 'Full SI' are marked in black and represent instances where SI reached its plateau on individual
265 psychometric curves. [Fig 2G](#) illustrates the linear relationship between SNR and rD within each
266 SI bin, revealing a progression from positive to neutral to negative correlation between rD and
267 SNR as SI increases. In summary, these results demonstrate that while SNR and SI are strongly
268 correlated, they have distinct and sometimes opposing effects on neural speech tracking, which
269 we will explore in greater detail in the next section

270 **3.2. Increased SNR Leads to Reduced Attentional Effort and Neural Speech Tracking**

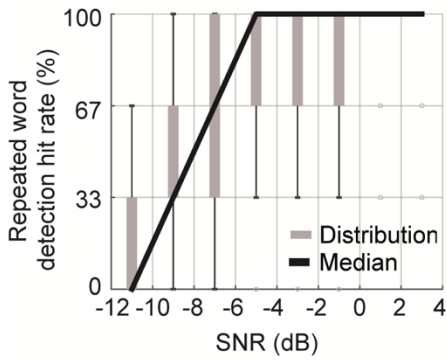
271 Our findings of a negative relationship between SNR and neural speech tracking (rD) under high
272 SI conditions suggest a secondary effect of SNR. One possibility is that increasing SNR may lead
273 to decreased attentional performance (AP) and/or attentional effort (AE), which could
274 consequently reduce rD . Specifically, AP is the actual performance outcome while AE refers
275 to the cognitive resources required to maintain attention, including motivation and resource
276 allocation (Pashler et al., 2001; Sarter et al., 2006). To investigate this, we used the hit rate of
277 word repetition detection task (HR) as an ongoing measure of attentional performance, where a
278 high HR indicates heightened attentional performance (Kirchner, 1958; Laffere et al., 2020;
279 Marinato and Baldauf, 2019). Additionally, we used gaze velocity (GV), to measure oculomotor
280 activity, as an indicator of ongoing attentional effort (AE); notably, a low GV suggests increased

281 attentional effort (Ala et al., 2020; Ciccarelli et al., 2019; Gopher, 1973). (see “**4. Materials and**
282 **Methods**” - “**2.4.3 Attention Measures**”). Fig 3A and 3B illustrate how SNR influences these
283 two attention-related metrics. In Fig 3A we observe that the median HR levels off after -5 dB,
284 suggesting that our attentional performance metric reaches a maximum beyond this SNR
285 threshold, limiting its ability to explain changes in rD in this range.

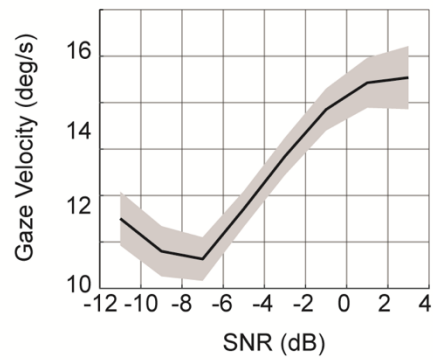
286
287 Conversely, Fig 3B shows that GV significantly increases at higher SNRs, where speech is highly
288 intelligible. This indicates that the attentional effort required to maintain focus on the target talker,
289 which is inversely related to GV, substantially decreases when SNR is sufficiently high. The
290 reduced attentional effort, reflected by increased ocular activity, correlates with the decline in rD
291 (Fig 3D, Pearson correlation test, $c = -0.15$, $p < 1e-5$).

292
293 We conducted a more detailed analysis to explore the impact of attention-related features on rD.
294 Fig 3C shows the modulation of rD by attentional performance, measured by HR. Higher HR was
295 found to be associated with higher rD. Moreover, this change is not continuous due to the discrete
296 nature of the behavioral response (three repeated words in each trial). The variation of rD with
297 attentional effort, measured by GV, is shown in Fig 3D. This analysis reveals that trials with lower
298 rD also exhibit more frequent gaze activity, suggesting that a decrease in attentional effort is
299 correlated with decreased neural speech tracking.

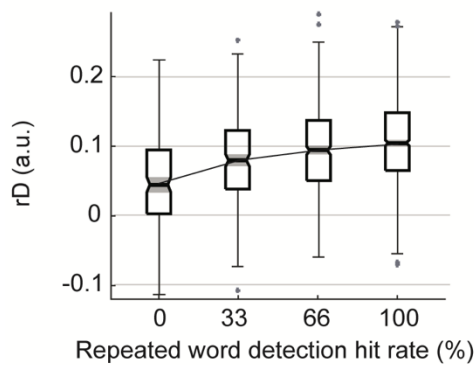
A. HR vs. SNR



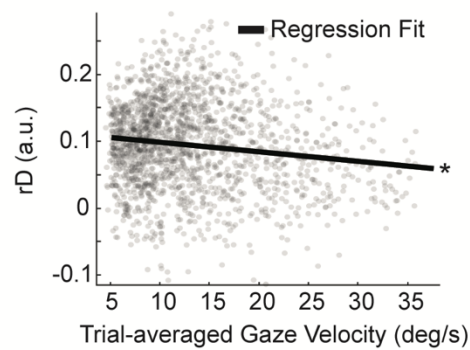
B. Gaze Velocity vs. SNR



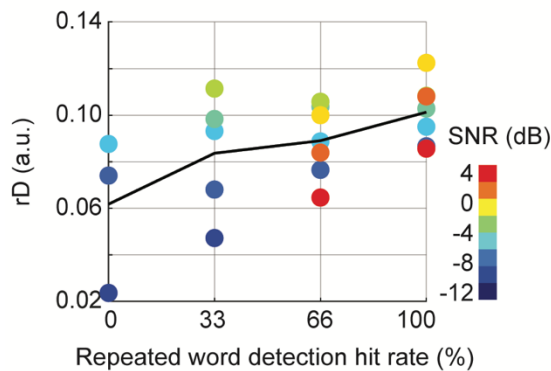
C. rD vs. HR



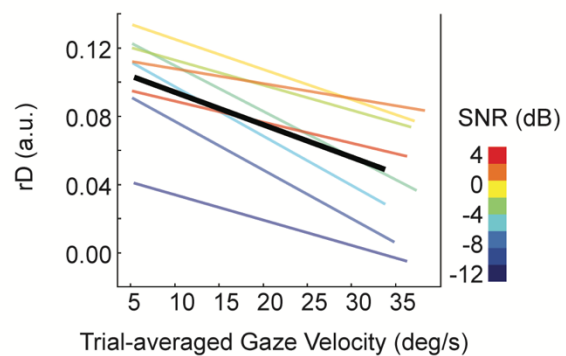
D. rD vs. Gaze Velocity



E. rD vs. HR (fixed SNR)



F. rD vs. Gaze Velocity (fixed SNR)



300

301 **Fig 3. Attentional performance and attentional effort.** (A) single-trial repeated word hit rate
 302 (HR) increases with SNR. (B) Gaze velocity (GV) increases with SNRs. (C) Interaction between
 303 rD and HR: distribution of rD across groups of HR. Medians and confidence intervals are marked
 304 with black lines and gray shades, respectively. Significant differences between groups are found
 305 ($p < 0.05$, Kruskal-Wallis test, Bonferroni-corrected). (D) Interaction between rD and GV: GV
 306 negatively correlates with rD ($c = -0.15$, $p < 1e-5$). The black line marks the regressed linear fit with
 307 significant slope and intercept: $rD = -0.0014 * GV + 0.1126$. (E) Interaction between rD and HR
 308 when fixing SNR. (F) Interaction between rD and GV when fixing SNR.

309

310 To ensure these results apply across the entire range of SNR values, we repeated the same
311 analyses separately for each SNR bin from -12 dB to 4 dB, as shown in [Fig 3E and 3F](#). The results
312 confirm that the observed positive correlation between rD and HR ([Fig 3E](#)) and the negative
313 correlation between rD and GV ([Fig 3F](#)) is consistent across different SNR values. Hence, while
314 the effect of attentional effort on rD becomes more visible in higher SNRs, these two show the
315 same relationship in all SNR conditions. In summary, [Fig 3](#) shows that attentional effort decreases
316 with SNR, meaning subjects exert less effort in easier trials, which also corresponds to decreased
317 target speech neural tracking.

318 **3.3. Modeling the Interactions Between SNR, Speech Intelligibility, and Attentional Effort** 319 **on Neural Speech Tracking**

320 Given that our previous results indicate that multiple interacting variables influence rD, we used
321 a computational model to elucidate these complex relationships. Specifically, we fitted a linear
322 model to predict rD for each trial from that trial's objective (SNR) and subjective (SI and GV)
323 measurements (Adjusted R^2 : 0.151; F-statistic vs. constant model: 48.2, $p < 0.001$). The main
324 effects and interaction terms are depicted in [Fig 4A](#). This analysis shows that SI positively
325 influences rD, while GV negatively affects rD. Interestingly, the direct influence of SNR on rD is
326 not significant when interaction terms are included. This suggests an indirect influence of rD by
327 SNR through the modulation of attentional effort and SI ([Fig 4A](#)). Specifically, increasing SNR
328 improves SI and reduces attentional effort. The opposing effects of SI and attentional effort on rD
329 could explain the non-linear relationship observed in [Fig 2E](#).

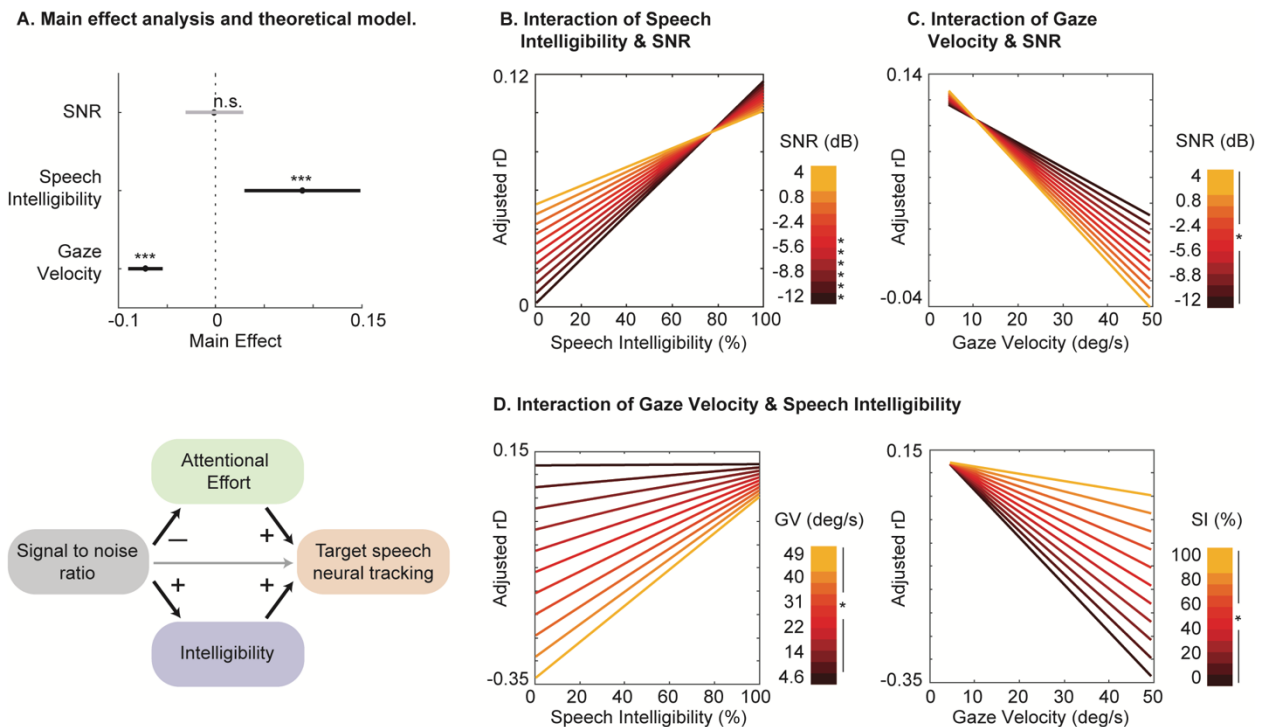
330

331 [Fig 4B-4D](#) further illustrate the interaction between these features. [Fig 4B](#) shows that SI has a
332 positive impact on rD irrespective of SNR levels. However, as SNR increases, the effect of SI and
333 its significance lessens, likely due to SI approaching the ceiling at higher SNRs. [Fig 4C](#)
334 demonstrates the interaction between SNR and GV, showing that GV negatively impacts rD

335 across all SNR levels. The change in the slope with increasing SNR suggests that rD is more
 336 sensitive to GV at higher SNRs. This increased sensitivity indicates that attentional effort plays a
 337 more significant role in shaping the neural target tracking at higher SNRs especially after
 338 maximum intelligibility.

339

340 Fig 4D shows how SI and GV differentially influence rD. Fig 4D (left) shows that the impact of SI
 341 on rD varies with different levels of attentional effort. With increased attentional effort (low GV),
 342 the influence of SI on rD decreases, highlighting the primary role of attention in shaping neural
 343 speech tracking. Conversely, Fig 4D (right) shows that the negative impact of increased GV on
 344 rD depends on SI. Attentional effort has the highest influence on rD in less intelligible conditions,
 345 where increased performance may attempt to compensate for the heightened difficulty of the
 346 listening task.



347

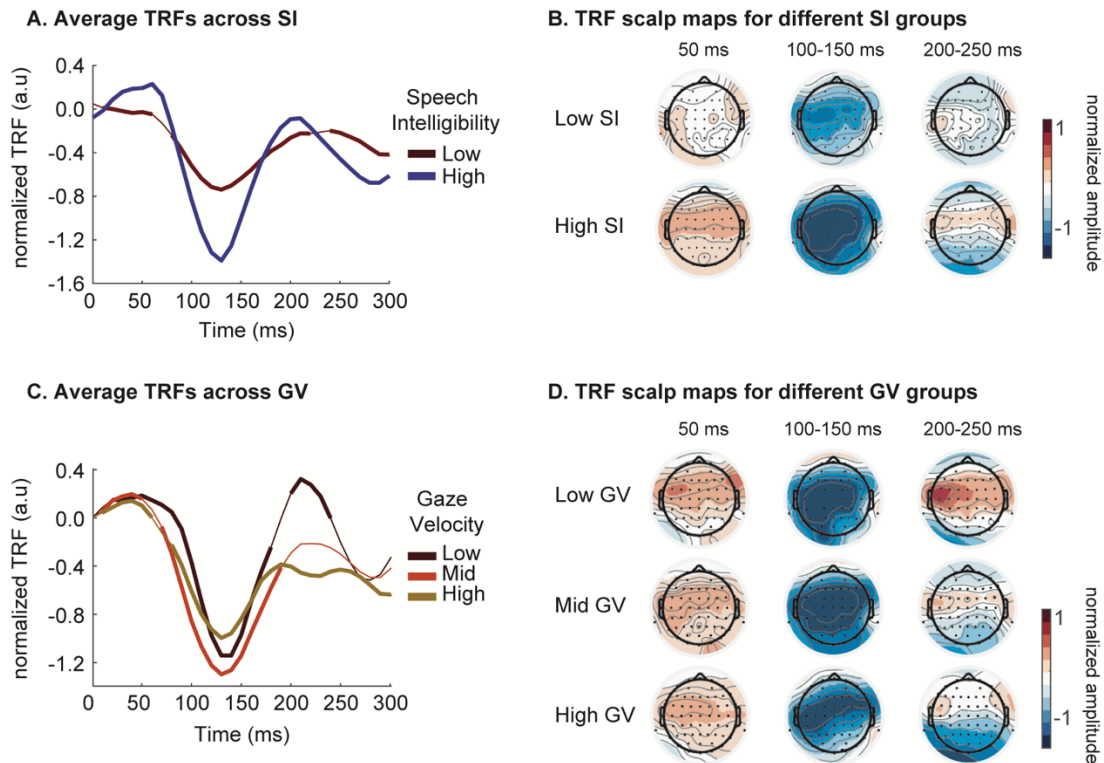
348 **Fig 4. Interaction of various factors.** (A). The main effect analysis of a linear model, and a
 349 hypothesized model of feature interactions. Speech intelligibility and gaze velocity exhibit a
 350 significant effect on rD ($p < 0.001$, t -test), while SNR does not. (B) Interaction effect between SNR

351 *and SI of the fitted linear model. (C) Interaction effect between SNR and GV of the fitted linear*
352 *model. (D) Interaction effect between SI and GV of the fitted linear model.*

353 **3.4 Temporal and Spatial Dynamics of Neural Responses Under Varying Speech** 354 **Intelligibility and Attentional Effort**

355 To investigate how the timing and spatial distribution of neural response patterns change under
356 different SI and GV conditions, we calculated the temporal response functions (TRFs) (Ding and
357 Simon, 2012b; Lalor et al., 2009) for target speech in different listening conditions. TRFs capture
358 the brain's temporal dynamics in response to continuous auditory stimuli, reflecting the
359 relationship between the EEG signal and the speech envelope over lags at different electrodes.
360 Normalized TRFs for target speech, averaged across all channels, are shown in [Fig 5A and 5B](#).
361 The group with low SI (< 50%) exhibits weaker early components TRF₅₀ (positive, around 50 ms)
362 across the scalp, especially in the temporal and central regions. Additionally, the low SI group
363 shows reduced attention-related TRF components TRF₁₀₀ (negative, around 100 ms) and TRF₂₀₀
364 (positive, around 200 ms), indicating reduced selectivity for target speech and less suppression
365 of masker speech components (Ding and Simon, 2012b; Fiedler et al., 2019).

366
367 Attentional effort, measured inversely by GV, also impacts the TRFs. While the acoustics-
368 modulated early components TRF₅₀ remain consistent across different GV, a significant difference
369 emerges for higher-level, attention-related components around TRF₁₀₀ and TRF₂₀₀ ([Fig 5C and](#)
370 [5D](#)). In trials with lower attentional effort (high GV, GV>66.7% percentile), TRF₁₀₀ responses
371 decrease in posterior electrodes ([Fig 5D](#)). For TRF₂₀₀, only the group with low GV (GV<33.3 %
372 percentile) shows strong activation in the anterior and central areas.



373

374 **Fig 5. Normalized Temporal Response Functions (TRFs) under different Speech**
375 **Intelligibility (SI) or Gaze Velocity (GV) levels. (A)** The normalized TRFs under low (<50%) and
376 **high (>50%) SI. TRFs are averaged across all EEG channels. Significant temporal components**
377 **(compared to the chance level, t -test, $p < 0.05$) are marked with thick lines. (B)** The topography of
378 **normalized TRFs amplitude at three critical time points (50ms, 100~150ms, and 200~250ms),**
379 **under different levels of SI. (C)** The normalized TRFs under low (<33%), mid (33%~67%), and
380 **high (>67%) GV. TRFs are averaged across all EEG channels. Significant temporal components**
381 **(compared to the chance level, t -test, $p < 0.05$) are marked in thick lines. (D)** The topography of
382 **normalized TRFs amplitude at three critical time points (50ms, 100~150ms, and 200~250ms),**
383 **under different levels of GV.**

384 4. Discussion

385 We demonstrate that neural tracking of target speech is influenced by both objective (signal to
386 noise ratio) and subjective (speech intelligibility, attentional performance and effort) factors in
387 distinct ways. As speech intelligibility increases, the positive effect of improving SNR on neural
388 tracking of target speech diminishes. Specifically, in conditions where speech is highly intelligible,
389 further increases in SNR decrease neural speech tracking. We propose that this decrease is
390 caused by the reduced attentional effort required to focus on the target speech. Our findings show
391 that gaze velocity, a measure proposed for quantifying attentional effort, effectively explains this

392 reduction in neural speech tracking accuracy. Together, our findings suggest a complex
393 interaction between speech intelligibility and attentional effort mediated by SNR in shaping the
394 neural representation of speech in noise.

395 **4.1. Distinct Contributions of SI and SNR to Neural Speech Tracking**

396 Despite the significant correlation between speech intelligibility (SI) and signal-to-noise ratio
397 (SNR), our findings reveal distinct impacts of each on neural tracking of target speech as
398 measured by EEG. While SI and SNR are often discussed together due to their high correlation,
399 they affect neural speech tracking through different mechanisms. Previous studies, including
400 those by (Das et al., 2018; Itzov and Parra, 2019; Vanthornhout et al., 2018), primarily
401 investigated the impact of SI on neural tracking at different SNRs. Our results, however,
402 underscore the importance of differentiating the effects of SI and SNR.

403 Objective acoustic features like SNR drive bottom-up processing in speech perception. In
404 contrast, SI is shaped by an individual's bottom-up perception and top-down processing
405 capabilities and strategies for allocating cognitive resources or restoring features masked by noise
406 (Raghavan et al., 2023). This distinction highlights how SNR and SI contribute differently to neural
407 processing. For example, the activity in brain regions responsible for top-down processing is
408 increased when bottom-up processing was impaired by degraded speech (lower SNRs) (Zekveld
409 et al., 2006). Several previous studies also suggested different modulation of neural speech
410 tracking by objective or perceptual speech attributes. Previous research has shown that while
411 SNR non-linearly modulates the amplitude of the temporal response function, changes in neural
412 latency align more closely with variations in SI (Yasmin et al., 2023). In contrast, the relationship
413 between neural speech tracking and SI does not exhibit such non-linearity (Decruy et al., 2020a).
414 Instead, the different metrics of neural speech tracking accuracy show slightly dissimilar
415 correlations with SNR and SI, but this mismatch has not been explained (Nogueira and

416 Dolhopiatenko, 2022). Our findings go beyond these observations by dissociating the interplay
417 between SI and SNR and their impact on neural speech tracking. We provide evidence that SNR
418 contributes indirectly to neural speech tracking by modulating attentional effort and speech
419 intelligibility. This supports the notion of an indirect contribution of SNR to neural speech tracking,
420 as suggested by (Etard and Reichenbach, 2019). It is important to note that our measure of neural
421 speech tracking is based on EEG recordings, which reflect scalp potentials from large populations
422 of neurons but do not provide the fine-grained detail available from invasive or single-neuron
423 recordings. Studies using invasive techniques in animals and humans have shown that noise-
424 invariant representations gradually develop along the auditory pathway (Kell and McDermott,
425 2019; Mesgarani et al., 2014; Rabinowitz et al., 2013), with lower areas representing the noise
426 and higher areas filtering it out. Our findings highlight how the combined effects of these
427 interactions manifest in scalp EEG signals which is critical as EEG is the most widely used
428 measure to study speech in noise in normal hearing, hearing impaired, and aging individuals (Di
429 Liberto et al., 2022; Fuglsang et al., 2020; Mesik et al., 2021).

430 **4.2. The Role of Attention in Neural Speech Tracking**

431 Attention plays a crucial role in how the brain tracks attended speech. The acoustic characteristics
432 of speech can influence attention levels and, consequently, the accuracy of neural speech
433 tracking (Ding and Simon, 2012a; lotzov and Parra, 2019; Mesgarani and Chang, 2012; Power et
434 al., 2012; Vanthornhout et al., 2019a; Zion Golumbic et al., 2013). Our study examined two
435 aspects of attention: attentional performance and attentional effort. Attentional effort refers to the
436 cognitive resources required to maintain attention, including motivation and resource allocation,
437 while attentional performance is the actual outcome (Pashler et al., 2001; Sarter et al., 2006). We
438 assessed attentional performance using repeated word hit rate (HR). We inferred attentional effort
439 by measuring gaze velocity (GV) (Ala et al., 2020; Ciccarelli et al., 2019; Gopher, 1973). In easier
440 listening conditions (SNR > -1.6 dB), we observed a significant reduction in neural speech tracking

441 accuracy with increasing SNR. This finding aligns with studies by (Das et al., 2018; Lesenfants et
442 al., 2019), which noted decreased neural speech tracking accuracy from mildly noisy to clean
443 conditions. Our study further investigates this paradoxical relationship by measuring ocular
444 activity as an approximation of attentional effort. The identified negative interaction between
445 ocular activity and task difficulty was also illustrated by (Contadini-Wright et al., 2023; Cui and
446 Herrmann, 2023; Herrmann and Ryan, 2024). In contrast, attentional performance, measured by
447 HR, shows a limited correlation with task difficulty. These results suggest that the reduction in
448 neural target speech tracking can be more accurately attributed to changes in attentional effort
449 rather than variations in attentional performance. Our findings also align with previous research
450 indicating that increased eye movement activity reflects less suppression of task-irrelevant
451 psychological activity, impairing information processing such as selective neural speech
452 perception (Abeles et al., 2020; Braga et al., 2016; Cui and Herrmann, 2023). More importantly,
453 we provide a potential explanation for the reduced neural speech tracking in easier listening
454 conditions, as also reported by (Das et al., 2018; Hauswald et al., 2022; Lesenfants et al., 2019).
455 Note that this decreasing effect exists across SNRs, not only in the high SNR listening conditions.
456 Attentional effort and attentional performance exhibit distinct characteristics despite their
457 interconnectedness (Bruya and Tang, 2018). Our study also supports differentiating the
458 modulation of neural entrainment between attentional effort and attentional performance, similar
459 to the findings of (Dai and Shinn-Cunningham, 2016), which showed that selective attention could
460 modulate the strength of cortical event-related potential but not change the attentional
461 performance. It is also worth mentioning studies that have demonstrated increased neural speech
462 tracking in older populations and subjects with hearing impairment (Decruy et al., 2020b, 2019).
463 Our study offers an explanation for these observations: increased task difficulty in these subject
464 populations elevates attentional effort, thereby enhancing neural speech tracking. To test this
465 hypothesis, we propose measuring differences in gaze velocity between populations or adjusting

466 the SNR to identify the threshold at which neural speech tracking declines relative to normal
467 hearing subjects, estimated in our study at approximately -1.6 dB.

468 **4.3 Modeling the Interplay of Speech Intelligibility and Attentional Effort on Neural Speech** 469 **Tracking**

470 In analyzing the interaction of various features on predicting neural speech tracking, we found
471 that both speech intelligibility (SI) and gaze velocity (GV) have significant effects, while signal-to-
472 noise ratio (SNR) does not. Supplementary analyses and comparisons of temporal response
473 functions (TRFs) for significant features (SI and GV) revealed that SI influences both acoustic-
474 related (TRF₅₀) (Ding and Simon, 2013) and attention-related components (TRF₁₀₀, TRF₂₀₀) (Ding
475 and Simon, 2012b; Fiedler et al., 2019) of neural speech tracking, consistent with previous studies
476 (Chen et al., 2023; Muncke et al., 2022). Notably, the modulation of early response (TRF₅₀) may
477 be attributed to the combined effect of SNR and SI, as shown in prior study, where a lower SNR
478 at the same SI resulted in reduced TRF₅₀ amplitude (Verschueren et al., 2020). In contrast, GV,
479 as an indicator of attentional effort, only modulates attention-related components, specifically the
480 activation area of TRF₁₀₀ and the intensity of TRF₂₀₀. These components are closely associated
481 with the top-down process of directing mental resources toward the target of interest (Fritz et al.,
482 2007; Kong et al., 2014; Vanthornhout et al., 2019b). These findings suggest that while SI affects
483 multiple aspects of neural speech processing, GV's influence is limited to the attentional
484 mechanisms.

485 From the detailed analyses of the interaction among features, we proposed a model based on our
486 finding that AE and SI show counterbalancing effect on neural speech tracking as SNR increases,
487 with the dominant factor shifting from SI to AE. The proposed model is able to explain the widely
488 observed non-linearity between task demands and neural speech tracking (Das et al., 2018;
489 Hauswald et al., 2022; Lesenfants et al., 2019), and also provides an explanation for the increased

490 speech tracking in hard of hearing and aging populations (Decruy et al., 2020b, 2019), for which
491 the increased task difficulty results in an increased attentional effort.

492 There are several limitations to consider while interpreting our results. As EEG signals provide
493 only a broad overview of cortical activity, complementary neuroimaging techniques would be
494 needed to fully characterize the encoding of noisy speech in various cortical and subcortical
495 auditory regions. Additionally, our measure of attentional effort is indirect. While used extensively
496 in the field (Ala et al., 2020; Ciccarelli et al., 2019; Gopher, 1973), gaze velocity is only an
497 approximation of the cognitive resources that are used to maintain focus. Finally, our measure of
498 attentional performance is sparse, as we cannot rule out the possibility that the listeners lose
499 focus in between the repeated words. Future research is needed to explore more direct methods
500 to measure cognitive load and attentional performance, and to expand these findings to aging
501 and hard of hearing population.

502 In summary, our study demonstrates that the neural tracking of target speech is influenced by
503 SNR, speech intelligibility, and attentional performance and attentional effort, with distinct and
504 sometimes opposing effects. By disentangling the roles of attentional performance and effort, we
505 provide a clearer understanding of how these factors interact to shape neural speech processing.
506 Beyond their scientific impact, these insights also have important implications for developing
507 auditory technologies and strategies to improve speech perception in noisy environments.

508 **Author Contributions**

509 XH, VR, and NM conceived the project. XH and NM designed the experiment and analyzed the
510 data. XH and NM wrote the manuscript, and all authors provided feedback and revisions.

511 **Reference**

- 512 Abeles, D., Amit, R., Tal-Perry, N., Carrasco, M., Yuval-Greenberg, S., 2020. Oculomotor
513 inhibition precedes temporally expected auditory targets. *Nat. Commun.* 11, 3524.
514 <https://doi.org/10.1038/s41467-020-17158-9>
- 515 Ala, T.S., Graversen, C., Wendt, D., Alickovic, E., Whitmer, W.M., Lunner, T., 2020. An
516 exploratory Study of EEG Alpha Oscillation and Pupil Dilation in Hearing-Aid Users During
517 Effortful listening to Continuous Speech. *PLOS ONE* 15, e0235782.
518 <https://doi.org/10.1371/journal.pone.0235782>
- 519 Barker, J., Marxer, R., Vincent, E., Watanabe, S., 2015. The third 'CHiME' speech separation and
520 recognition challenge: Dataset, task and baselines, in: 2015 IEEE Workshop on Automatic
521 Speech Recognition and Understanding (ASRU). Presented at the 2015 IEEE Workshop
522 on Automatic Speech Recognition and Understanding (ASRU), pp. 504–511.
523 <https://doi.org/10.1109/ASRU.2015.7404837>
- 524 Braga, R.M., Fu, R.Z., Seemungal, B.M., Wise, R.J.S., Leech, R., 2016. Eye Movements during
525 Auditory Attention Predict Individual Differences in Dorsal Attention Network Activity. *Front.*
526 *Hum. Neurosci.* 10.
- 527 Bruenech, J.R., 2008. Age-Related Changes in the Oculomotor System, in: Cavallotti, C.A.P.,
528 Cerulli, L. (Eds.), *Age-Related Changes of the Human Eye*. Humana Press, Totowa, NJ,
529 pp. 343–373. https://doi.org/10.1007/978-1-59745-507-7_20
- 530 Brungart, D.S., Simpson, B.D., Ericson, M.A., Scott, K.R., 2001. Informational and energetic
531 masking effects in the perception of multiple simultaneous talkers. *J. Acoust. Soc. Am.*
532 110, 2527–2538. <https://doi.org/10.1121/1.1408946>
- 533 Bruya, B., Tang, Y.-Y., 2018. Is Attention Really Effort? Revisiting Daniel Kahneman's Influential
534 1973 Book *Attention and Effort*. *Front. Psychol.* 9, 1133.
535 <https://doi.org/10.3389/fpsyg.2018.01133>

- 536 Chen, Y.-P., Schmidt, F., Keitel, A., Rösch, S., Hauswald, A., Weisz, N., 2023. Speech
537 intelligibility changes the temporal evolution of neural speech tracking. *NeuroImage* 268,
538 119894. <https://doi.org/10.1016/j.neuroimage.2023.119894>
- 539 Ciccarelli, G., Nolan, M., Perricone, J., Calamia, P.T., Haro, S., O’Sullivan, J., Mesgarani, N.,
540 Quatieri, T.F., Smalt, C.J., 2019. Comparison of Two-Talker Attention Decoding from EEG
541 with Nonlinear Neural Networks and Linear Methods. *Sci. Rep.* 9, 11538.
542 <https://doi.org/10.1038/s41598-019-47795-0>
- 543 Contadini-Wright, C., Magami, K., Mehta, N., Chait, M., 2023. Pupil Dilation and Microsaccades
544 Provide Complementary Insights into the Dynamics of Arousal and Instantaneous
545 Attention during Effortful Listening. *J. Neurosci.* 43, 4856–4866.
546 <https://doi.org/10.1523/JNEUROSCI.0242-23.2023>
- 547 Corbetta, M., Akbudak, E., Conturo, T.E., Snyder, A.Z., Ollinger, J.M., Drury, H.A., Linenweber,
548 M.R., Petersen, S.E., Raichle, M.E., Essen, D.C.V., Shulman, G.L., 1998. A Common
549 Network of Functional Areas for Attention and Eye Movements. *Neuron* 21, 761–773.
550 [https://doi.org/10.1016/S0896-6273\(00\)80593-0](https://doi.org/10.1016/S0896-6273(00)80593-0)
- 551 Cox, R.M., Alexander, G.C., Gilmore, C., 1987. Development of the Connected Speech Test
552 (CST): *Ear Hear.* 8, 119s. <https://doi.org/10.1097/00003446-198710001-00010>
- 553 Crosse, M.J., Butler, J.S., Lalor, E.C., 2015. Congruent Visual Speech Enhances Cortical
554 Entrainment to Continuous Auditory Speech in Noise-Free Conditions. *J. Neurosci.* 35,
555 14195–14204. <https://doi.org/10.1523/JNEUROSCI.1829-15.2015>
- 556 Crosse, M.J., Di Liberto, G.M., Bednar, A., Lalor, E.C., 2016. The Multivariate Temporal
557 Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to
558 Continuous Stimuli. *Front. Hum. Neurosci.* 10. <https://doi.org/10.3389/fnhum.2016.00604>
- 559 Cui, M.E., Herrmann, B., 2023. Eye Movements Decrease during Effortful Speech Listening. *J.*
560 *Neurosci.* 43, 5856–5869. <https://doi.org/10.1523/JNEUROSCI.0240-23.2023>

- 561 Dähne, S., Bießmann, F., Samek, W., Haufe, S., Goltz, D., Gundlach, C., Villringer, A., Fazli, S.,
562 Müller, K.-R., 2015. Multivariate Machine Learning Methods for Fusing Multimodal
563 Functional Neuroimaging Data. *Proc. IEEE* 103, 1507–1530.
564 <https://doi.org/10.1109/JPROC.2015.2425807>
- 565 Dai, L., Shinn-Cunningham, B.G., 2016. Contributions of Sensory Coding and Attentional Control
566 to Individual Differences in Performance in Spatial Auditory Selective Attention Tasks.
567 *Front. Hum. Neurosci.* 10. <https://doi.org/10.3389/fnhum.2016.00530>
- 568 Das, N., Bertrand, A., Francart, T., 2018. EEG-based auditory attention detection: boundary
569 conditions for background noise and speaker positions. *J. Neural Eng.* 15, 066017.
570 <https://doi.org/10.1088/1741-2552/aae0a6>
- 571 de Cheveigné, A., Wong, D.D.E., Di Liberto, G.M., Hjortkjær, J., Slaney, M., Lalor, E., 2018.
572 Decoding the auditory brain with canonical component analysis. *NeuroImage* 172, 206–
573 216. <https://doi.org/10.1016/j.neuroimage.2018.01.033>
- 574 Decruy, L., Lesenfants, D., Vanthornhout, J., Francart, T., 2020a. Top-down modulation of neural
575 envelope tracking: The interplay with behavioral, self-report and neural measures of
576 listening effort. *Eur. J. Neurosci.* 52, 3375–3393. <https://doi.org/10.1111/ejn.14753>
- 577 Decruy, L., Vanthornhout, J., Francart, T., 2020b. 🌟Hearing impairment is associated with
578 enhanced neural tracking of the speech envelope. *Hear. Res.* 393, 107961.
579 <https://doi.org/10.1016/j.heares.2020.107961>
- 580 Decruy, L., Vanthornhout, J., Francart, T., 2019. Evidence for enhanced neural tracking of the
581 speech envelope underlying age-related speech-in-noise difficulties. *J. Neurophysiol.* 122,
582 601–615. <https://doi.org/10.1152/jn.00687.2018>
- 583 Delorme, A., Makeig, S., 2004. EEGLAB: an open source toolbox for analysis of single-trial EEG
584 dynamics including independent component analysis. *J. Neurosci. Methods* 134, 9–21.
585 <https://doi.org/10.1016/j.jneumeth.2003.10.009>

- 586 Devocht, E.M.J., Janssen, A.M.L., Chalupper, J., Stokroos, R.J., George, E.L.J., 2017. The
587 Benefits of Bimodal Aiding on Extended Dimensions of Speech Perception: Intelligibility,
588 Listening Effort, and Sound Quality. *Trends Hear.* 21, 2331216517727900.
589 <https://doi.org/10.1177/2331216517727900>
- 590 Di Liberto, G.M., Hjortkjær, J., Mesgarani, N., 2022. Editorial: Neural Tracking: Closing the Gap
591 Between Neurophysiology and Translational Medicine. *Front. Neurosci.* 16.
592 <https://doi.org/10.3389/fnins.2022.872600>
- 593 Diaz, G., Cooper, J., Kit, D., Hayhoe, M., 2013. Real-time recording and classification of eye
594 movements in an immersive virtual environment. *J. Vis.* 13, 5.
595 <https://doi.org/10.1167/13.12.5>
- 596 Dimitrijevic, A., Smith, M.L., Kadis, D.S., Moore, D.R., 2019. Neural indices of listening effort in
597 noisy environments. *Sci. Rep.* 9, 11278. <https://doi.org/10.1038/s41598-019-47643-1>
- 598 Ding, N., Simon, J.Z., 2013. Adaptive Temporal Encoding Leads to a Background-Insensitive
599 Cortical Representation of Speech. *J. Neurosci.* 33, 5728–5735.
600 <https://doi.org/10.1523/JNEUROSCI.5297-12.2013>
- 601 Ding, N., Simon, J.Z., 2012a. Emergence of neural encoding of auditory objects while listening to
602 competing speakers. *Proc. Natl. Acad. Sci.* 109, 11854–11859.
603 <https://doi.org/10.1073/pnas.1205381109>
- 604 Ding, N., Simon, J.Z., 2012b. Neural coding of continuous speech in auditory cortex during
605 monaural and dichotic listening. *J. Neurophysiol.* 107, 78–89.
606 <https://doi.org/10.1152/jn.00297.2011>
- 607 Etard, O., Reichenbach, T., 2019. Neural Speech Tracking in the Theta and in the Delta
608 Frequency Band Differentially Encode Clarity and Comprehension of Speech in Noise. *J.*
609 *Neurosci.* 39, 5750–5759. <https://doi.org/10.1523/JNEUROSCI.1828-18.2019>

- 610 Feng, Y., Chen, F., 2022. Nonintrusive objective measurement of speech intelligibility: A review
611 of methodology. *Biomed. Signal Process. Control* 71, 103204.
612 <https://doi.org/10.1016/j.bspc.2021.103204>
- 613 Fiedler, L., Wöstmann, M., Herbst, S.K., Obleser, J., 2019. Late cortical tracking of ignored
614 speech facilitates neural selectivity in acoustically challenging conditions. *NeuroImage*
615 186, 33–42. <https://doi.org/10.1016/j.neuroimage.2018.10.057>
- 616 Fritz, J.B., Elhilali, M., David, S.V., Shamma, S.A., 2007. Auditory attention — focusing the
617 searchlight on sound. *Curr. Opin. Neurobiol., Sensory systems* 17, 437–455.
618 <https://doi.org/10.1016/j.conb.2007.07.011>
- 619 Fuglsang, S.A., Märcher-Rørsted, J., Dau, T., Hjortkjær, J., 2020. Effects of Sensorineural
620 Hearing Loss on Cortical Synchronization to Competing Speech during Selective Attention.
621 *J. Neurosci.* 40, 2562–2572. <https://doi.org/10.1523/JNEUROSCI.1936-19.2020>
- 622 Gopher, D., 1973. Eye-movement patterns in selective listening tasks of focused attention.
623 *Percept. Psychophys.* 14, 259–264. <https://doi.org/10.3758/BF03212387>
- 624 Hauswald, A., Keitel, A., Chen, Y.-P., Rösch, S., Weisz, N., 2022. Degradation levels of
625 continuous speech affect neural speech tracking and alpha power differently. *Eur. J.*
626 *Neurosci.* 55, 3288–3302. <https://doi.org/10.1111/ejn.14912>
- 627 Herrmann, B., Ryan, J.D., 2024. Pupil Size and Eye Movements Differently Index Effort in Both
628 Younger and Older Adults. *J. Cogn. Neurosci.* 36, 1325–1340.
629 https://doi.org/10.1162/jocn_a_02172
- 630 Horton, C., Srinivasan, R., D’Zmura, M., 2014. Envelope responses in single-trial EEG indicate
631 attended speaker in a ‘cocktail party.’ *J. Neural Eng.* 11, 046015.
632 <https://doi.org/10.1088/1741-2560/11/4/046015>
- 633 Horwitz-Martin, R.L., Quatieri, T.F., Godoy, E., Williamson, J.R., 2016. A vocal modulation model
634 with application to predicting depression severity, in: 2016 IEEE 13th International
635 Conference on Wearable and Implantable Body Sensor Networks (BSN). Presented at the

- 636 2016 IEEE 13th International Conference on Wearable and Implantable Body Sensor
637 Networks (BSN), pp. 247–253. <https://doi.org/10.1109/BSN.2016.7516268>
- 638 Iotzov, I., Parra, L.C., 2019. EEG can predict speech intelligibility. *J. Neural Eng.* 16, 036008.
639 <https://doi.org/10.1088/1741-2552/ab07fe>
- 640 Kang, S.S., Lano, T.J., Sponheim, S.R., 2015. Distortions in EEG interregional phase synchrony
641 by spherical spline interpolation: causes and remedies. *Neuropsychiatr. Electrophysiol.* 1,
642 9. <https://doi.org/10.1186/s40810-015-0009-5>
- 643 Kell, A.J.E., McDermott, J.H., 2019. Invariance to background noise as a signature of non-primary
644 auditory cortex. *Nat. Commun.* 10, 3958. <https://doi.org/10.1038/s41467-019-11710-y>
- 645 Kirchner, W.K., 1958. Age differences in short-term retention of rapidly changing information. *J.*
646 *Exp. Psychol.* 55, 352–358. <https://doi.org/10.1037/h0043688>
- 647 Kong, Y.-Y., Mullangi, A., Ding, N., 2014. Differential modulation of auditory responses to
648 attended and unattended speech in different listening conditions. *Hear. Res.* 316, 73–81.
649 <https://doi.org/10.1016/j.heares.2014.07.009>
- 650 Krueger, M., Schulte, M., Zokoll, M.A., Wagener, K.C., Meis, M., Brand, T., Holube, I., 2017.
651 Relation Between Listening Effort and Speech Intelligibility in Noise. *Am. J. Audiol. Online*
652 26, 378–392. https://doi.org/10.1044/2017_AJA-16-0136
- 653 Laffere, A., Dick, F., Tierney, A., 2020. Effects of auditory selective attention on neural phase:
654 individual differences and short-term training. *NeuroImage* 213, 116717.
655 <https://doi.org/10.1016/j.neuroimage.2020.116717>
- 656 Lalor, E.C., Power, A.J., Reilly, R.B., Foxe, J.J., 2009. Resolving Precise Temporal Processing
657 Properties of the Auditory System Using Continuous Stimuli. *J. Neurophysiol.* 102, 349–
658 359. <https://doi.org/10.1152/jn.90896.2008>
- 659 Lesenfants, D., Vanthornhout, J., Verschueren, E., Decruy, L., Francart, T., 2019. Predicting
660 individual speech intelligibility from the cortical tracking of acoustic- and phonetic-level

- 661 speech representations. *Hear. Res.* 380, 1–9.
662 <https://doi.org/10.1016/j.heares.2019.05.006>
- 663 Marinato, G., Baldauf, D., 2019. Object-based attention in complex, naturalistic auditory streams.
664 *Sci. Rep.* 9, 2854. <https://doi.org/10.1038/s41598-019-39166-6>
- 665 Mesgarani, N., Chang, E.F., 2012. Selective cortical representation of attended speaker in multi-
666 talker speech perception. *Nature* 485, 233–236. <https://doi.org/10.1038/nature11020>
- 667 Mesgarani, N., David, S.V., Fritz, J.B., Shamma, S.A., 2014. Mechanisms of noise robust
668 representation of speech in primary auditory cortex. *Proc. Natl. Acad. Sci.* 111, 6792–
669 6797. <https://doi.org/10.1073/pnas.1318017111>
- 670 Mesgarani, N., David, S.V., Fritz, J.B., Shamma, S.A., 2009. Influence of Context and Behavior
671 on Stimulus Reconstruction From Neural Activity in Primary Auditory Cortex. *J.*
672 *Neurophysiol.* 102, 3329–3339. <https://doi.org/10.1152/jn.91128.2008>
- 673 Mesik, J., Ray, L., Wojtczak, M., 2021. Effects of Age on Cortical Tracking of Word-Level Features
674 of Continuous Competing Speech. *Front. Neurosci.* 15.
675 <https://doi.org/10.3389/fnins.2021.635126>
- 676 Muncke, J., Kuruvila, I., Hoppe, U., 2022. Prediction of Speech Intelligibility by Means of EEG
677 Responses to Sentences in Noise. *Front. Neurosci.* 16.
678 <https://doi.org/10.3389/fnins.2022.876421>
- 679 Nilsson, M., Soli, S.D., Sullivan, J.A., 1994. Development of the Hearing in Noise Test for the
680 measurement of speech reception thresholds in quiet and in noise. *J. Acoust. Soc. Am.*
681 95, 1085–1099. <https://doi.org/10.1121/1.408469>
- 682 Nogueira, W., Dolhopiatenko, H., 2022. Predicting speech intelligibility from a selective attention
683 decoding paradigm in cochlear implant users. *J. Neural Eng.* 19, 026037.
684 <https://doi.org/10.1088/1741-2552/ac599f>

- 685 Oord, A. van den, Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner,
686 N., Senior, A., Kavukcuoglu, K., 2016. WaveNet: A Generative Model for Raw Audio.
687 <https://doi.org/10.48550/arXiv.1609.03499>
- 688 O’Sullivan, J.A., Power, A.J., Mesgarani, N., Rajaram, S., Foxe, J.J., Shinn-Cunningham, B.G.,
689 Slaney, M., Shamma, S.A., Lalor, E.C., 2015. Attentional Selection in a Cocktail Party
690 Environment Can Be Decoded from Single-Trial EEG. *Cereb. Cortex* N. Y. N 1991 25,
691 1697–1706. <https://doi.org/10.1093/cercor/bht355>
- 692 Pashler, H., Johnston, J.C., Ruthruff, E., 2001. Attention and Performance. *Annu. Rev. Psychol.*
693 52, 629–651. <https://doi.org/10.1146/annurev.psych.52.1.629>
- 694 Perrin, F., Pernier, J., Bertrand, O., Echallier, J.F., 1989. Spherical splines for scalp potential and
695 current density mapping. *Electroencephalogr. Clin. Neurophysiol.* 72, 184–187.
696 [https://doi.org/10.1016/0013-4694\(89\)90180-6](https://doi.org/10.1016/0013-4694(89)90180-6)
- 697 Power, A.J., Foxe, J.J., Forde, E.-J., Reilly, R.B., Lalor, E.C., 2012. At what time is the cocktail
698 party? A late locus of selective attention to natural speech. *Eur. J. Neurosci.* 35, 1497–
699 1503. <https://doi.org/10.1111/j.1460-9568.2012.08060.x>
- 700 Rabinowitz, N.C., Willmore, B.D.B., King, A.J., Schnupp, J.W.H., 2013. Constructing Noise-
701 Invariant Representations of Sound in the Auditory Pathway. *PLOS Biol.* 11, e1001710.
702 <https://doi.org/10.1371/journal.pbio.1001710>
- 703 Raghavan, V.S., O’Sullivan, J., Bickel, S., Mehta, A.D., Mesgarani, N., 2023. Distinct neural
704 encoding of glimpsed and masked speech in multitalker situations. *PLOS Biol.* 21,
705 e3002128. <https://doi.org/10.1371/journal.pbio.3002128>
- 706 Sarter, M., Gehring, W.J., Kozak, R., 2006. More attention must be paid: The neurobiology of
707 attentional effort. *Brain Res. Rev.* 51, 145–160.
708 <https://doi.org/10.1016/j.brainresrev.2005.11.002>

- 709 Sharma, D., Naylor, P.A., Brookes, M., 2013. Non-intrusive speech intelligibility assessment, in:
710 21st European Signal Processing Conference (EUSIPCO 2013). Presented at the 21st
711 European Signal Processing Conference (EUSIPCO 2013), pp. 1–5.
- 712 Spahr, A.J., Dorman, M.F., Litvak, L.M., Van Wie, S., Gifford, R.H., Loizou, P.C., Loiseau, L.M.,
713 Oakes, T., Cook, S., 2012. Development and validation of the AzBio sentence lists. *Ear*
714 *Hear.* 33, 112–117. <https://doi.org/10.1097/AUD.0b013e31822c2549>
- 715 Strauss, D.J., Francis, A.L., 2017. Toward a taxonomic model of attention in effortful listening.
716 *Cogn. Affect. Behav. Neurosci.* 17, 809–825. <https://doi.org/10.3758/s13415-017-0513-0>
- 717 Vanthornhout, J., Decruy, L., Francart, T., 2019a. Effect of Task and Attention on Neural Tracking
718 of Speech. *Front. Neurosci.* 13.
- 719 Vanthornhout, J., Decruy, L., Francart, T., 2019b. Effect of Task and Attention on Neural Tracking
720 of Speech. *Front. Neurosci.* 13. <https://doi.org/10.3389/fnins.2019.00977>
- 721 Vanthornhout, J., Decruy, L., Wouters, J., Simon, J.Z., Francart, T., 2018. Speech Intelligibility
722 Predicted from Neural Entrainment of the Speech Envelope. *J. Assoc. Res. Otolaryngol.*
723 *JARO* 19, 181–191. <https://doi.org/10.1007/s10162-018-0654-z>
- 724 Verschueren, E., Vanthornhout, J., Francart, T., 2020. The Effect of Stimulus Choice on an EEG-
725 Based Objective Measure of Speech Intelligibility. *Ear Hear.* 41, 1586.
726 <https://doi.org/10.1097/AUD.0000000000000875>
- 727 Wichmann, F.A., Hill, N.J., 2001a. The psychometric function: I. Fitting, sampling, and goodness
728 of fit. *Percept. Psychophys.* 63, 1293–1313. <https://doi.org/10.3758/BF03194544>
- 729 Wichmann, F.A., Hill, N.J., 2001b. The psychometric function: II. Bootstrap-based confidence
730 intervals and sampling. *Percept. Psychophys.* 63, 1314–1329.
731 <https://doi.org/10.3758/BF03194545>
- 732 Yasmin, S., Irsik, V.C., Johnsrude, I.S., Herrmann, B., 2023. The effects of speech masking on
733 neural tracking of acoustic and semantic features of natural speech. *Neuropsychologia*
734 186, 108584. <https://doi.org/10.1016/j.neuropsychologia.2023.108584>

735 Zekveld, A.A., Heslenfeld, D.J., Festen, J.M., Schoonhoven, R., 2006. Top-down and bottom-up
736 processes in speech comprehension. *NeuroImage* 32, 1826–1836.
737 <https://doi.org/10.1016/j.neuroimage.2006.04.199>

738 Zion Golumbic, E.M., Ding, N., Bickel, S., Lakatos, P., Schevon, C.A., McKhann, G.M., Goodman,
739 R.R., Emerson, R., Mehta, A.D., Simon, J.Z., Poeppel, D., Schroeder, C.E., 2013.
740 Mechanisms Underlying Selective Neuronal Tracking of Attended Speech at a ‘Cocktail
741 Party.’ *Neuron* 77, 980–991. <https://doi.org/10.1016/j.neuron.2012.12.037>

742