

# 1 An integrative view on attentional 2 modulation in naturalistic speech

3 Jiawei Li<sup>1,3§</sup>, Bo Hong<sup>2,3</sup>, Guido Nolte<sup>4</sup>, Andreas K. Engel<sup>4</sup>, Dan Zhang<sup>1,3\*</sup>

\*For correspondence:

[dzhang@tsinghua.edu.cn](mailto:dzhang@tsinghua.edu.cn) (FMS)

Present address: <sup>§</sup>Department of  
Education and Psychology, Freie  
Universität Berlin, Berlin, Germany

4 <sup>1</sup>Department of Psychology, School of Social Sciences, Tsinghua University, Beijing,  
5 China; <sup>2</sup>Department of Biomedical Engineering, School of Medicine, Tsinghua  
6 University, Beijing, China; <sup>3</sup>Tsinghua Laboratory of Brain and Intelligence, Tsinghua  
7 University, Beijing, China; <sup>4</sup>Department of Neurophysiology and Pathophysiology,  
8 University Medical Center Hamburg Eppendorf, Hamburg, Germany

---

10 **Abstract** Attending to a speaker is a complex process: to hear sound waves that represent  
11 *acoustic features*; to understand the meaning of words that represent *semantic features*; and the  
12 listener and speaker need to be aligned to form a common ground, which represents *inter-brain*  
13 *features*. Little is known about how attention modulates these features from the speaker in an  
14 integrative way. Adopting naturalistic speech, combing with natural language processing models  
15 and inter-brain EEG analysis methods, we measured how listener responses to different  
16 information from the attended speaker simultaneously. Our result reveals that: the sound is the  
17 first to be processed; the meaning of the attended speech is parsed after that. The listener's mind  
18 aligned to the speaker even seconds before the speech begins. Together, our results illustrated  
19 how our brain is selectively entrained to different types of information from the speaker in an  
20 integrative view.

---

## 22 Introduction

23 *"Men do not understand one another by actually exchanging signs for things... they do it by striking the*  
24 *same note on their mental instruments."*

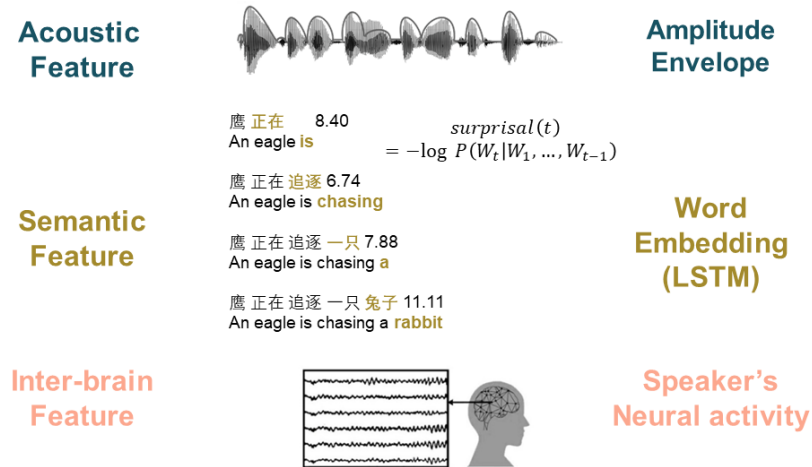
25 *Wilhelm von Humboldt*

26 We don't always hear toasts in the daily life. There are many speakers in a "cocktail party" sit-  
27 uation. It is not a difficult task for most people to follow the person that they pay attention to  
28 and ignore others. However, the neural process behind this "cocktail party" situation (*Cherry, 1953*;  
29 *McCarthy and Nobre, 1993*; *Middlebrooks et al., 2017*) is complex: the listener needs to trans-  
30 form the attended sound wave to the meaning (*Hasson et al., 2012*; *Heilbron et al., 2022*). More-  
31 over, the listener needs to actively perceive the speech and form a common grounding with the  
32 speaker (*Friston, 2009*; *Jiang et al., 2021, 2012*; *Pulvermüller and Fadiga, 2010*; *Stolk et al., 2016*;  
33 *Yeshurun et al., 2021*). We still do not know when we pay attention to someone how our brain in-  
34 tegrates different types of information from the speaker. Previous studies mainly focused on the  
35 modulation mechanism of the speech itself. The auditory scene analysis studies mainly focused  
36 on acoustic features of the speech information in the cocktail party problem (*Bregman, 1990*; *Brod-*  
37 *beck and Simon, 2022*; *Ding and Simon, 2012b,a*; *Shamma et al., 2011*; *Shinn-Cunningham, 2008*;  
38 *Teoh et al., 2022*; *Wang et al., 2019*). These studies demonstrated that the attentional modulation  
39 of processing of acoustic features mainly occurs 100-250 ms after the speech onset on the low-  
40 frequency bands (e.g., 2-8 Hz) (*Broderick et al., 2021*; *Mesik et al., 2021*; *Weissbart et al., 2019*).  
41 Recent studies revealed that processing of semantic or linguistic features could also be modulated

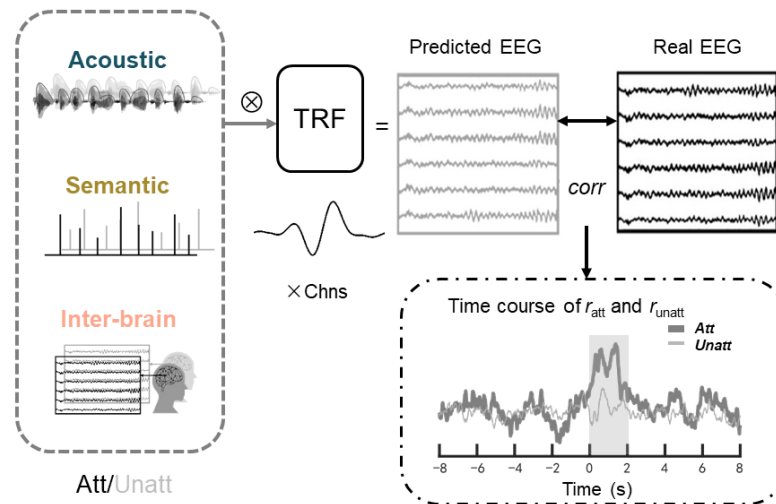
42 by attention, for the features in the attended stream could also better understood than the unat-  
43 tended stream(Broderick *et al.*, 2018; Connolly *et al.*, 1990; Dai *et al.*, 2022; Har-shai Yahav and  
44 Zion Golumbic, 2021; Heil *et al.*, 2004). However, when we pay attention to someone, the speaker  
45 itself is also important. As famous German philosopher Gadamer (1975) once said, “there can be  
46 no speech that does not bind the speaker and the person spoken to.” The inter-brain studies pro-  
47 vided another angle to the traditional attention studies and included the speaker in the scene: the  
48 listeners’ neural activity could entrain not only the sound wave but also coupled with the speaker’s  
49 neural activity (Pérez *et al.*, 2017; Stephens *et al.*, 2010; Yeshurun *et al.*, 2021). One study reported  
50 that interpersonal neural synchronization (INS) between the listener and the attended speaker was  
51 selectively enhanced using fNIRS recording(Dai *et al.*, 2018). It indicates the listener could actively  
52 “understand” the speaker not only through the speech itself but rely on the “beyond the stimu-  
53 lus” grounding (Bashivan *et al.*, 2019; Hartley and Poeppel, 2020; Hasson *et al.*, 2012; Jiang *et al.*,  
54 2021; Redcay and Schilbach, 2019; Stolk *et al.*, 2016). While attentional modulation of different  
55 features was explored separately, to our knowledge, there has not been a single study to explore  
56 the attentional modulation of multiple types of information from the speaker at the same time.

57 An integrative view enables us to investigate the concurrent attentional modulation of differ-  
58 ent features, and it may answer some critical questions about attentional modulation. The first  
59 question concerns the temporal dynamics of attentional modulation of different features. Does  
60 the attention modulate all the features at the same time? Or does the attention modulate them  
61 in sequential order? It was hard for the previous studies to answer because they mostly focused  
62 on one or two types of speech features(Broderick *et al.*, 2021, 2019; Mesik *et al.*, 2021), and they  
63 didn’t consider the speaker. The time range of attentional modulations varied in the previous stud-  
64 ies. Most single brain studies chose the time range within 1 second after speech onset(Ding and  
65 Simon, 2012b,a; Lalor and Foxe, 2010; Power *et al.*, 2012; Teoh *et al.*, 2022; Zion Golumbic *et al.*,  
66 2013). However, the inter-brain studies revealed that the listener could selectively tune in to the  
67 attended speaker up to 5 seconds before the speech onset(Dai *et al.*, 2018). This phenomenon  
68 was rarely observed in the high temporal resolution EEG studies(Kuhlen *et al.*, 2012), due to the  
69 different selection of the time windows. We still do not know which speech processing stages at-  
70 tention could modulate. The second question centers on the relationship between different types  
71 of attentional modulations. Previous studies revealed that the entrainment to acoustic features  
72 and semantic features could interact with each other during a monologue condition: The semantic  
73 feature was reported to enhance the entrainment to the acoustic feature(Anderson *et al.*, 2019;  
74 Gillis *et al.*, 2021; Heilbron *et al.*, 2022). There were few studies providing evidence about how  
75 these three features interact with each other in a “cocktail party situation”(Dai *et al.*, 2018; Pérez  
76 *et al.*, 2017, 2019). The correlation between acoustic features, semantic features, and inter-brain  
77 features remains elusive. The present study aimed to reveal the neural mechanism of attentional  
78 modulation in an integrative view, implying that the attentional modulation of the three types of  
79 information from the speaker would be explored simultaneously. Naturalistic speech was used  
80 as stimulus material, which contains much richer information than either the sound sequence or  
81 single word stimulation employed in previous studies(Broderick *et al.*, 2018; Hamilton *et al.*, 2018;  
82 Hartley and Poeppel, 2020; Nastase *et al.*, 2021; Sonkusare *et al.*, 2019; Willems *et al.*, 2020). The  
83 sound, the meaning, and the speaker could appear at the same time in an ecological situation and  
84 enables us to investigate different types of information from the speaker simultaneously(Hasson  
85 *et al.*, 2012; Willems *et al.*, 2020). As Figure 1(A) illustrates, the amplitude envelope was the rep-  
86 resentation of the acoustic feature. A natural language processing model (an LSTM model) was  
87 applied to extract the semantic feature in the text. We chose surprisal as the representation of the  
88 semantic feature in the previous studies, which is one of the most important semantic features that  
89 have received sufficient investigation in previous neuroscience studies(Brodbeck and Simon, 2022;  
90 Frank and Willems, 2017; Willems and Jacobs, 2016). The early event-related potentials studies also  
91 used the congruent or incongruent words to explore the neural mechanism of semantic process-  
92 ing, which is an early version of surprisal(Kutas and Hillyard, 1984; Lau *et al.*, 2008). A sequential

A.



B.



**Figure 1.** The stimuli, the experimental paradigm and the analysis process. (A) The demonstration of acoustic, semantic, and inter-brain features. The envelope of the speech represents the acoustic feature. The surprisal index, which was calculated by an LSTM model, was used as the semantic feature. The speaker's neural activity is the inter-brain feature. (B) A "cocktail party" selective attention paradigm was used, in which the listener was asked to pay attention to one side of the speech stream and ignore the other. While listening to the speech stream, the listener's EEG signals were recorded. The Encoding  $r$  for the attended feature and the unattended feature was calculated time point by time point by applying the TRF method.

93 dual-brain approach was used, and the electroencephalogram (EEG) of both speaker and listener  
94 was recorded. The speaker's neural activity represented the inter-brain feature (*Hasson et al., 2012*;  
95 *Jiang et al., 2021*; *Leong et al., 2017*; *Pérez et al., 2017, 2019*; *Stolk et al., 2016*). A temporal response  
96 function (TRF) method was used to measure the difference between the entrainment to attended  
97 and unattended features. Based on previous studies, we hypothesized that attention modulates  
98 the different features in distinct ways, which means the modulation to different features would  
99 happen on different frequency bands and time ranges. For the frequency bands, we hypothesize  
100 that the delta band reflects attentional modulation of the semantic feature (*Dai et al., 2022*; *Teoh*  
101 *et al., 2022*; *Yu et al., 2022*), and the theta band represents the acoustic feature (*Ding et al., 2014*;  
102 *Etard et al., 2019*). The entrainment to the speaker's neural feature needs further exploration. For  
103 the time course, we assume that the attentional modulation of the acoustic feature occurs at first  
104 (*Ding and Simon, 2012b,a*; *Hillyard et al., 1973*; *O'Sullivan et al., 2015*; *Power et al., 2011*; *Teoh*  
105 *et al., 2022*), and the attention effect of the semantic feature lasts longer (*Broderick et al., 2018*;  
106 *Dai et al., 2022*).

107 The listener's neural activity is aligned with the speaker's neural activity in a broader time range (*Dai*  
108 *et al., 2018*; *Kuhlen et al., 2012*). We also hypothesize that the attentional modulation of processing  
109 of the acoustic feature and the semantic feature correlates to each other (*Heilbron et al., 2022*),  
110 but the inter-brain feature has a distinct pattern, which is independent of the acoustic and seman-  
111 tic feature. Together, our study adopted an integrative view to investigate three different types of  
112 features when attending to a speaker, which would further our understanding of the attentional  
113 modulation from the speech to the speaker.

## 114 Results

### 115 Behavioral Performance of the Listeners

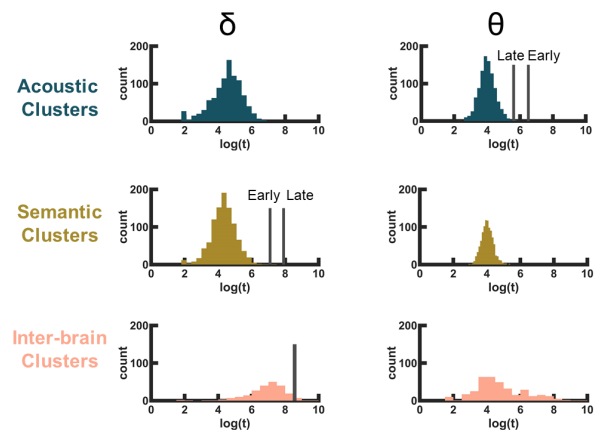
116 The average comprehension performance was significantly better for the 28 attended stories than  
117 for the 28 unattended stories ( $67.0 \pm 2.5\%$  (standard error) vs.  $36.0 \pm 1.6\%$ ,  $t(19) = 10.948$ ,  $p < .001$ ;  
118 the four-choice chance level: 25%). The participants reported a moderate level of attention ( $8.146$   
119  $\pm 0.343$  on a 10-point Likert scale) and attention difficulties ( $2.039 \pm 0.530$  on a 10-point Likert scale).  
120 The accuracy for the attended story was significantly correlated with both the self-reported atten-  
121 tion level ( $r(18) = .476$ ,  $p = .043$ ) and attention difficulty ( $r(18) = -.677$ ,  $p = .001$ ). The self-reported  
122 story familiarity level was low for all the participants ( $0.860 \pm 0.220$  on a 10-point Likert scale) and  
123 was not correlated with comprehension performance ( $r(18) = -.224$ ,  $p = .342$ ). These results sug-  
124 gest that participants' selective attention was effectively manipulated, and the measurement of  
125 comprehension performance was reliable. The response accuracy varied from 25.0% to 51.8% for  
126 unattended stories.

### 127 The theta band and the delta band reflect distinct attentional to different features

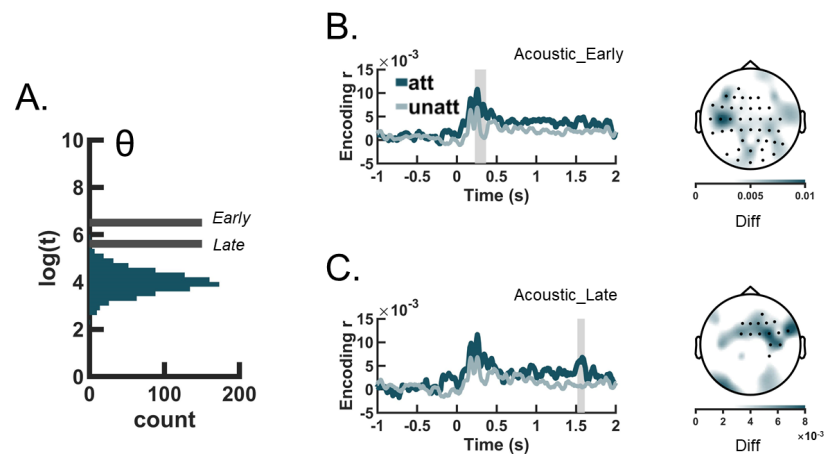
128 The delta and theta bands have different functional roles in attentional modulation, as Figure 2  
129 illustrates. A cluster-based permutation (*Maris et al., 2007*) was conducted to reveal the difference  
130 between the encoding of the attended and the unattended features and control for multiple com-  
131 parisons. The theta band only modulates the processing of the acoustic feature. There are two  
132 significant acoustic clusters that appear in the theta band. They were designated as Acoustic-Early  
133 (cluster-based  $p < .001$ ) and Acoustic-Late (cluster-based  $p = .005$ ) depending on the time they oc-  
134 curred. In contrast to the acoustic feature, semantic clusters were found in the delta band. Two  
135 clusters illustrated the difference in entrainment to the attended and unattended semantic fea-  
136 tures. They were labeled Semantic-early (cluster-based  $p = .002$ ) and Semantic-late (cluster-based  
137  $p < .001$ ). One inter-brain cluster was also found in the delta band (cluster-based  $p < .001$ ).

### 138 The attentional modulation of different features unfolds in a distinct time range

139 As shown in Figure 3, the Acoustic-Early cluster involved the left-lateralized fronto-central and oc-  
140 cipital electrodes (cluster-based permutation  $p < .001$ ) at a latency of 0.219-0.359 s after the onset



**Figure 2.** Different attentional modulation roles of different bands. The null distribution of the  $t$ -statistics of every feature in the delta band(left) and theta band(right). The grey lines indicate significant clusters.



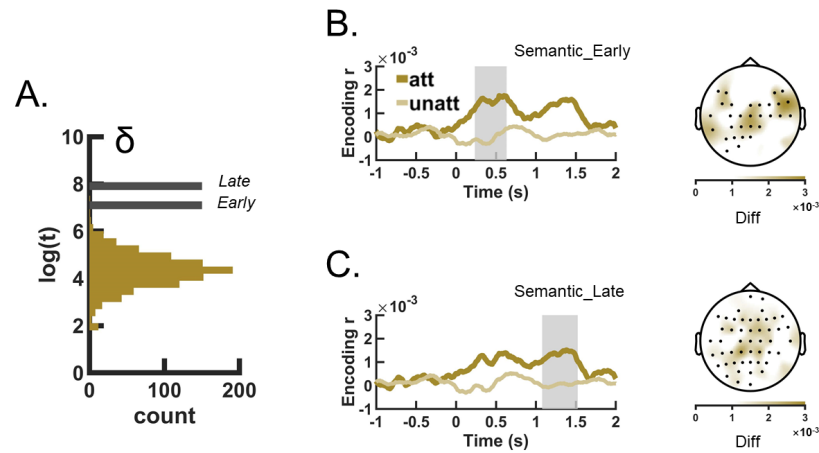
**Figure 3.** The attentional modulation of the acoustic feature. (A) The null distribution of the cluster-based  $t$ -statistics of the acoustic feature. (B) and (C) The time course and the topo-plot of the significant acoustic cluster. The dark blue line represents the Encoding  $r_{att}$ , and the light line represents the Encoding  $r_{unatt}$ . The shaded region depicts a significant difference in the time window. The topo-plot of the average difference between in Encoding  $r_{att}$  and Encoding  $r_{unatt}$  cluster. The black dots indicate the channels in the cluster.

141 of speech. The Acoustic-Late cluster had a later latency of 1.508-1.602 s with the electrodes in the  
 142 right-frontal regions (cluster-based  $p = .005$ ).

143 As Figure 4 indicates, the Semantic-Early cluster occurred at 0.227-0.621 s covering the elec-  
 144 trodes in frontal and central regions. The Semantic-Late cluster was found at 1.073-1.516 s involv-  
 145 ing the wide distribution of the electrodes. There was only one cluster inter-brain cluster. Unlike  
 146 the acoustic clusters and semantic clusters, the inter-brain cluster had a wide time range of -4.836  
 147 to -0.539 s with the electrodes in the left frontal region, as shown in Figure 5.

### 148 **The entrainment to the inter-brain feature is independent of the acoustic and se-** 149 **semantic features**

150 As Figure 6 indicates, the average Encoding  $r_{att}$  in Semantic-Early cluster and the average Encoding  
 151  $r_{att}$  in Acoustic-Early clusters were highly correlated ( $r(18) = .786, p < .001$ , FDR-corrected). The  
 152 average Encoding  $r_{att}$  in Semantic-late cluster and the average Encoding  $r_{att}$  in Acoustic-Early  
 153 cluster were also highly correlated ( $r(18) = .565, p = .045$ , FDR-corrected). There were no other  
 154 significant correlations between other clusters ( $ps > .05$ ).



**Figure 4.** The attentional modulation of the semantic feature. (A) The null distribution of the cluster-based  $t$ -statistics of the semantic feature. (B) and (C) The time course and the topo-plot of the significant semantic cluster. The dark brown line represents the Encoding  $r_{att}$ , and the light brown represents the Encoding  $r_{unatt}$ . The shaded region depicts a significant difference in the time window. The topo-plot of the average difference between in Encoding  $r_{att}$  and Encoding  $r_{unatt}$  cluster. The black dots indicate the channels in the cluster.

155 **The entrainment to the inter-brain feature was correlated to the comprehension**  
 156 **performance**

157 The partial correlation between the behavioral performance and the coefficients in was calculated  
 158 to reveal the unique contribution of a certain feature to the comprehension performance. Specifi-  
 159 cally, the correlations between the mean Encoding  $r_{att}$  in the cluster and the accuracy of the ques-  
 160 tions were calculated while controlling other features. As Table 1 illustrated, only the inter-brain  
 161 cluster has a significant partial correlation with the behavioral performance, partial correlation  $r$   
 162 (18) = -.769,  $p = .002$  (FDR-corrected). All the other clusters didn't reveal significant correlations. We  
 163 further analyzed the Encoding  $r_{att}$  in the inter-brain cluster and difficulty. We found significant  
 164 positive correlation,  $r(18) = .499$ ,  $p = .025$  as shown in Figure 5(c).

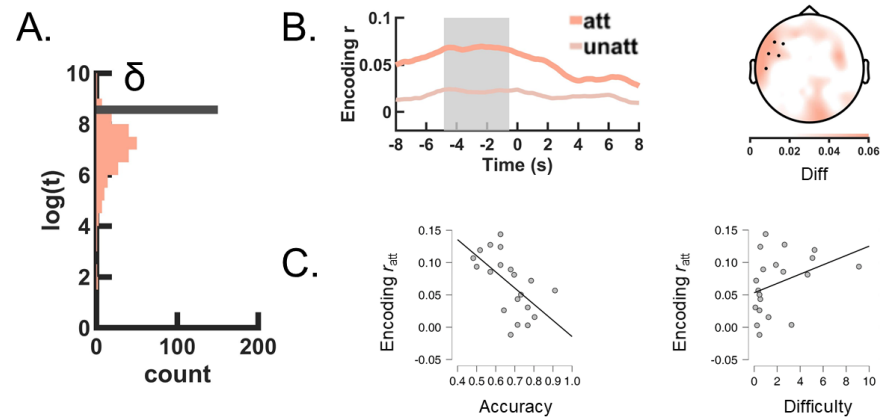
**Table 1.** The partial correlation between the entrainment coefficients and behavioral performance.

Clusters	$\rho$	$p$ -values(FDR corrected)
Acoustic-early	0.512	.089
Acoustic-late	-0.306	.291
Semantic-early	-0.412	.170
Semantic-late	0.098	.710
Inter-brain	-0.769**	.002

\*\*  $p < .01$  (FDR corrected)

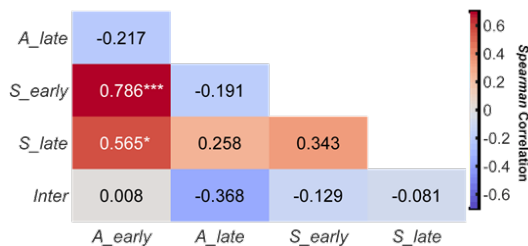
165 **Discussion**

166 Our study provided an integrative view of how our brain allocates attention to the different types of  
 167 information from the speaker in the “cocktail party problem”, as Figure 7 illustrates: the attention  
 168 modulated the sound firstly in the theta band at 200-350 ms. The meaning of speech was modu-  
 169 lated later, in the delta band at 200-600 ms. The listeners aligned to the speaker's neural activity 5  
 170 s before the speech onset. The entrainment to the acoustic feature and semantic feature were cor-  
 171 related, but the entrainment to the speaker's neural activity is independent of the speech stimuli.  
 172 Noticeably, only the entrainment to the speaker's neural activity has a negative correlation with the

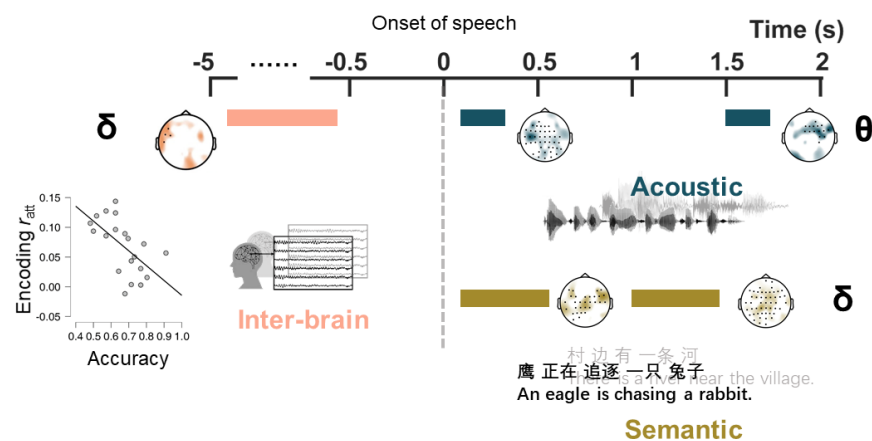


**Figure 5.** The attentional modulation of the inter-brain feature. (A) The null distribution of the cluster-based  $t$ -statistics of the inter-brain feature. (B) The time course and the topo-plot of the significant inter-brain cluster. The light pink line represents the Encoding  $r_{unatt}$ , and the dark pink represents the Encoding  $r_{att}$ . The shaded region depicts a significant difference in the time window. The topo-plot of the average difference between in Encoding  $r_{att}$  and Encoding  $r_{unatt}$  cluster. The black dots indicate the channels in the cluster. (C) The correlation between the Encoding  $r_{att}$  and the behavioral results.

173 comprehension score, which indicates the compensation role of the entrainment to the speaker's  
 174 neural activity. Our study clearly illustrated the temporal dynamics of attentional modulation of  
 175 the different features in an integrative view. The *acoustic feature* was modulated occurred at first,  
 176 with a latency of 200-350 ms. The time range was in line with the classical attention studies applying  
 177 naturalistic speech as the stimulus (*Ding and Simon, 2012b; Lalor and Foxe, 2010; O'Sullivan et al.,*  
 178 *2015; Wang et al., 2012*). The attention effect of the semantic feature lasted longer, from 200-600  
 179 ms. The N400 effect happened around 400 ms after the speech onset, and it was the most crucial  
 180 neural signature of the semantic process (*Kutas and Hillyard, 1989, 1984; Lau et al., 2008*). Our  
 181 result demonstrated that the attentional modulation of the *semantic feature* lasted longer than  
 182 the acoustic feature, and the time range was consistent with the classical N400 (*Brodbeck et al.,*  
 183 *2018*). We also found two late clusters for the acoustic feature and the semantic feature, which  
 184 have rarely been reported in previous studies. It may indicate the attentional modulation on the  
 185 sentence level (e.g., the terminal of the sentence), as the early ERP studies suggested (*Connolly*  
 186 *et al., 1990; Connolly and Phillips, 1994; Sanders and Neville, 2003*). The time range of the *inter-*  
 187 *brain feature* revealed a different pattern. The listener aligned to the speaker's neural activity 5  
 188 s before the speech onset, which was much earlier than the other two features. This result repli-  
 189 cates previous results using fMRI or fNIRS on high temporal resolution EEG signals (*Dai et al., 2018;*  
 190 *Jiang et al., 2012; Liu et al., 2020b; Stephens et al., 2010*), indicating that attentional modulation  
 191 may occur much earlier than we expect. Our result further suggests that the entrainment to the



**Figure 6.** The correlations between Encoding  $r_{att}$  in different clusters. A-early stands for Acoustic-Early, A-late stands for Acoustic-Late, S-early is short for Semantic-early, S-late stands for Semantic-Late, and Inter is short for Inter-brain.



**Figure 7.** An integrative view of the attentional modulations of different features in speech.

192 speaker's neural activity was different from other entrainments, which emphasizes the importance  
 193 of the speaker (Dai et al., 2018; Pérez et al., 2017, 2019).

194 Our study reveals the distinct roles of delta and theta bands in attentional modulation. The  
 195 theta band modulates with the acoustic feature, and the delta band modulates with the semantic  
 196 feature and inter-brain feature. This result possibly extends our understanding of the functional  
 197 roles of the two bands. It is consistent with the previous finding that the theta band processes  
 198 stimulus-linked features, like the acoustic feature (Ding et al., 2014; Etard et al., 2019; Li et al., 2022;  
 199 Lu et al., 2022). In particular, the modulation of the inter-brain feature was also found in the delta  
 200 band, which may indicate that the delta bands reflects comprehension-related functions, which  
 201 was rarely reported in the previous studies. Our study extends to the attended target from the  
 202 speech to the speaker by adopting the inter-brain feature. The entrainment to the speaker's neural  
 203 activity was the only predictor of comprehension performance, and it didn't correlate with other  
 204 features. Noticeably, the entrainment to the speaker's neural activity was a negative correlation  
 205 with the comprehension performance. We further found that the encoding index was positively  
 206 correlated with the perceived task difficulty reported by the listeners. We called that a "compensa-  
 207 tion" mechanism: when the listeners find it hard to complete the task, they start to guess what the  
 208 speaker may want to say. However, guessing is not always correct. Therefore, the comprehension  
 209 performance is decreasing. The spatial distribution of inter-brain clusters is also different from  
 210 the distribution of the acoustic and semantic clusters. While the central electrodes are primarily  
 211 involved in the acoustic and semantic modulations, only the left-frontal electrodes are recruited in  
 212 the inter-brain modulation. The left frontal regions play a critical region in the language process  
 213 (Har-shai Yahav and Zion Golumbic, 2021; Hickok and Poeppel, 2007) and are a 'high-order' area in  
 214 attention selection (Zion Golumbic et al., 2013). The left-frontal electrodes may indicate a unique  
 215 contribution of IFG when listeners are under adverse listening conditions in previous inter-brain  
 216 studies (Dai et al., 2018; Li and Pylkkänen, 2021; Liu et al., 2020a). Our study highlighted the crit-  
 217 ical role of the speaker in the attention process. While speech itself may only serve as a trigger  
 218 and an entrainment signal, the attended speaker and the listener aligned their mind "beyond the  
 219 speech stimulus" (Hartley and Poeppel, 2020) even before the speech onset. To our knowledge,  
 220 our study is the first study to combine the NLP method and the inter-brain method to extract the  
 221 different levels of features in speech and investigate their attentional modulation, which gives rise  
 222 to an integrative and "beyond the stimulus" perspective (Hartley and Poeppel, 2020). In our study,  
 223 the speech was separated into three levels features: the acoustic feature, the semantic feature,  
 224 and the inter-brain feature. With the help of the NLP models (Armeni et al., 2019; Brodbeck and  
 225 Simon, 2022; Broderick et al., 2018; Kingma and Ba, 2015), we could calculate semantic features in  
 226 the text to which the listeners attended. Inspired by the inter-brain studies (Dai et al., 2018; Jiang



227 *et al., 2015, 2012; Leong et al., 2017; Stephens et al., 2010*), how the listeners pay attention to the  
228 “hidden meaning” behind the text was also analyzed. Meanwhile, we applied the TRF method to  
229 describe the precise neural activity towards the attended and unattended features and simultane-  
230 ously compared entrainment to the different features. In conclusion, our study used the attention  
231 as a spotlight and revealed that the listener would strike the different neural notes at distinct stages  
232 in an integrative way: the acoustic note is struck on the theta bands at first, and the semantic note  
233 comes later and lasts longer on the delta band. The striking on the mental instrument, which is  
234 achieved by the inter-brain coupling, appears even before the speech onset. Our study depicts the  
235 temporal dynamics of the attentional modulation and the functional roles of different frequency  
236 bands, which contributes to the old “cocktail party” a new integrative perspective.

## 237 **Methods and Materials**

### 238 **Ethics statement**

239 The study was conducted in accordance with the Declaration of Helsinki and was approved by the  
240 local Ethics Committee of Tsinghua University. Written informed consent was obtained from all  
241 participants.

### 242 **Participants**

243 Two participants (both male, aged 26 and 24 years) were recruited for this study as speakers. Both  
244 speakers were from the broadcasting station of Tsinghua University and had experience related  
245 to broadcasting and hosting. Twenty college students (10 females; mean age: 24.7 years; range:  
246 20–43 years) from Tsinghua University participated in the study as paid volunteers for listeners.  
247 All participants were native Chinese speakers and reported having normal hearing and normal or  
248 corrected-to-normal vision. The sample size ( $N = 20$ ) was decided empirically following previous  
249 TRF-based studies on human speech processing (*Broderick et al., 2018; Di Liberto et al., 2015; Li*  
250 *et al., 2022; Mirkovic et al., 2015*).

### 251 **Experimental procedure for the speakers**

252 A sequential inter-brain approach was adopted by the present study (*Redcay and Schilbach, 2019*),  
253 in which the neural activities of the speakers were recorded prior to the listeners. The sequen-  
254 tial design was more appropriate for this study than the real-time interactive design because the  
255 speakers’ audio and neural activity remained consistent for all listeners (*Leong et al., 2017; Liu et al.,*  
256 *2017; Stephens et al., 2010*). In this experiment, each speaker participated in 30 trials, each of  
257 which was approximately 51–76 seconds in length, while the speakers’ audio signals and EEG sig-  
258 nals were recorded. The experimenter selected 28 trials for the listener’s experiment, excluding  
259 the two most unqualified trials. The speaker first read the relevant material on the screen. There  
260 was a wide variety of content to be covered, including one’s hometown, a recent book, a fable, etc.  
261 The speaker could decide how long they wanted to spend on preparation and start talking when  
262 they were fully prepared (the length of preparation was usually 3 minutes). When the speaker was  
263 prepared, they would press the space bar on the computer keyboard, and the recording would  
264 begin. When the space bar was pressed, three 1000 Hz pure tone cues were triggered (duration:  
265 1000 ms; cue interval: 1500 ms). The cues were presented as event markers, synchronized with  
266 the sound in the listener’s experiment to ensure that the neural signals of the speaker and listener  
267 remain aligned with the sound stimuli. The speaker was asked to start speaking immediately after  
268 the end of the third beep (within approximately 3 s). A fixation and a countdown timer appeared on  
269 the screen during the talking part. The speaker was asked to stare at the fixation and to complete  
270 the speaking as clearly, completely, and naturally as possible. During the recording process, the  
271 experimenter listened to the speaker’s narration simultaneously and controlled the quality. The  
272 experimenter had the right to ask the speaker to retell the clip if there was a reason for the lack  
273 of fluency, length, etc., that might affect the listeners’ perception. The materials of both speakers’

274 content were varied. Between each trial, the speakers were allowed to rest on their own. During  
275 the experiment, the speakers were asked to control their head movements and facial muscles to  
276 obtain better quality EEG signals. The speech stimuli were recorded from two male speakers using  
277 the microphone of an iPad2 mini (Apple Inc., Cupertino, CA) at a sampling rate of 44,100 Hz.

### 278 **Experimental procedure for the listeners**

279 The experiment consisted of four blocks, each containing seven trials. Two speech streams were  
280 presented simultaneously in each trial, one to the left ear and the other to the right ear. Two  
281 speech streams of the same trial matched the volume, i.e., the root mean squared intensity of  
282 the amplitude of the speech streams in the same trial were the same. The participants were in-  
283 structed to attend to one spatial side according to the hints on the screen (“Please pay attention to  
284 the [LEFT/RIGHT]”). Considering the possible duration difference between the two audio streams,  
285 the researchers set the end of the trial after the longer speech audio had ended. Each trial be-  
286 gan when participants pressed the SPACE key on the computer keyboard. A white fixation cross  
287 was also displayed throughout the trial. The speech stimuli were played immediately after the key-  
288 press and were preceded by the three beep sounds. At the end of each trial, four multiple-choice  
289 questions (two for the attended story and the other two for the unattended story) were presented  
290 sequentially in random order on the computer screen. Each question had four options, and par-  
291 ticipants entered the letter of the correct option as their answer. The listeners were not explicitly  
292 informed about the correspondence between the questions and the stories. For instance, one  
293 question following a story about one’s hometown was, “What is the most dissatisfying thing about  
294 the speaker’s hometown? (推测讲述人对于家乡最不满意的的地方在于?)”, and the four choices were A)  
295 There is no heating in winter; B) There are no hot springs in summer; C) There is no fruit in autumn;  
296 D) There are no flowers in spring (A. 冬天没暖气; B. 夏天没温泉; C. 秋天没水果; D. 春天没鲜花). The  
297 single-trial comprehension accuracy could be 0% (two wrong answers), 50% (one correct answer),  
298 or 100% (two correct answers) for both the attended and the unattended stories. No feedback on  
299 whether the questions were answered correctly or not. After completing these questions, partici-  
300 pants rated their concentration level of the attended stream, the experienced difficulty performing  
301 the attention task, and the familiarity with the attended material using three 10-point Likert scales.  
302 Throughout the trial, participants were required to maintain visual fixation on the fixation cross  
303 while listening to the speech. Meanwhile, they were asked to minimize eye blinks and all other  
304 motor activities. The participants were recommended to take a short break (around 1 min) after  
305 every trial within one block and a long break (no longer than 10 min) between blocks. In each  
306 block, the side being attended to was fixed (two blocks for attending to the left side and two for  
307 attending to the right side). Within each block, the identity of the speaker is kept constant on the  
308 left and right sides. The to-be-attended spatial side and the corresponding speaker identity were  
309 balanced within the participant, with seven trials per side for both speakers. The assignment of the  
310 stories to the four blocks was randomized across the participants. The experiment was conducted  
311 in a sound-attenuated, dimly lit, and electrically shielded room. The participants were seated in  
312 a comfortable chair in front of a 19.7-inch LCD monitor (Lenovo LT2013s). The viewing distance  
313 was approximately 60 cm. The experimental procedure was programmed in MATLAB using the  
314 Psychophysics Toolbox 3.0 extensions (Brainard and Brainard, 1997). The speech stimuli were de-  
315 livered binaurally via an air-tube earphone (Etymotic ER2, Etymotic Research, Elk Grove Village, IL,  
316 USA) to avoid possible electromagnetic interference from auditory devices. The volume of the au-  
317 dio stimuli was adjusted to be at a comfortable level (70 dB SPL) that was well above the auditory  
318 threshold. The average presentation level was measured with a BK (Brüel Kjær, Nærum, Denmark)  
319 Sound Level Meter (Type 2250 Investigator) with a 1-inch Free-field Microphone (Type 4144) and  
320 an Artificial Ear (Type 4152).

## 321 **Data acquisition and pre-processing**

322 EEG was recorded from 60 electrodes (FP1/2, FPZ, AF3/4, F7/8, F5/6, F3/4, F1/2, FZ, FT7/8, FC5/6,  
323 FC3/4, FC1/2, FCZ, T7/8, C5/6, C3/4, C1/2, CZ, TP7/8, CP5/6, CP3/4, CP1/2, CPZ, P7/8, P5/6, P3/4,  
324 P1/2, PZ, PO7/8, PO5/6, PO3/4, POZ, Oz, and O1/2), which were referenced to an electrode be-  
325 tween Cz and CPz, with a forehead ground at Fz. A NeuroScan amplifier (SynAmp II, NeuroScan,  
326 Compumedics, USA) was used to record EEG at a sampling rate of 1000 Hz. Electrode impedances  
327 were kept below ten kOhm for all electrodes. The recorded EEG data were first notch filtered to  
328 remove the 50 Hz powerline noise and then subjected to an artifact rejection procedure using inde-  
329 pendent component analysis. Independent components (ICs) with large weights over the frontal or  
330 temporal areas, together with a corresponding temporal course showing eye movement or muscle  
331 movement activities, were removed. The remaining ICs were then back-projected onto the scalp  
332 EEG channels, reconstructing the artifact-free EEG signals. While the relatively long duration of the  
333 speech trials in the present study (about 1 minute per story, see Experimental procedure) has made  
334 it more difficult for the participants to avoid inducing movement-related artifacts as compared to  
335 the classical ERP-based studies, a temporally continuous, non-interrupted EEG segment per trial  
336 was preferred for the employment of the CCA method. Therefore, any ICs with artifact-like EEG ac-  
337 tivities for more than 20% of the trial time (i.e., about 12 sec) were rejected, leading to around 4–11  
338 ICs rejected per participant. The cleaned EEG data were used for the mTRF analysis without any  
339 further artifact rejection procedures. Next, the EEG data were segmented into 28 trials according  
340 to the markers representing speech onsets. The analysis window for each trial was extended from  
341 10 to 55 s (duration: 45 s) to avoid the onset and the offset of the stories. The pre-processed EEG  
342 signals were re-referenced to the average of all scalp channels and then downsampled to 128 Hz  
343 before the modeling. Then, the EEG data were filtered in delta (1–4 Hz) and theta (4–8 Hz) (filter or-  
344 der: 64, one-pass forward filter). The use of a causal FIR filter ensured that filtered EEG signals were  
345 decided only by the current and previous data samples (*de Cheveigné and Nelken, 2019*), which is  
346 essential for accurate time-course analysis. The filter order of 64 was chosen to keep a balance  
347 of temporal resolution and filter performance: the filtered EEG signals were therefore calculated  
348 based on the preceding 500 ms data (64 at 128 Hz).

## 349 **Speech Representations**

### 350 Acoustic Features

351 The amplitude envelope of the speech represented the acoustic features of the speech. It was  
352 obtained using a Hilbert transform and then down-sampled to the same sampling rate of 128 Hz.

### 353 Semantic Features

354 The original audio recorded by the speaker during the EEG recording was converted to the text  
355 firstly automatically by *Iflyrec* software (Iflytek Co., Ltd, Hefei, Anhui) and then double checked  
356 manually. The onset time of every word was extracted during this process. The recent emergence  
357 of Natural Process Language (NLP) models has enabled the description of the semantic features  
358 in speech (*Brookshire, 2022; Broderick et al., 2021, 2018*). Next word prediction is one of the fun-  
359 damental NLP tasks using the semantic information in the texts (*Schrimpf et al., 2021; Vaswani*  
360 *et al., 2017*). The goal of the task is to predict the next word when given a sequence of words  
361  $W_1, W_2, \dots, W_{t-1}$  which was consistent with the human understanding process. The probability of  
362 the next word is  $P(W_t | W_1, \dots, W_{t-1})$  and can be calculated by varied NLP models. The surprisal of  
363 the word was defined as follow, which reflected how surprised the next word (*Willems and Jacobs,*  
364 *2016*):

$$surprisal(t) = \log P(W_t | W_1, \dots, W_{t-1}) \quad (1)$$

365 The index was calculated based on ADAM, a widely accepted classical natural language process  
366 model (*Bengio et al., 2003; Kingma and Ba, 2015*). The model was trained on a couple of People's

367 Daily. There were 534,246 words involved in the model training. 66,781 words were in the cross-  
368 validation set, and 66,781 words used as a test set. The details of the model are described in Table  
369 S1. After calculating the surprisal index of every word, we generated a “semantic vector” at the  
370 same sampling rate as the EEG data (*Broderick et al., 2018*). The vectors contained the time-aligned  
371 impulses at the start of each word of the surprisal value for every audio clip.

### 372 Inter-brain Features

373 The Inter-brain recording method enables us to study how the listeners are aligned with the at-  
374 tended speaker (*Dai et al., 2018; Hasson et al., 2012; Jiang et al., 2021; Stephens et al., 2010*). We  
375 used the speaker’s neural activity as the representation of the inter-brain feature. The speaker’s  
376 EEG served as the inter-brain feature. It followed the same pre-processing procedure as the lis-  
377 tener’s EEG.

### 378 Temporal response function modeling

379 The analysis workflow for the analysis related to the attended speech stream is shown in Figure  
380 1. The neural responses to the three different features were characterized using a temporal re-  
381 sponse function (TRF)-based modeling method (*Crosse et al., 2016, 2021*). Three different features  
382 mentioned above are the input signal required by TRF. The corresponding neural response  $r(t, n)$   
383 can be formulated as follows:

$$r(t, n) = \sum_{\tau} w(\tau, n)S(t - \tau) + \varepsilon(t, n) \quad (2)$$

384 where  $r(t, n)$  is the actual EEG response at every channel  $n$ ;  $n = 1, \dots, T$  time point;  $S(t - \tau)$  means  
385 the multivariate stimulus representation;  $w(n, \tau)$  the channel specific TRF at lag and  $\varepsilon(t, n)$  is the  
386 residual. The TRF is estimated by minimizing the mean square error between the actual neural  
387 response  $r(t, n)$  and the neural response predicted by the model  $\hat{r}(t, n)$ . The Pearson’s correlation  
388 between the actual neural response and predicted neural response was referred as Encoding  $r$ .  
389 The mTRF toolbox (*Crosse et al., 2016*) was used to estimate the  $TRF(w)$  as follow:

$$w = (S^T S + \lambda I)^{-1} S^T r \quad (3)$$

390 where the lambda( $\lambda$ ) is the ridge regression parameter,  $I$  is the identity matrix, and the matrix  
391  $S$  is the stimulus matrix. The lambda varied from  $10^{-1}$  to  $10^8$  ( $\lambda = 10^{-1}, 10^0, \dots, 10^8$ ) to make  
392 the model optimal (*Crosse et al., 2021*). The lambda value, which produces the highest encoding  $r$ ,  
393 averaged across trials and channels, was selected as the regularization parameter for all trials per  
394 participant (*Broderick et al., 2019*). The cross-validation procedure was implemented in a leave-  
395 one-trial-out manner: the TRFs were trained based on data from 27 trials and tested on the left-  
396 out trial each time. The TRF was trained at individual time lags of -8 s to 8 s to investigate the  
397 specific interval of attentional modulation of each feature. At a sampling rate of 128 Hz, there  
398 are 2049 individual time-lag intervals of 7.625 ms. The TRF calculation procedure was performed  
399 for the EEG signals from each EEG channel filtered at the four frequency bands. Only attended  
400 features are used as input to the model, and TRFs trained by the attended features and the neural  
401 response were applied to the tests of attended and unattended features, referring to Encoding  $r_{att}$   
402 and Encoding  $r_{unatt}$ , respectively.

### 403 Quantification and statistical analysis

404 The paired  $t$ -tests were performed to investigate the attentional modulation of different features,  
405 contrasting the encoding  $r$  of the attended speech versus the unattended speech at each time lag.  
406 Encoding  $r$  was normalized using the Fisher-  $z$  transform (*Corey et al., 1998*). A nonparametric  
407 cluster-based permutation analysis was applied to account for multiple comparisons (*Maris et al.,*  
408 *2007*). In this procedure, neighboring channel-latency bins with uncorrected  $t$ -tests  $p$ -value below  
409 0.05 were combined into clusters, for which the sum of the correlational  $t$ -statistics corresponding

410 to the  $t$ -tests was obtained. The combing process was initially automated by the toolbox and then  
411 manually double-checked. Two clusters were combined if they shared a similar spatial distribu-  
412 tion or time lag. A null distribution was created through permutations of data across participants  
413 ( $n = 1,000$  permutations), which defined the maximum cluster-level test statistics and corrected  
414  $p$ -values for each cluster. Clusters with  $p$ -values below 0.01 based on clusters were selected for  
415 further analysis. The above statistical analysis followed the standard cluster-based permutation  
416 procedure as employed in classical ERP and related studies (*Arnal et al., 2011; Henry and Obleser,*  
417 *2012; Zhang et al., 2012*). Note that the reported  $p$ -values were only corrected for the tests per-  
418 formed within each frequency band by using cluster-based permutation tests. No multiple com-  
419 parison correction was employed across different frequency bands.

#### 420 **Correlation between clusters**

421 The *Spearman* correlation of Encoding  $r$ -att in each cluster was calculated for each pair of clusters  
422 to analyze the correlation between them.

#### 423 **Partial correlation**

424 The partial correlation between every cluster and the comprehension performance was calculated  
425 to investigate the unique contribution of each cluster to the behavioral performance.

#### 426 **Acknowledgments**

427 This work was supported by the National Science Foundation of China (NSFC) and the German Re-  
428 search Foundation (DFG) in project Crossmodal Learning (grant number: NSFC 61621136008/DFG  
429 TRR-169/C1, B1), and the National Natural Science Foundation of China (grant number: 61977041).  
430 The authors would like to thank Prof. Dr. Xiaoqin Wang and Dr. Yue Ding for providing the shielded  
431 room for the experiment as well as necessary technical support. The authors would like to acknowl-  
432 edge Prof. Dr. Zhiyuan Liu and members of his lab for computing the NLP models.

#### 433 **References**

- 434 **Anderson AJ**, Binder JR, Fernandino L, Humphries CJ, Conant LL, Raizada RDS, Lin F, Lalor EC. An Integrated  
435 Neural Decoder of Linguistic and Experiential Meaning. *The Journal of Neuroscience*. 2019 nov; 39(45):8969-  
436 8987. <http://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.2575-18.2019>, doi: 10.1523/JNEUROSCI.2575-  
437 18.2019.
- 438 **Armeni K**, Willems RM, van den Bosch A, Schoffelen JM. Frequency-specific brain dynamics related  
439 to prediction during language comprehension. *NeuroImage*. 2019 sep; 198:283-295. <https://doi.org/10.1016/j.neuroimage.2019.04.083><https://linkinghub.elsevier.com/retrieve/pii/S1053811919304057>, doi:  
440 10.1016/j.neuroimage.2019.04.083.
- 442 **Arnal LH**, Wyart V, Giraud AL. Transitions in neural oscillations reflect prediction errors generated in audiovisual  
443 speech. *Nature Neuroscience*. 2011; 14(6):797-801. doi: 10.1038/nn.2810.
- 444 **Bashivan P**, Kar K, DiCarlo JJ. Neural population control via deep image synthesis. *Science*. 2019; 364(6439).  
445 <http://www.ncbi.nlm.nih.gov/pubmed/31048462>, doi: 10.1126/science.aav9436.
- 446 **Bengio Y**, Ducharme R, Vincent P. A neural probabilistic language model. *Journal of Machine Learning Research*.  
447 2003; 3:1137-1155.
- 448 **Bregman AS**. Auditory scene analysis: The perceptual organization of sound. Cambridge, MA, US: The MIT  
449 Press; 1990.
- 450 **Brodbeck C**, Presacco A, Simon JZ. Neural source dynamics of brain responses to continuous stimuli: Speech  
451 processing from acoustics to comprehension. *NeuroImage*. 2018; 172(January):162-174. [https://doi.org/10.](https://doi.org/10.1016/j.neuroimage.2018.01.042)  
452 [1016/j.neuroimage.2018.01.042](https://doi.org/10.1016/j.neuroimage.2018.01.042), doi: 10.1016/j.neuroimage.2018.01.042.
- 453 **Brodbeck C**, Simon JZ. Cortical tracking of voice pitch in the presence of multiple speakers depends on selective  
454 attention. *Frontiers in Neuroscience*. 2022 aug; 16(August):1-11. [https://www.frontiersin.org/articles/10.3389/fnins.2022.828546/](https://www.frontiersin.org/articles/10.3389/fnins.2022.828546/full)  
455 [full](https://www.frontiersin.org/articles/10.3389/fnins.2022.828546/full), doi: 10.3389/fnins.2022.828546.

- 456 **Broderick MP**, Anderson AJ, Di Liberto GM, Crosse MJ, Lalor EC. Electrophysiological Correlates of Semantic  
457 Dissimilarity Reflect the Comprehension of Natural, Narrative Speech. *Current Biology*. 2018 mar; 28(5):803–  
458 809.e3. <http://linkinghub.elsevier.com/retrieve/pii/S0960982218301465><https://linkinghub.elsevier.com/retrieve/pii/S0960982218301465>, doi: 10.1016/j.cub.2018.01.080.
- 460 **Broderick MP**, Anderson AJ, Lalor EC. Semantic Context Enhances the Early Auditory Encoding of Natural  
461 Speech. *The Journal of Neuroscience*. 2019 sep; 39(38):7564–7575. <http://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.0584-19.2019>, doi: 10.1523/JNEUROSCI.0584-19.2019.
- 463 **Broderick MP**, Di Liberto GM, Anderson AJ, Rofes A, Lalor EC. Dissociable electrophysiological measures of nat-  
464 ural language processing reveal differences in speech comprehension strategy in healthy ageing. *Scientific*  
465 *Reports*. 2021; 11(1):1–12. <https://doi.org/10.1038/s41598-021-84597-9>, doi: 10.1038/s41598-021-84597-9.
- 466 **Brookshire G**. Putative rhythms in attentional switching can be explained by aperiodic temporal structure. *Nature*  
467 *Human Behaviour*. 2022 jun; <https://www.nature.com/articles/s41562-022-01364-0>, doi: 10.1038/s41562-  
468 022-01364-0.
- 469 **Cherry EC**. Some experiments on the recognition of speech, with one and with two ears. *The Journal of the*  
470 *Acoustical Society of America*. 1953; 25(5):975–979.
- 471 **de Cheveigné A**, Nelken I. Filters: When, Why, and How (Not) to Use Them. *Neuron*. 2019 apr; 102(2):280–293.  
472 <https://linkinghub.elsevier.com/retrieve/pii/S0896627319301746>, doi: 10.1016/j.neuron.2019.02.039.
- 473 **Connolly JF**, Phillips NA. Event-related potential components reflect phonological and semantic processing  
474 of the terminal word of spoken sentences. *Journal of Cognitive Neuroscience*. 1994; 6(3):256–266. doi:  
475 10.1162/jocn.1994.6.3.256.
- 476 **Connolly JF**, Stewart SH, Phillips NA. The effects of processing requirements on neurophysiological responses  
477 to spoken sentences. *Brain and language*. 1990; 39(2):302–318.
- 478 **Corey DM**, Dunlap WP, Burke MJ. Averaging correlations: Expected values and bias in combined pear-  
479 son rs and fisher’s z transformations. *Journal of General Psychology*. 1998; 125(3):245–261. doi:  
480 10.1080/00221309809595548.
- 481 **Crosse MJ**, Di Liberto GM, Lalor EC. Eye can hear clearly now: Inverse effectiveness in natural audiovisual  
482 speech processing relies on long-term crossmodal temporal integration. *Journal of Neuroscience*. 2016;  
483 36(38):9888–9895. doi: 10.1523/JNEUROSCI.1396-16.2016.
- 484 **Crosse MJ**, Zuk NJ, Di Liberto GM, Nidiffer AR, Molholm S, Lalor EC. Linear Modeling of Neurophysiological  
485 Responses to Speech and Other Continuous Stimuli: Methodological Considerations for Applied Research.  
486 *Frontiers in Neuroscience*. 2021 nov; 15. <https://www.frontiersin.org/articles/10.3389/fnins.2021.705621/full>,  
487 doi: 10.3389/fnins.2021.705621.
- 488 **Dai B**, Chen C, Long Y, Zheng L, Zhao H, Bai X, Liu W, Zhang Y, Liu L, Guo T, Ding G, Lu C. Neural mechanisms for  
489 selectively tuning in to the target speaker in a naturalistic noisy situation. *Nature Communications*. 2018 dec;  
490 9(1):1–12. <http://dx.doi.org/10.1038/s41467-018-04819-z><https://www.nature.com/articles/s41467-018-04819-z>,  
491 doi: 10.1038/s41467-018-04819-z.
- 492 **Dai B**, McQueen JM, Terporten R, Hagoort P, Kösem A. Distracting linguistic information impairs neural tracking  
493 of attended speech. *Current Research in Neurobiology*. 2022; 3:100043. doi: 10.1016/j.crneur.2022.100043.
- 494 **Di Liberto GM**, O’Sullivan JA, Lalor EC. Low-Frequency Cortical Entrainment to Speech Reflects Phoneme-Level  
495 Processing. *Current Biology*. 2015 oct; 25(19):2457–2465. <http://dx.doi.org/10.1016/j.cub.2015.08.030><https://linkinghub.elsevier.com/retrieve/pii/S0960982215010015>, doi: 10.1016/j.cub.2015.08.030.
- 497 **Ding N**, Chatterjee M, Simon JZ. Robust cortical entrainment to the speech envelope relies on the spectro-  
498 temporal fine structure. *NeuroImage*. 2014; 88:41–46. <http://dx.doi.org/10.1016/j.neuroimage.2013.10.054>,  
499 doi: 10.1016/j.neuroimage.2013.10.054.
- 500 **Ding N**, Simon JZ. Emergence of neural encoding of auditory objects while listening to competing speakers.  
501 *Proceedings of the National Academy of Sciences of the United States of America*. 2012; 109(29):11854–9.  
502 <http://www.pnas.org/content/109/29/11854.full>, doi: 10.1073/pnas.1205381109.
- 503 **Ding N**, Simon JZ. Neural coding of continuous speech in auditory cortex during monaural and dichotic listening.  
504 *Journal of Neurophysiology*. 2012 jan; 107(1):78–89. <https://www.physiology.org/doi/10.1152/jn.00297.2011>,  
505 doi: 10.1152/jn.00297.2011.

- 506 **Etard O**, Kegler M, Braiman C, Forte AE, Reichenbach T. Decoding of selective attention to continuous speech  
507 from the human auditory brainstem response. *NeuroImage*. 2019; 200(November 2018):1–11. <https://doi.org/10.1016/j.neuroimage.2019.06.029>, doi: 10.1016/j.neuroimage.2019.06.029.
- 509 **Frank SL**, Willems RM. Word predictability and semantic similarity show distinct patterns of brain activity  
510 during language comprehension. *Language, Cognition and Neuroscience*. 2017; 32(9):1192–1203. doi:  
511 10.1080/23273798.2017.1323109.
- 512 **Friston K**. The free-energy principle: a rough guide to the brain? *Trends in Cognitive Sciences*. 2009; 13(7):293–  
513 301. doi: 10.1016/j.tics.2009.04.005.
- 514 **Gadamer HG**. *Truth and Method*. A Continuum book, Seabury Press; 1975. [https://books.google.co.jp/books?](https://books.google.co.jp/books?id=zQnXAAAAMAAJ)  
515 [id=zQnXAAAAMAAJ](https://books.google.co.jp/books?id=zQnXAAAAMAAJ).
- 516 **Gillis M**, Vanthornhout J, Simon JZ, Francart T, Brodbeck C. Neural Markers of Speech Comprehension: Mea-  
517 suring EEG Tracking of Linguistic Speech Representations, Controlling the Speech Acoustics. *The Journal of*  
518 *neuroscience*. 2021; 41(50):10316–10329. doi: 10.1523/JNEUROSCI.0812-21.2021.
- 519 **Hamilton LS**, Edwards E, Chang EF. A Spatial Map of Onset and Sustained Responses to Speech in the Human  
520 Superior Temporal Gyrus. *Current Biology*. 2018; 28(12):1860–1871.e4. [https://doi.org/10.1016/j.cub.2018.](https://doi.org/10.1016/j.cub.2018.04.033)  
521 [04.033](https://doi.org/10.1016/j.cub.2018.04.033), doi: 10.1016/j.cub.2018.04.033.
- 522 **Har-shai Yahav P**, Zion Golumbic E. Linguistic processing of task-irrelevant speech at a cocktail party. *eLife*.  
523 2021; 10:1–24. doi: 10.7554/eLife.65096.
- 524 **Hartley CA**, Poeppel D. Beyond the Stimulus: A Neurohumanities Approach to Language, Mu-  
525 sic, and Emotion. *Neuron*. 2020; 108(4):597–599. <https://doi.org/10.1016/j.neuron.2020.10.021>, doi:  
526 10.1016/j.neuron.2020.10.021.
- 527 **Hasson U**, Ghazanfar AA, Galantucci B, Garrod S, Keysers C. Brain-to-brain coupling: A mechanism for creating  
528 and sharing a social world. *Trends in Cognitive Sciences*. 2012; 16(2):114–121. [http://dx.doi.org/10.1016/j.](http://dx.doi.org/10.1016/j.tics.2011.12.007)  
529 [tics.2011.12.007](http://dx.doi.org/10.1016/j.tics.2011.12.007), doi: 10.1016/j.tics.2011.12.007.
- 530 **Heil M**, Rolke B, Pecchinenda A, Heil M, Rolke B, Pecchinenda A. Automatic Semantic Activation Is No Myth.  
531 *Psychological Science*. 2004; 15(12):852–857. doi: [doi.org/10.1111](https://doi.org/10.1111).
- 532 **Heilbron M**, Armeni K, Schoffelen JM, Hagoort P, de Lange FP. A hierarchy of linguistic predictions dur-  
533 ing natural language comprehension. *Proceedings of the National Academy of Sciences*. 2022 aug;  
534 119(32):2020.12.03.410399. <https://doi.org/10.1101/2020.12.03.410399>[https://pnas.org/doi/full/10.1073/pnas.](https://pnas.org/doi/full/10.1073/pnas.2201968119)  
535 [2201968119](https://pnas.org/doi/full/10.1073/pnas.2201968119), doi: 10.1073/pnas.2201968119.
- 536 **Henry MJ**, Obleser J. Frequency modulation entrains slow neural oscillations and optimizes human listening  
537 behavior. *Proceedings of the National Academy of Sciences*. 2012; 109(49):20095–20100. [http://www.pnas.](http://www.pnas.org/cgi/doi/10.1073/pnas.1213390109)  
538 [org/cgi/doi/10.1073/pnas.1213390109](http://www.pnas.org/cgi/doi/10.1073/pnas.1213390109), doi: 10.1073/pnas.1213390109.
- 539 **Hickok G**, Poeppel D. The cortical organization of speech processing. *Nature Reviews Neuroscience*.  
540 2007 may; 8(5):393–402. <http://dx.doi.org/10.1038/nrn2113><http://www.nature.com/articles/nrn2113>, doi:  
541 10.1038/nrn2113.
- 542 **Hillyard SA**, Hink RF, Schwent VL, Picton TW. Electrical signs of selective attention in the human brain. *Science*.  
543 1973; 182(4108):177–180. doi: 10.1126/science.182.4108.177.
- 544 **Jiang J**, Chen C, Dai B, Shi G, Ding G, Liu L, Lu C. Leader emergence through interpersonal neural synchronization.  
545 *Proceedings of the National Academy of Sciences*. 2015; 112(14):4274–4279. [http://www.pnas.org/lookup/](http://www.pnas.org/lookup/doi/10.1073/pnas.1422930112)  
546 [doi/10.1073/pnas.1422930112](http://www.pnas.org/lookup/doi/10.1073/pnas.1422930112), doi: 10.1073/pnas.1422930112.
- 547 **Jiang J**, Dai B, Peng D, Zhu C, Liu L, Lu C. Neural synchronization during face-to-face communication. *Journal*  
548 *of Neuroscience*. 2012; 32(45):16064–16069. doi: 10.1523/JNEUROSCI.2926-12.2012.
- 549 **Jiang J**, Zheng L, Lu C. A hierarchical model for interpersonal verbal communication. *Social Cognitive and Affec-*  
550 *tive Neuroscience*. 2021 jan; 16(1-10):246–255. <https://academic.oup.com/scan/article/16/1-2/246/5956560>,  
551 doi: 10.1093/scan/nsaa151.
- 552 **Kingma DP**, Ba JL. Adam: A method for stochastic optimization. 3rd International Conference on Learning  
553 Representations, ICLR 2015 - Conference Track Proceedings. 2015; p. 1–15.

- 554 **Kuhlen AK**, Allefeld C, Haynes JD. Content-specific coordination of listeners' to speakers' EEG during com-  
555 munication. *Frontiers in human neuroscience*. 2012; 6(October):1–15. [http://www.pubmedcentral.nih.gov/](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3461523&tool=pmcentrez&rendertype=abstract)  
556 [articlerender.fcgi?artid=3461523&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3461523&tool=pmcentrez&rendertype=abstract), doi: 10.3389/fnhum.2012.00266.
- 557 **Kutas M**, Hillyard SA. Brain potentials during reading reflect word expectancy and semantic association. *Nature*.  
558 1984; 307(5947):161–163.
- 559 **Kutas M**, Hillyard SA. An Electrophysiological Probe of Incidental Semantic Association. *Journal*  
560 *of Cognitive Neuroscience*. 1989 jan; 1(1):38–49. [https://direct.mit.edu/jocn/article/1/1/38/2956/](https://direct.mit.edu/jocn/article/1/1/38/2956/An-Electrophysiological-Probe-of-Incidental)  
561 [An-Electrophysiological-Probe-of-Incidental](https://direct.mit.edu/jocn/article/1/1/38/2956/An-Electrophysiological-Probe-of-Incidental), doi: 10.1162/jocn.1989.1.1.38.
- 562 **Lalor EC**, Foxe JJ. Neural responses to uninterrupted natural speech can be extracted with precise temporal  
563 resolution. *European Journal of Neuroscience*. 2010; 31(1):189–193. doi: 10.1111/j.1460-9568.2009.07055.x.
- 564 **Lau EF**, Phillips C, Poeppel D. A cortical network for semantics: (de)constructing the N400. *Nature Reviews*  
565 *Neuroscience*. 2008; 9(12):920–933. doi: Doi 10.1038/Nrn2532.
- 566 **Leong V**, Byrne E, Clackson K, Georgieva S, Lam S, Wass S. Speaker gaze increases information coupling between  
567 infant and adult brains. *Proceedings of the National Academy of Sciences of the United States of America*.  
568 2017; 114(50):13290–13295. doi: 10.1073/pnas.1702493114.
- 569 **Li J**, Pykkänen L. Disentangling semantic composition and semantic association in the left temporal lobe. *Jour-*  
570 *nal of Neuroscience*. 2021; 41(30):6526–6538. doi: 10.1523/JNEUROSCI.2317-20.2021.
- 571 **Li X**, Huang L, Yao P, Hyönä J. Universal and specific reading mechanisms across different writing systems.  
572 *Nature Reviews Psychology*. 2022 feb; 0123456789. <https://www.nature.com/articles/s44159-022-00022-6>,  
573 doi: 10.1038/s44159-022-00022-6.
- 574 **Liu L**, Zhang Y, Zhou Q, Garrett DD, Lu C, Chen A, Qiu J, Ding G. Auditory-Articulatory Neural Alignment between  
575 Listener and Speaker during Verbal Communication. *Cerebral Cortex*. 2020; 30(3):942–951. doi: 10.1093/cer-  
576 [cor/bhz138](http://cercor.oup.com/bhz138).
- 577 **Liu Y**, Li M, Zhang X, Lu Y, Gong H, Yin J, Chen Z, Qian L, Yang Y, Andolina IM, Shipp S, Mcloughlin N, Tang S, Wang  
578 W. Hierarchical Representation for Chromatic Processing across Macaque V1, V2, and V4. *Neuron*. 2020 nov;  
579 108(3):538–550.e5. <https://doi.org/10.1016/j.neuron.2020.07.037>[https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S089662732030581X)  
580 [S089662732030581X](https://linkinghub.elsevier.com/retrieve/pii/S089662732030581X), doi: 10.1016/j.neuron.2020.07.037.
- 581 **Liu Y**, Piazza EA, Simony E, Shewokis PA, Onaral B, Hasson U, Ayaz H. Measuring speaker–listener neural  
582 coupling with functional near infrared spectroscopy. *Scientific Reports*. 2017 apr; 7(1):43293. [http:](http://www.nature.com/articles/srep43293)  
583 [://www.nature.com/articles/srep43293](http://www.nature.com/articles/srep43293), doi: 10.1038/srep43293.
- 584 **Lu Y**, Jin P, Ding N, Tian X. Delta-band neural tracking primarily reflects rule-based chunking instead of se-  
585 mantic relatedness between words. *Cerebral Cortex*. 2022 sep; p. 1–11. [https://academic.oup.com/cercor/](https://academic.oup.com/cercor/advance-article/doi/10.1093/cercor/bhac354/6702814)  
586 [advance-article/doi/10.1093/cercor/bhac354/6702814](https://academic.oup.com/cercor/advance-article/doi/10.1093/cercor/bhac354/6702814), doi: 10.1093/cercor/bhac354.
- 587 **Maris E**, Schoffelen JM, Fries P. Nonparametric statistical testing of coherence differences. *Journal of Neuro-*  
588 *science Methods*. 2007; 163(1):161–175. doi: 10.1016/j.jneumeth.2007.02.011.
- 589 **McCarthy G**, Nobre AC. Modulation of semantic processing by spatial selective attention. *Electroen-*  
590 *cephalography and Clinical Neurophysiology/ Evoked Potentials*. 1993; 88(3):210–219. doi: 10.1016/0168-  
591 5597(93)90005-A.
- 592 **Mesik J**, Ray L, Wojtczak M. Effects of Age on Cortical Tracking of Word-Level Features of Continuous Competing  
593 Speech. *Frontiers in Neuroscience*. 2021; 15(April):1–21. doi: 10.3389/fnins.2021.635126.
- 594 **Middlebrooks JC**, Simon JZ, Popper AN, Fay RR. The Auditory System at the Cocktail Party. Middlebrooks JC, Si-  
595 mon JZ, Popper AN, Fay RR, editors, *Springer Handbook of Auditory Research*, Cham: Springer International  
596 Publishing; 2017. [http://download.springer.com/static/pdf/730/bok%253A978-3-319-51662-2.pdf?originUrl=](http://download.springer.com/static/pdf/730/bok%253A978-3-319-51662-2.pdf?originUrl=http%3A%2F%2Flink.springer.com%2Fbook%2F10.1007%2F978-3-319-51662-2&token2=exp=1497617002%5Csim%5Cac1=%2Fstatic%2Fpdf%2F730%2Fbok%25253A978-3-319-51662-2.pdf%3ForiginUrl%3Dhttp%25)  
597 [http%3A%2F%2Flink.springer.com%2Fbook%2F10.1007%2F978-3-319-51662-2&token2=exp=1497617002%5C](http://download.springer.com/static/pdf/730/bok%253A978-3-319-51662-2.pdf?originUrl=http%3A%2F%2Flink.springer.com%2Fbook%2F10.1007%2F978-3-319-51662-2&token2=exp=1497617002%5Csim%5Cac1=%2Fstatic%2Fpdf%2F730%2Fbok%25253A978-3-319-51662-2.pdf%3ForiginUrl%3Dhttp%25)  
598 [sim%5Cac1=%2Fstatic%2Fpdf%2F730%2Fbok%25253A978-3-319-51662-2.pdf%3ForiginUrl%3Dhttp%25](http://download.springer.com/static/pdf/730/bok%253A978-3-319-51662-2.pdf?originUrl=http%3A%2F%2Flink.springer.com%2Fbook%2F10.1007%2F978-3-319-51662-2&token2=exp=1497617002%5Csim%5Cac1=%2Fstatic%2Fpdf%2F730%2Fbok%25253A978-3-319-51662-2.pdf%3ForiginUrl%3Dhttp%25), doi:  
599 10.1007/978-3-319-51662-2.
- 600 **Mirkovic B**, Debener S, Jaeger M, De Vos M. Decoding the attended speech stream with multi-channel EEG:  
601 implications for online, daily-life applications. *Journal of Neural Engineering*. 2015 aug; 12(4):046007. [https:](https://iopscience.iop.org/article/10.1088/1741-2560/12/4/046007)  
602 [://iopscience.iop.org/article/10.1088/1741-2560/12/4/046007](https://iopscience.iop.org/article/10.1088/1741-2560/12/4/046007), doi: 10.1088/1741-2560/12/4/046007.



- 603 **Nastase SA**, Liu YF, Hillman H, Zadbood A, Hasenfratz L, Keshavarzian N, Chen J, Honey CJ, Yeshurun Y, Regev  
604 M, Nguyen M, Chang CHC, Baldassano C, Lositsky O, Simony E, Chow MA, Leong YC, Brooks PP, Micciche E,  
605 Choe G, et al. The “Narratives” fMRI dataset for evaluating models of naturalistic language comprehension.  
606 *Scientific Data*. 2021; 8(1):1–22. <http://dx.doi.org/10.1038/s41597-021-01033-3>, doi: 10.1038/s41597-021-  
607 01033-3.
- 608 **O’Sullivan JA**, Power AJ, Mesgarani N, Rajaram S, Foxe JJ, Shinn-Cunningham BG, Slaney M, Shamma SA, Lalor  
609 EC. Attentional Selection in a Cocktail Party Environment Can Be Decoded from Single-Trial EEG. *Cerebral*  
610 *Cortex*. 2015; 25(7):1697–1706. doi: 10.1093/cercor/bht355.
- 611 **Pérez A**, Carreiras M, Duñabeitia JA. Brain-To-brain entrainment: EEG interbrain synchronization while speak-  
612 ing and listening. *Scientific Reports*. 2017; 7(1):1–12. doi: 10.1038/s41598-017-04464-4.
- 613 **Pérez A**, Dumas G, Karadag M, Duñabeitia JA. Differential brain-to-brain entrainment while speaking and list-  
614 ening in native and foreign languages. *Cortex*. 2019; 111:303–315. doi: 10.1016/j.cortex.2018.11.026.
- 615 **Power AJ**, Foxe JJ, Forde EJ, Reilly RB, Lalor EC. At what time is the cocktail party? A late locus of selective  
616 attention to natural speech. *European Journal of Neuroscience*. 2012; 35(9):1497–1503. doi: 10.1111/j.1460-  
617 9568.2012.08060.x.
- 618 **Power AJ**, Lalor EC, Reilly RB. Endogenous auditory spatial attention modulates obligatory sensory activity in  
619 auditory cortex. *Cerebral Cortex*. 2011; 21(6):1223–1230. doi: 10.1093/cercor/bhq233.
- 620 **Pulvermüller F**, Fadiga L. Active perception: Sensorimotor circuits as a cortical basis for language. *Nature*  
621 *Reviews Neuroscience*. 2010; 11(5):351–360. <http://dx.doi.org/10.1038/nrn2811>, doi: 10.1038/nrn2811.
- 622 **Redcay E**, Schilbach L. Using second-person neuroscience to elucidate the mechanisms of social interac-  
623 tion. *Nature Reviews Neuroscience*. 2019; 20(8):495–505. <http://dx.doi.org/10.1038/s41583-019-0179-4>, doi:  
624 10.1038/s41583-019-0179-4.
- 625 **Sanders LD**, Neville HJ. An ERP study of continuous speech processing. *Cognitive Brain Research*. 2003;  
626 15(3):228–240. doi: 10.1016/s0926-6410(02)00195-7.
- 627 **Schrimpf M**, Blank IA, Tuckute G, Kauf C, Hosseini EA, Kanwisher N, Tenenbaum JB, Fedorenko E. The neural  
628 architecture of language: Integrative modeling converges on predictive processing. *Proceedings of the Na-*  
629 *tional Academy of Sciences*. 2021 nov; 118(45):e2105646118. [http://www.pnas.org/lookup/doi/10.1073/pnas.](http://www.pnas.org/lookup/doi/10.1073/pnas.2105646118)  
630 [2105646118](http://www.pnas.org/lookup/doi/10.1073/pnas.2105646118), doi: 10.1073/pnas.2105646118.
- 631 **Shamma SA**, Elhilali M, Micheyl C. Temporal coherence and attention in auditory scene analysis. *Trends in Neu-*  
632 *rosiences*. 2011; 34(3):114–123. <http://dx.doi.org/10.1016/j.tins.2010.11.002>, doi: 10.1016/j.tins.2010.11.002.
- 633 **Shinn-Cunningham BG**. Object-based auditory and visual attention. *Trends in Cognitive Sciences*. 2008;  
634 12(5):182–186. doi: 10.1016/j.tics.2008.02.003.
- 635 **Sonkusare S**, Breakspear M, Guo C. Naturalistic Stimuli in Neuroscience: Critically Acclaimed. *Trends in Cogni-*  
636 *tive Sciences*. 2019 aug; 23(8):699–714. <https://linkinghub.elsevier.com/retrieve/pii/S1364661319301275>, doi:  
637 10.1016/j.tics.2019.05.004.
- 638 **Stephens GJ**, Silbert LJ, Hasson U. Speaker–listener neural coupling underlies successful commu-  
639 nication. *Proceedings of the National Academy of Sciences*. 2010 aug; 107(32):14425–14430.  
640 <http://www.pnas.org/content/107/32/14425><http://www.pnas.org/cgi/doi/10.1073/pnas.1008662107><https://pnas.org/doi/full/10.1073/pnas.1008662107>, doi: 10.1073/pnas.1008662107.
- 642 **Stolk A**, Verhagen L, Toni I. Conceptual Alignment: How Brains Achieve Mutual Understanding.  
643 *Trends in Cognitive Sciences*. 2016; 20(3):180–191. <http://dx.doi.org/10.1016/j.tics.2015.11.007>, doi:  
644 10.1016/j.tics.2015.11.007.
- 645 **Teoh ES**, Ahmed F, Lalor EC. Attention Differentially Affects Acoustic and Phonetic Feature Encoding in a Mul-  
646 tispeaker Environment. *The Journal of Neuroscience*. 2022 jan; 42(4):682–691. [https://www.jneurosci.org/](https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.1455-20.2021)  
647 [lookup/doi/10.1523/JNEUROSCI.1455-20.2021](https://www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.1455-20.2021), doi: 10.1523/JNEUROSCI.1455-20.2021.
- 648 **Vaswani A**, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I. Attention is all you  
649 need. *Advances in Neural Information Processing Systems*. 2017; 2017-Decem(Nips):5999–6009.
- 650 **Wang Y**, Ding N, Ahma N, Xiang J, Poeppel D, Simon JZ. Sensitivity to temporal modulation rate and spectral  
651 bandwidth in the human auditory system: MEG evidence. *Journal of Neurophysiology*. 2012; 107(8):2033–  
652 2041. doi: 10.1152/jn.00310.2011.

- 653 **Wang Y**, Zhang J, Zou J, Luo H, Ding N. Prior Knowledge Guides Speech Segregation in Human Auditory Cortex.  
654 *Cerebral Cortex*. 2019; 29(4):1561–1571. doi: 10.1093/cercor/bhy052.
- 655 **Weissbart H**, Kandylaki KD, Reichenbach T. Cortical tracking of surprisal during continuous speech compre-  
656 hension. *Journal of Cognitive Neuroscience*. 2019; 32(1):155–166. doi: 10.1162/jocn\_a\_01467.
- 657 **Willems RM**, Jacobs AM. Caring About Dostoyevsky: The Untapped Potential of Studying Literature.  
658 *Trends in Cognitive Sciences*. 2016; 20(4):243–245. <http://dx.doi.org/10.1016/j.tics.2015.12.009>, doi:  
659 [10.1016/j.tics.2015.12.009](https://doi.org/10.1016/j.tics.2015.12.009).
- 660 **Willems RM**, Nastase SA, Milivojevic B. Narratives for Neuroscience. *Trends in Neurosciences*. 2020  
661 may; 43(5):271–273. <https://doi.org/10.1016/j.tins.2020.03.003>[https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0166223620300497)  
662 [S0166223620300497](https://doi.org/10.1016/j.tins.2020.03.003), doi: 10.1016/j.tins.2020.03.003.
- 663 **Yeshurun Y**, Nguyen M, Hasson U. The default mode network: where the idiosyncratic self meets the  
664 shared social world. *Nature Reviews Neuroscience*. 2021 mar; 22(3):181–192. [http://dx.doi.org/10.1038/](http://dx.doi.org/10.1038/s41583-020-00420-w)  
665 [s41583-020-00420-w](https://doi.org/10.1038/s41583-020-00420-w)<http://www.nature.com/articles/s41583-020-00420-w>, doi: 10.1038/s41583-020-00420-w.
- 666 **Yu Q**, Bi Z, Jiang S, Yan B, Chen H, Wang Y, Miao Y, Li K, Wei Z, Xie Y, Tan X, Liu X, Fu H, Cui L, Xing L, Weng S, Wang  
667 X, Yuan Y, Zhou C, Wang G, et al. Visual cortex encodes timing information in humans and mice. *Neuron*.  
668 2022 oct; p. 1–18. <https://doi.org/10.1016/j.neuron.2022.09.008>[https://linkinghub.elsevier.com/retrieve/pii/](https://linkinghub.elsevier.com/retrieve/pii/S0896627322008133)  
669 [S0896627322008133](https://doi.org/10.1016/j.neuron.2022.09.008), doi: 10.1016/j.neuron.2022.09.008.
- 670 **Zhang ZG**, Hu L, Hung YS, Mouraux A, Iannetti GD. Gamma-band oscillations in the primary somatosen-  
671 sory cortex-A direct and obligatory correlate of subjective pain intensity. *Journal of Neuroscience*. 2012;  
672 32(22):7429–7438. doi: 10.1523/JNEUROSCI.5877-11.2012.
- 673 **Zion Golumbic E**, Cogan GB, Schroeder CE, Poeppel D. Visual input enhances selective speech envelope  
674 tracking in auditory cortex at a "cocktail party". *The Journal of neuroscience*. 2013; 33(4):1417–1426.  
675 <http://www.jneurosci.org/content/33/4/1417.full>, doi: 10.1523/JNEUROSCI.3675-12.2013.

**Appendix 0—table A1.** Details of Neural Language Processing Model.

Type	Parameter
model type	LSTM
embedding size	200
hidden units per layer	200
number of layers	2
initial learning rate	3
gradient clipping	0.25
sequence length	35
drop out	0.2
epoch	50
batch size	3