# Amplified cortical tracking of word-level features of continuous competing speech in older adults

Juraj Mesik[1*], Lucia A. Ray[1], and Magdalena Wojtczak[1]

[1] Department of Psychology, University of Minnesota, Minneapolis, MN, USA

**\* Correspondence:**
Juraj Mesik
mesik002@umn.edu

**Abstract**

Speech-in-noise comprehension difficulties are common among the elderly population, yet traditional objective measures of speech perception are largely insensitive to this deficit, particularly in the absence of clinical hearing loss. In recent years, a growing body of research in young normal-hearing adults has demonstrated that high-level features related to speech semantics and lexical predictability elicit strong centro-parietal negativity in the EEG signal around 400 ms following the word onset. Here we investigate effects of age on cortical tracking of these word-level features within a two-talker speech mixture, and their relationship with self-reported difficulties with speech-in-noise understanding. While undergoing EEG recordings, younger and older adult participants listened to a continuous narrative story in the presence of a distractor story. We then utilized forward encoding models to estimate cortical tracking of three speech features: 1) "semantic" dissimilarity of each word relative to the preceding context, 2) lexical surprisal for each word, and 3) overall word audibility. Our results revealed robust tracking of all three features for attended speech, with surprisal and word audibility showing significantly stronger contributions to neural activity than dissimilarity. Additionally, older adults exhibited significantly stronger tracking of surprisal and audibility than younger adults, especially over frontal electrode sites, potentially reflecting increased listening effort. Finally, neuro-behavioral analyses revealed trends of a negative relationship between subjective speech-in-noise perception difficulties and the model goodness-of-fit for attended speech, as well as a positive relationship between task performance and the goodness-of-fit, indicating behavioral relevance of these measures. Together, our results demonstrate the utility of modeling cortical responses to multi-talker speech using complex, word-level features and the potential for their use to study changes in speech processing due to aging and hearing loss.

**Keywords**: speech perception; aging; speech-in-noise; electroencephalography; temporal response function; lexical surprisal; semantic processing

**1. Introduction**

42  Speech perception is fundamentally important for human communication. While speech signals
43  are often embedded in complex sound mixtures that can interfere with speech perception via
44  energetic and informational masking, the auditory system is remarkably adept at utilizing
45  attentional mechanisms to suppress distractor information and enhance representations of the
46  target speech (e.g., Ding and Simon, 2012a; Mesgarani and Chang, 2012; O'Sullivan et al.,
47  2019). However, the robustness of speech perception, particularly in the presence of noise, is
48  vulnerable to deterioration through both noise-induced and age-related hearing loss (Dubno et
49  al., 1984; Helfer and Wilber, 1990; Fogerty et al., 2015, 2020) as well as age-related cognitive
50  decline (van Rooij and Plomp, 1990; Akeroyd, 2008; Dryden et al., 2017). Additionally, a small
51  but significant portion of the population experiences speech-in-noise (SIN) perception
52  difficulties, without exhibiting clinical hearing loss (Saunders, 1989; Zhao and Stephens, 2007;
53  Tremblay et al., 2015). Together, these SIN perception difficulties can lead to significant
54  impairment in quality of life (Dalton et al., 2003; Chia et al., 2007), and in older adults they may
55  result in increased social isolation (Chia et al., 2007; Mick et al., 2014; Pronk et al., 2014),
56  potentially exacerbating loss of cognitive function (Loughrey et al., 2018; Ray et al., 2018).
57  Although subjective SIN perception difficulties are relatively common in older
58  individuals, objective tests for quantifying these deficits, such as identification of words or
59  sentences in noise (e.g., QuickSin; Killion et al., 2004), often do not strongly correlate with the
60  degree of subjective deficit (Phatak et al., 2018), particularly in cases with little-to-no clinical
61  hearing loss. Smith and colleagues (2019) recently reported that only 8% of their sample of 194
62  listeners exhibited deficits in objective SIN tasks, while 42% of listeners indicated experiencing
63  subjective SIN perception difficulties. A likely reason for this mismatch is that objective speech
64  perception tests do not accurately reflect real world scenarios where SIN difficulties arise. For
65  example, while existing tests generally require identification of isolated words or sentences
66  embedded in noise (e.g., speech-shaped noise or a competing talker), real world speech
67  perception often requires real-time comprehension of multi-sentence expressions, embedded
68  in a reverberant environment, in the presence of multiple competing speakers at different
69  spatial positions. In these scenarios, listeners who need to expend additional time and cognitive
70  resources to identify the meaning of the incoming speech may "fall behind" in comprehension
71  of later parts of the utterance. Moreover, even if the listener can correctly piece together the
72  meaning of the utterance, their subjective confidence may be diminished, potentially "blurring"
73  the predictive processes thought to facilitate perception of upcoming speech (Pickering and
74  Gambi, 2018). As such, behavioral measures that more accurately reflect subjective SIN
75  perception difficulties may require utilization of more realistic, narrative stimuli, and focus on
76  quantifying comprehension, as opposed to simple word or sentence identification (e.g., Xia et
77  al., 2017).
78  While development of behavioral paradigms focusing on characterizing SIN perception
79  difficulties is an important goal, a complementary and potentially more sensitive approach to
80  quantifying these deficits may be provided by neural measures of continuous-speech tracking.
81  In recent years, non-invasive methodologies for measurement of neural representations of
82  continuous speech in humans have become increasingly popular (Lalor and Foxe, 2010; Crosse

83    et al., 2016), particularly in application to young normal-hearing (YNH) populations. One
84    important result of this work has been the demonstration of profound attentional modulation
85    of speech whereby temporal dynamics of neural responses to attended and ignored speech
86    differ considerably, both in representation of lower-level features such as the speech envelope
87    (Ding and Simon, 2012; Power et al., 2012; Kong et al., 2014; Fiedler et al., 2019), and higher-
88    level features related to lexical and semantic content of speech (Brodbeck et al., 2018;
89    Broderick et al., 2018). Indeed, while lower-level features produce robust responses even when
90    speech is ignored, features related to linguistic representations only show robust responses for
91    attended speech, suggesting that they are tightly linked with speech comprehension.
92    Responses to higher-level features may therefore be particularly sensitive to SIN perception
93    difficulties, which are likely associated with impaired comprehension performance. In fact, SIN
94    perception difficulties could potentially manifest themselves not only in terms of poorer
95    tracking of higher-level features in attended speech, but also in increased tracking of features in
96    ignored speech, when facing difficulties with suppression of distractor information.
97            Changes in neural processing of continuous speech in aging populations, compared to
98    young adults, are relatively poorly understood. Several studies have utilized magneto- and
99    electroencephalography (M/EEG) to address this question. Studies comparing envelope-related
100   cortical responses have revealed a pattern of amplified envelope representations in older
101   populations (Presacco et al., 2016; Decruy et al., 2019; Zan et al., 2020), potentially reflecting
102   changes in the utilization of cognitive resources during speech comprehension. More recently,
103   Broderick et al. (2020) compared higher-level representations of speech in younger and older
104   populations. They estimated EEG responses to 5-gram surprisal, reflecting the predictability of
105   words given the preceding sequence of four words, as well as semantic dissimilarity, reflecting
106   the contribution of each word to the semantic content of a sentence. While younger listeners
107   showed strong responses to both of these features, older adults exhibited a delayed surprisal
108   response and a near-absent response to semantic dissimilarity. These findings demonstrate
109   that representations of higher-level features of speech may indeed reveal robust effects of age.
110   However, because Broderick et al. (2020) did not report behavioral measures related to speech
111   comprehension, nor measures of subjective speech perception difficulties among their
112   participants, it is unclear whether these metrics would correlate with the reported EEG-based
113   findings. Moreover, participants in that study were presented with clear speech without any
114   distractors (e.g., competing speakers), making it unclear how speech representations differ in
115   complex listening scenarios where speech perception difficulties are most commonly reported.
116           The goal of this study was to compare higher-level neural representations of two-talker
117   speech mixtures between younger and older adults, and to explore how these measures relate
118   to comprehension performance and self-reported SIN perception difficulties. In particular, we
119   examined representations related to word dissimilarity relative to short-term preceding
120   context, lexical surprisal based on multi-sentence context, and word-level audibility. We chose
121   to pursue this paradigm for several reasons. First, a multi-talker paradigm was chosen because
122   subjective SIN perception difficulties commonly arise in aging listeners in the context of
123   competing speech. If age-related changes in neural representations are confirmed, then these

124 neural signatures could potentially be further explored as a candidate objective correlate for
125 subjective SIN difficulties. Second, we chose to characterize responses to word-level features
126 linked to meaning and lexical predictability because existing evidence indicates that responses
127 to higher-level features are tightly linked to speech comprehension (Broderick et al., 2018). As
128 such, we anticipated that responses to these features are more likely to exhibit differences as a
129 function of age and SIN perception difficulties. Although neural representations reflecting the
130 end-goal of speech perception may allow for only limited inference about the underlying causes
131 of SIN perception difficulties, which can range from peripheral changes in acoustic
132 representations to more central changes in cognitive processes, these representations may
133 offer increased sensitivity due to capturing the combined effects of the various etiologies
134 underlying the deficit.
135
136 **2. Materials and Methods**
137
138 **2.1 Participants**
139
140 In total, 45 adult volunteers completed the experiment, and data from 41 participants were
141 used due to a methodological change implemented early in data collection. The participant
142 pool was divided into two groups, younger adults (YA) and older adults (OA), with participants
143 who were 18-39 years included in the former, and participants who were 40-70 years included
144 in the latter. The YA group consisted of 20 participants (6 male, 14 female; mean ± s.d. age:
145 29.40 ± 6.40 years), while the OA group included 21 participants (9 male, 12 female; mean ±
146 s.d. age: 53.48 ± 8.68 years). Participants were recruited via email advertisement from a pool of
147 students, staff, and alumni of the University of Minnesota. All participants provided informed
148 written consent and received either course credit or monetary compensation for their
149 participation. The procedures were approved by the Institutional Review Board of the
150 University of Minnesota.
151
152 **2.2 Audiometry**
153
154 An air-conduction audiogram was measured in each ear for each participant prior to beginning
155 the EEG procedures. Detection thresholds were measured at octave frequencies in the 250 –
156 8000 Hz range, and frequencies for which thresholds exceeded 20 dB HL were deemed to be
157 affected by hearing loss (HL). This procedure resulted in the detection of 2 participants in the
158 YA group, and 16 participants in the OA group as having mild-to-moderate high-frequency HL.
159 The skewed distribution of HL towards the older population was expected, as peripheral
160 frequency sensitivity naturally diminishes with age (see reviews by Huang and Tang, 2010;
161 Yamasoba et al., 2013).
162

163 For participants with any hearing loss, all experimental audio materials were amplified in the
164 frequency regions of hearing loss, as described in section 2.4 below. Under these conditions, we
165 observed no association between task performance and high-frequency hearing loss.
166
167 **2.3 Modified SSQ questionnaire**
168
169 Prior to the EEG procedures, all participants completed a modified version of a subset of
170 Speech, Spatial and Qualities of Hearing Scale ($SSQ_m$). The original version of SSQ (Gatehouse
171 and Noble, 2004) was designed to measure subjective hearing challenges faced by listeners in
172 various situations of daily life. In our version, we specifically probed participants about
173 difficulties with and frustrations related to hearing speech in noisy situations, such as cafes and
174 social gatherings. Each of the 14 items was presented on a computer screen along with four
175 graded choices of frequency, difficulty, or discomfort related to the presented listening
176 scenarios. E.g.,
177
178 Item 1:
179 I find it difficult to talk with staff in places such as shops, cafes, or banks, due to struggling to
180 hear what they are saying.
181
182 Item 10:
183 In group conversations I worry about mishearing people and responding based on incorrect
184 information.
185
186 Response choices:
187     1) Not at all
188     2) Rarely
189     3) Often
190     4) Very often
191
192
193 **2.4 Stimuli**
194
195 Stimuli were four public domain short story audiobooks (*Summer Snow Storm* by Adam Chase;
196 *Mr. Tilly's Seance* by Edward F. Benson; *A Pail of Air* by Fritz Leiber; *Home Is Where You Left It*
197 by Adam Chase; source: LibriVox.org), spoken by two male speakers (two stories per speaker).
198 Each story was about 25 min in duration and was pre-processed to truncate any silences
199 between words that exceeded a 500-ms interval to 500 ms. On a block-by-block basis (see
200 section 2.5 below), each audiobook was root-mean-square (RMS) normalized and scaled to 65
201 dB SPL. Stimuli were presented to participants using ER1 Insert Earphones (Etymotic Research,
202 Elk Grove Village, IL), shielded with copper foil to prevent electrical artifacts in the EEG data.
203

204    In order to minimize the odds of finding age-related differences in neural responses that could
205    be attributed to reduced audibility in participants with hearing loss, all audio materials were
206    custom-filtered for each participant with HL using a FIR filter implemented in MATLAB
207    (Mathworks, Natick, MA) via the *designfilt* and *filter* functions. The filter was designed to apply
208    half gain, amplifying all frequency bands by half the amount of the hearing loss:

209

210    $A(f) = 0.5 \times (T(f) - 20)$          when T(*f*) > 20 dB HL,
211    $A(f) = 0$                              otherwise,

212

213    where T(*f*) is the detection threshold in dB HL at frequency *f*. Note that half gain amplification is
214    a commonly used strategy to mitigate reduced audibility due to hearing loss, while preventing
215    discomfort from loudness recruitment, whereby loudness growth for frequencies affected by
216    cochlear hearing loss is steeper than that observed in normal hearing (Fowler, 1936; Steinberg
217    and Gardner, 1937).

218

219    **2.5 Experimental procedures**

220

221    The experimental setup was implemented using the Psychophysics Toolbox (Brainard, 1997;
222    Pelli, 1997; Kleiner et al., 2007) in MATLAB. Two experimental runs were completed by each
223    study participant. In each run, a pair of audiobooks read by different male speakers (Fig. 1A)
224    was presented diotically (the mixture of the two audiobooks in each ear) to the participant. One
225    of the stories served as the *attended* story, while the other was the *ignored* story, with these
226    designations being counter-balanced across participants. A run was broken up into 24-27 blocks
227    (variation was due to small differences in durations of audiobooks used in each of the two
228    runs). Each block contained a roughly 1-minute segment of audio, followed by a series of
229    questions, detailed below. Block duration was allowed to exceed 1 minute in order to ensure
230    that each block concluded at the end of a sentence in the attended story. The attended story
231    remained the same throughout the run. To cue the participants to follow the correct story, the
232    audio of the attended story started 1 sec prior to the onset of the ignored story. This was
233    further aided by making this initial 1-sec portion of the attended story in each block (except
234    block #1) correspond to the final 1-sec of the attended story from the previous block. These
235    repeated segments with the attended story alone were excluded from statistical analyses.
236    Throughout each block, participants were instructed to stay as still as possible, and to keep
237    their gaze on a central fixation marker presented on a computer display in front of the
238    participant. The purpose of this was to minimize EEG artifacts caused by muscle activity.
239         Following each block, participants were presented on a display with a series of Yes/No
240    questions about the audio from that block, including:

241

242    1) Four comprehension questions about the contents of the attended story
243    2) Confidence ratings for each of the comprehension questions
244    3) Intelligibility judgment about the attended speaker
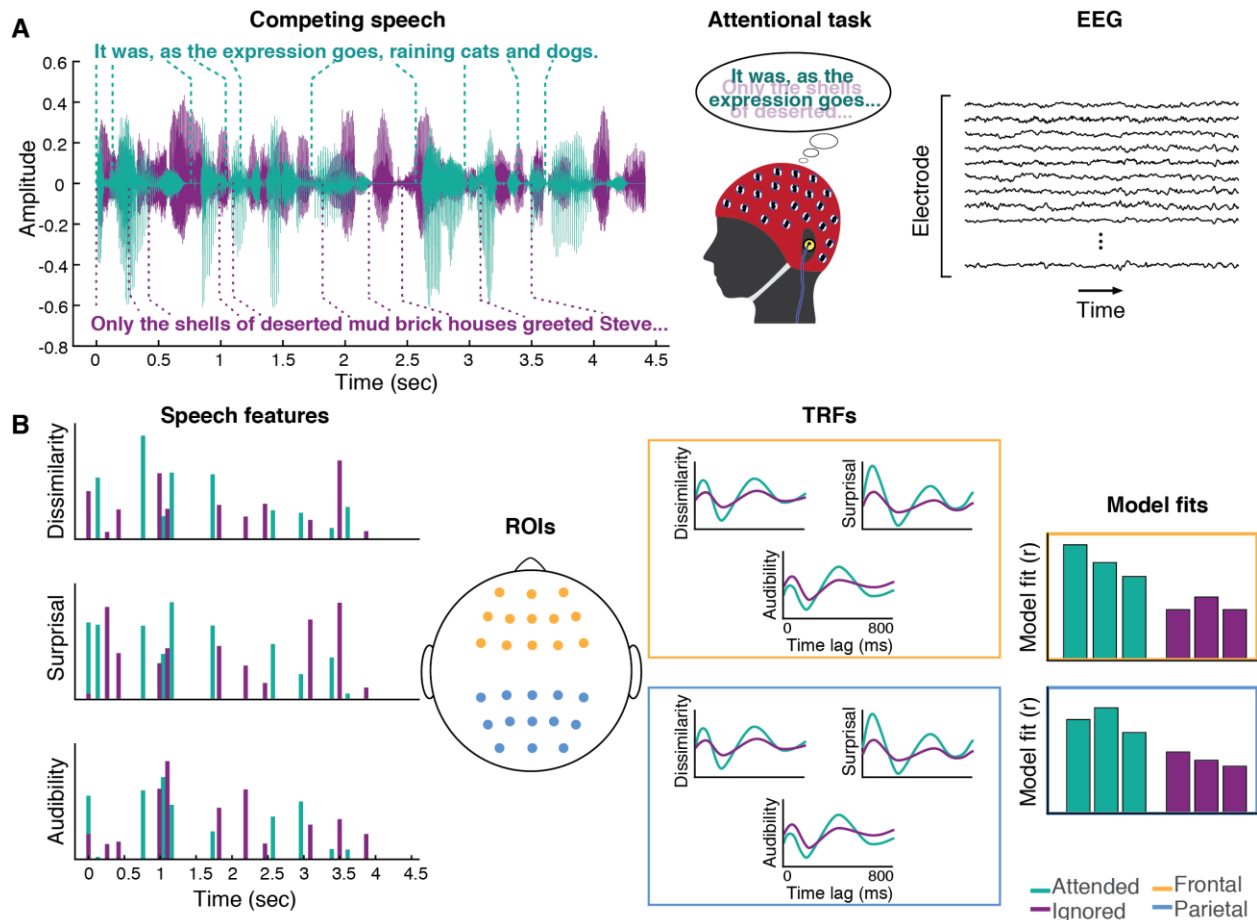
245    4) Subjective attentiveness rating

246

247        As each behavioral question had binary answer choices (e.g., for attentiveness,
248    participants answered "Were you able to stay focused on the target story?" Yes/No), the main
249    purpose of these questions was to gather information about participants' comprehension and
250    subjective experience throughout the run, and to make sure that they were attending to the
251    correct story.
252        Participants were given 10 seconds to answer each question using a key press. If 10
253    seconds elapsed without a response, the question was marked as no-response. After answering
254    each block's questions, participants were allowed to request a short break to ensure that they
255    remained comfortable throughout the experiment. These breaks were limited to up to two
256    minutes, during which participants remained seated. The next block started as soon as the
257    break was terminated by the participant with a key press, or two minutes elapsed.
258    Furthermore, between the two experimental runs, participants were offered an extended break
259    inside the booth. The EEG cap and the insert phones were not removed during the breaks.
260        The second experimental run was procedurally identical to the first one, except a
261    different pair of stories was presented, neither of which was used in the first run. Additionally,
262    the attended and ignored speakers were switched, so that the speaker that narrated the
263    ignored story in the first run was attended in the second run, while the attended speaker from
264    the first run became the ignored speaker in the second run. Participants were explicitly
265    informed of this switch, and the purpose of this was to balance any possible speaker effects on
266    each participant's EEG data.

267

Figure 1. Experimental procedures. (A) Participants listened to a mixture of two speakers, while attending to one of them. Meanwhile, 64-channel EEG was recorded from their scalp. (B) Three word-level features (dissimilarity, surprisal, and audibility) were extracted from the speech for both the attended and ignored stories, and used to generate regressors containing impulses that were time-aligned to the word onsets scaled by the amplitude of each feature. These features were regressed against the EEG signals recorded during the experiment, resulting in TRF and model fit contributions for each of the features. These TRFs and goodness-of-fit values were averaged across groups of frontal (yellow) and parietal (blue) electrodes for use in group-level analyses.

## 2.6 EEG procedures

While engaging in the experimental task described above, each participant's EEG activity was sampled at 4096 Hz from their scalp using a Biosemi ActiveTwo system (BioSemi B.V., Amsterdam, The Netherlands), with 64 channels positioned according to the international 10-20 system (Klem et al., 1999). Additional external electrodes were placed on the left and right mastoids, and above and below the right eye (vertical electro-oculogram, VEOG). Prior to the beginning of the recording, and between the two runs, the experimenter visually inspected

287  signals in all electrodes, and for any electrodes with DC offsets exceeding ± 20 mV, the contact
288  between the electrode and scalp was readjusted until the offset fell below ± 20 mV.
289
290  **2.7 EEG preprocessing**
291  All pre-processing analyses were implemented via the EEGLAB toolbox (Delorme and Makeig,
292  2004) for MATLAB, unless otherwise stated. To reduce computational load, the raw EEG data
293  were initially downsampled to 256 Hz, and band-pass filtered between 1 and 80 Hz using a
294  Hamming windowed sinc FIR filter implemented in the *pop_eegfiltnew* function of EEGLAB.
295  Subsequently, data were pre-processed using the PREP pipeline (Bigdely-Shamlo et al., 2015).
296  These steps included line noise removal, detection of disproportionately noisy channels via an
297  iterative robust referencing procedure, interpolation of noisy channels, and referencing the
298  data using the final "clean" estimate of the global mean activation. The benefit of this
299  procedure is that it minimizes the risk of signal contamination from electrodes with abnormal
300  signals (e.g., due to faulty hardware) during the referencing stage.
301       Next, activations from all experimental blocks were epoched and independent
302  component analysis (ICA; Jutten and Herault, 1991; Comon, 1994) was applied to the data using
303  the infomax ICA algorithm (Bell and Sejnowski, 1995) implementation in EEGLAB. This
304  procedure decomposes the EEG signal into statistically independent sources of activation, some
305  of which reflect sensory and cognitive processes, while others capture muscle-related signal
306  contributions and other sources of noise. We removed all components that matched eye-blink
307  related activity in component topography, amplitude, and temporal characteristics, as well as
308  other high-amplitude artifacts that reflected muscle activity. This, on average, led to the
309  removal of 2.52 (SD: 0.97) components.
310       The cleaned EEG signals were then band-pass filtered between 1 and 8 Hz with a
311  Chebyshev type 2 filter designed using MATLAB's *designfilt* function (optimized to achieve 80
312  dB attenuation below 0.5 Hz and above 9 Hz, with pass-band ripple of 1 dB), and applied to the
313  data using the *filtfilt* function. Afterwards, the data were z-scored in order to control for inter-
314  subject variability in the overall signal amplitude due to nuisance factors such as skull thickness
315  or scalp conductivity, as well as to improve efficiency in the cross-validated regression and ridge
316  parameter search for deriving the temporal response function (TRF), described below (section
317  2.9.1). Finally, because run duration varied slightly due to unequal lengths of the two pairs of
318  audiobooks (i.e. 24-27 minutes), in order to equalize contributions from each run to the overall
319  analysis results, only blocks 2-23 from each run were used in the remaining analyses. The first
320  block was excluded in order to minimize effects of initial errors in attending to the target story,
321  which happened to a very small number of participants (less than 5), but was quickly corrected
322  after initial comprehension questions were presented.
323
324  **2.8 Word timing estimation**
325  Word onset timings for all words within each story were estimated using the Montreal Forced
326  Aligner (McAuliffe et al., 2017). Prior to running the aligner, the audiobook text was
327  preprocessed to remove punctuation, typographic errors and abbreviations, and both the text

328 and audio were divided into roughly 30-sec segments. This segmented alignment approach was
329 used in order to prevent accumulation of alignment errors for later portions of the audio. All
330 alignments were subsequently manually inspected for timing errors, and when noticeable
331 alignment errors were detected, the aligner was re-run on further-shortened (15 sec) segments
332 of the affected audio. While forced alignment routinely results in some degree of timing errors,
333 these are typically small, with a median of about 15 ms for the aligner used here. As such, only
334 a small degree of temporal smearing of estimated neural responses should occur due to these
335 errors.
336
337 **2.9 Data analysis**
338
339 **2.9.1 TRF analyses**
340 Time courses of cortical responses to different speech features, known as the TRFs, were
341 extracted from preprocessed EEG activity using cross-validated regularized linear regression,
342 implemented via the mTRF toolbox (Crosse et al., 2016). Briefly, deconvolution of a TRF for a
343 given feature from the EEG signal is accomplished by first constructing a regressor containing a
344 time series, sampled at a rate matching the EEG signal, of that feature's amplitudes. By
345 including multiple time-lagged copies of the regressor for each feature, the effect of a given
346 feature on the neural activity at different latencies relative to the word onset can be estimated,
347 resulting in a time course of neural response. Regressors for all features are combined into a
348 full design matrix, and this matrix is then regressed against the EEG signal to yield the impulse
349 responses (i.e., TRFs) for each of the included features at each electrode site.
350      In practice, this procedure was implemented through 11-fold cross-validation, with each
351 fold involving three steps. First, the data and regressors were split into a training set, composed
352 of 40 blocks of the data (~40 minutes), and a testing set, containing the remaining 4 blocks of
353 the data (~4 minutes). Next, the training set was used to determine the ridge parameter, λ, by
354 iteratively fitting the cortical-response model using a range of ridge parameters. The TRF
355 estimates were obtained for the λ parameter that produced the best model fit to the training
356 data, as determined by the highest Pearson's correlation coefficient between the predicted and
357 actual EEG signal. The TRF estimates were then used to assess the model fit for the test data.
358 This was done by convolving the estimated TRFs with the corresponding word-feature
359 regressors for the test data set, and computing the Pearson's correlation between the
360 predicted and actual test data. Following cross-validation, average TRFs for each feature and an
361 average model goodness-of-fit were computed from results of all cross-validation folds for use
362 in group-level analyses.
363
364 **2.9.1.1 Regression features**
365 Word features used in the regression analyses included semantic dissimilarity, surprisal, and
366 word audibility (Fig. 1B).
367
368 **2.9.1.1.1 Semantic dissimilarity**

369 Semantic dissimilarity, reflecting approximately the degree to which each word adds new
370 information to a sentence, was computed as described in Broderick et al., (2018). Briefly, we
371 used Google's pre-trained *word2vec* neural network (Mikolov et al., 2013a, 2013b),
372 implemented using the Gensim library (Rehurek and Sojka, 2010) for Python, to compute a 300-
373 dimensional vector representation (otherwise known as an embedding) of each word within
374 our stimuli. An important property of these vector representations is that in the 300-
375 dimensional vector space, vectors of words with similar meanings point in similar directions.
376 Computing correlation between vectors representing any two words approximates their
377 semantic similarity. Because EEG response to incongruent words has been shown to elicit a
378 strong N400 component (Kutas and Hillyard, 1980), for regression purposes these similarity
379 values were subtracted from 1 to convert them to dissimilarity.
380 To construct semantic dissimilarity regressors, we computed the dissimilarity between
381 each word's vector, and the average of vectors for all preceding words in a given sentence. In
382 the case of the first word in a sentence, we computed dissimilarity from the average vector for
383 words in the previous sentence. These dissimilarity values were then used to construct the
384 regressor consisting of unit-length impulses aligned to word onsets that were scaled by each
385 word's dissimilarity value and zeros between these impulses. Although neural responses to
386 semantic content of words may not be strictly time-locked to word onsets, potentially leading
387 to some degree of temporal smearing in the estimated TRFs, word onset timings have been
388 successfully used as timestamps for characterizing higher-order lexical and semantic processes
389 (e.g., Broderick et al., 2018; Weissbart et al., 2019).
390

391 **2.9.1.1.2 Lexical surprisal**
392 Surprisal regressors were constructed in an identical way to dissimilarity, except the feature
393 values were computed using OpenAI's GPT-2 (Radford et al., 2019; 12-layer, 117M parameter
394 version) artificial neural network (ANN), similar to the approach demonstrated by Heilbron et
395 al. (2019). These procedures were implemented in Python using the Transformers library (Wolf
396 et al., 2020) for PyTorch (Paszke et al., 2019). GPT-2 is a transformer-based (Vaswani et al.,
397 2017) ANN that, using a "self-attention" mechanism, is capable of effectively using hundreds of
398 words worth of preceding context in order to generate seemingly realistic sequences of text. As
399 a result, it can be used as a proxy for computing the predictability of words within a sequence.
400 Surprisal is calculated based on a much longer time scale (a large number of words in the
401 preceding context) than semantic dissimilarity. Specifically, by providing GPT-2 with a segment
402 of text and then generating the distribution over the next word, it is possible to assess the
403 relative probability of the actual next word within GPT-2's distribution of possibilities.
404 Generation of all probabilities involves iteratively adding words into the context, and computing
405 the probability of each successive word. In practice, GPT-2 utilizes a tokenized representation of
406 text, whereby GPT-2's vocabulary corresponds to a combination of whole words (particularly in
407 the case of shorter words) and word fragments.
408 As a result, the probability of the i-th word $w_i$ was computed as a product of conditional
409 probabilities of the constituent word tokens $t$, with each token's probability being computed

410 with the model's knowledge of the preceding tokens (i.e. preceding text plus current word's
411 tokens whose probabilities were already estimated):

412

413
$$p(w_i) = \prod_{j=1}^{n} p\big(t_{k+j} \,\big|\, t_{k+j-512}, \dots t_{k+j-1}\big),$$

414

415 where $j$ indexes the $n$ tokens of word $w_i$, $k$ is the absolute index of the last token in the
416 preceding word (relative to text beginning), and 512 is the maximum number of tokens utilized
417 for prediction. For token indices less than 512 (i.e., early portions of the text), all of the
418 available context was used. Furthermore, in cases where one or more tokens from the word at
419 the far boundary of the context window did not fit into the 512 token limit, that word's tokens
420 were excluded from being used for prediction. Note that although GPT-2 is capable of utilizing
421 up to 1024 tokens for prediction, we utilized a context length of 512 tokens due to limited
422 computational resources. Across the 4 stories, when full predictive context was utilized for
423 prediction, it contained on average 393.3 [s.d. = 31.1] words.

424 Because brain mechanisms underlying lexical prediction respond more to unexpected
425 than to expected words (Kutas and Hillyard, 1984), surprisal was computed by taking the
426 negative log of the conditional probabilities of each word, leading to less expected words
427 receiving higher surprisal values:

428

429 $S(w_i) = -\log(p(w_i))$

430

431 **2.9.1.1.3 Audibility**
432 Word audibility regressors were constructed separately for the attended and ignored stories to
433 capture the degree of masking of each word in one story by the speaker of the other story. In
434 contrast to dissimilarity and surprisal, this value reflects the information at the shortest, word-
435 by-word time scale, with higher signal-to-noise ratio (SNR) values reflecting greater peripheral
436 fidelity of target speech, leading to lower uncertainty in speech identification on the basis of
437 the bottom-up signal. For each word $w_i$ in a given story, its audibility was defined in dB SNR
438 units:

439

440
$$Aud(w_i) = 20 \log \frac{RMS(y(w_i))}{RMS(z(w_i))},$$

441

442 where y($w_i$) is the acoustic waveform of a word $w_i$ spoken by one speaker, and z($w_i$) is the
443 acoustic waveform of the other speaker at the same time. Because neural responses have
444 limited dynamic range while the audibility measure ranged from –inf to inf, the audibility values
445 were rescaled to range from 0 to 1. In order to do this, audibility values were first clipped above
446 10 dB and below -10 dB, and then scaled to the 0-1 range by:

447

448
$$Aud_{scaled} = \frac{Aud + 10}{20}$$

449

450       Finally, because the distributions of regressor values had distinct means for different
451 features, we normalized each feature's non-zero regressor values to have an RMS of 1. Bringing
452 different features into similar amplitude ranges was done in order to make the amplitudes of
453 corresponding TRFs more similar to each other, thus improving regularization performance.

454       It is notable that although neither dissimilarity, nor surprisal correlated with audibility (r
455 = 0.03 and -0.02, respectively), there was a modest correlation between dissimilarity and
456 surprisal (r = 0.22), suggesting that both features captured some aspects of speech
457 predictability. Nevertheless, the fact that the correlation was relatively low suggests that much
458 of the variance in each of the two features captured distinct aspects of the linguistic content in
459 the speech stimuli.

460

461 **2.9.2 Feature-specific model performance**

462

463 After fitting the full three-feature model as described above, we computed the unique
464 contribution of each feature to the overall model fit using procedures described in Broderick et
465 al. (2020). Briefly, on each cross-validation fold, we estimated each feature's contribution to the
466 overall fit by comparing the goodness-of-fit for the full model to a null model, in which that
467 feature's contribution was eliminated. This was done by permuting regressor values of that
468 feature, while maintaining their original timing. For all other features, the original regressors
469 were used. Null model fits were computed by convolving the estimated TRFs with these
470 regressors and correlating the predicted EEG waveform with the test data. This procedure was
471 repeated 10 times to estimate the average null-model performance. Each feature's model
472 contribution was then computed as the difference between the goodness-of-fit metrics for the
473 full model and its null model.

474

475 **2.9.3 Regions of interest**

476

477 To strengthen our statistical analyses in light of inter-subject variability due to nuisance
478 variables such as head shape and electrode cap placement, all analyses were performed on two
479 regions of interest (ROI) derived by averaging model goodness-of-fit and TRFs from subsets of
480 frontal and parietal electrodes (Fig. 1B). The parietal ROI was chosen because of prior evidence
481 that responses to higher-level features such as dissimilarity or surprisal tend to peak over
482 parietal sites near electrode Pz (e.g., Broderick et al., 2018; Weissbart et al., 2019). The frontal
483 ROI was included because we hypothesized that recruitment of frontal regions may aid
484 prediction and disambiguation of the speech signals, particularly in challenging listening
485 scenarios such as in the presence of a competing speaker.

486

487

488 **2.9.4 Statistical analysis**

489

490  Group-level statistical analyses were applied to pooled outputs of single subject TRF analyses.
491  Prior to performing statistical tests, outliers were detected using a two-stage approach, applied
492  separately to samples from each age group to minimize the influence of true between-group
493  differences on this procedure. First, full model goodness-of-fit values that were more than 1.5
494  inter-quartile ranges (IQR) below the goodness-of-fit corresponding to the lower quartile, or 1.5
495  IQR above the value corresponding to the upper quartile were detected as outliers. No
496  participant met this criterion. Second, for each feature's TRF for the attended stories (which
497  were generally more robust compared to the ignored stories), we used the same 1.5 IQR
498  criterion to detect outliers at each time point of the TRF. Subsequently, we computed the
499  proportion of outlier time points for each subject. We set the outlier-proportion criterion to
500  0.15, so that participants with more than 15% of outlier time points were detected as outliers.
501  This led to the exclusion of 2 participants (1 YA, and 1 OA), leaving a total of 39 participants (19
502  YA and 20 OA, including 17 with HL) in the analysis.

503  A mixed-design ANOVA with a between-subjects factor of age group (YA vs. OA), and
504  within-subject factors of ROI (frontal vs. posterior), model feature (dissimilarity, surprisal, and
505  audibility), and attention (attended vs. ignored story) was used to assess how these factors
506  related to the feature-specific contributions to the model fit. Post-hoc tests were conducted
507  using two-tailed t-tests or the analogous non-parametric test, depending on the outcome of an
508  Anderson-Darling test of normality on the data.

509  Comparisons of TRFs for the attended and ignored stories were performed for each time
510  point of the TRFs using two-tailed, paired-samples t-tests. Because this involved hundreds of
511  statistical comparisons, we applied the *false discovery rate* (FDR; Benjamini and Hochberg,
512  1995) correction to control for the proportion of false positives among all significant
513  discoveries. Similarly, between-group comparisons (i.e., younger vs. older adults) were
514  performed on TRF time courses, with two-sample t-tests applied separately to the attended and
515  ignored TRFs and corrected using the FDR method.

516  Finally, exploratory correlation analyses were performed on different combinations of
517  neural (e.g., full model goodness-of-fit, feature-wise model contributions, TRF amplitudes) and
518  behavioral metrics (e.g., comprehension, confidence, and $SSQ_m$ scores). In these analyses we
519  corrected each set of correlations using the Bonferroni correction. Importantly, we used less
520  stringent multiple comparisons correction (i.e., not correcting by the total number of
521  comparisons across all combinations of correlated variables), because of the large number of
522  comparisons performed.

523
524  **3. Results**

525
526  **3.1 Behavioral measures of speech understanding**

527
528  Following each 1-minute block of listening to a two-talker speech mixture, participants
529  responded to four true/false questions about the content of the attended story and indicated
530  their confidence about their response. The average performance on this comprehension task

531 was 83.2% (SD: 6.8%, 65.9 - 94.2% range), significantly above the 50% chance level [$t(38)$ =
532 30.48, $p < 0.001$], indicating that participants were successfully able to attend to the target
533 speaker and comprehend the content of the story. We found a significant effect of age on
534 performance [$t(37)$ = -3.04, $p = 0.004$], with older participants performing better than younger
535 participants (YA: mean ± s.d. = 80.1 ± 7.5%, OA: 86.1% ± 4.6%). A correlation analysis with age
536 used as a continuous variable showed the same association with the proportion of correct
537 responses ($r = 0.33$, $p = 0.043$). Confidence measures showed the same general pattern of
538 results as the comprehension scores and the two measures were positively correlated [$r = 0.69$,
539 $p < 0.001$], indicating that participants had good awareness of their performance.
540        Because hearing loss was more common among the older participants, and we
541 compensated for it by amplifying the audio in frequency ranges of elevated thresholds (see
542 Methods), we assessed whether this amplification could account for the difference in
543 performance. As expected, in the portion of participants who received amplification ($n = 17$),
544 there was no relationship between average high-frequency audiogram (2-8 kHz range), and
545 comprehension-performance ($r = 0.06$, $p = 0.81$) or confidence ($r = 0.2$, $p = 0.44$) measures. The
546 same pattern was observed when using the average of the entire 0.25-8 kHz range of
547 audiometry. As such, there was no evidence that amplification had an impact on performance,
548 or that it could account for between-group differences in performance.
549        Prior to the experimental session, each participant filled out a modified subset of the
550 SSQ ($SSQ_m$) questionnaire to assess their subjective difficulties with speech-in-noise perception.
551 We found no difference in these measures between younger and older participants ($z = -0.42$, $p$
552 = 0.67, Mann-Whitney U-test), and no correlation between $SSQ_m$ score and the proportion of
553 correct responses from the behavioral task ($r = -0.17$, $p = 0.29$), or between $SSQ_m$ and high-
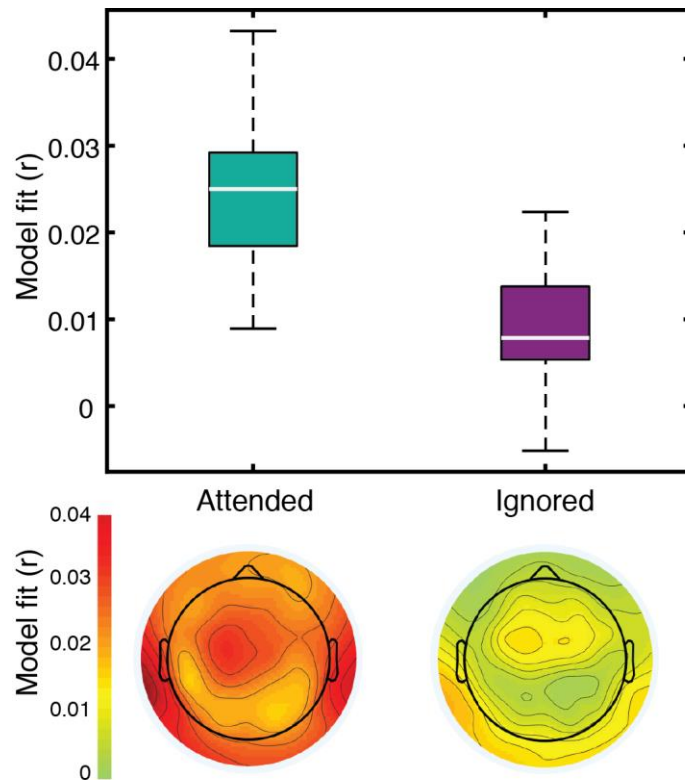554 frequency hearing loss ($r = 0.03$, $p = 0.91$).
555
556 **3.2 Cortical measures of speech-mixture processing**
557
558 In order to characterize cortical responses to semantic content of speech, we applied
559 computational models to EEG responses measured while participants listened to a mixture of
560 two distinct narrative stories, while attending to one of them. The features included in the
561 model were word audibility reflecting word-by-word fidelity of the incoming acoustic signal,
562 semantic dissimilarity reflecting short-term (sentence timescale) dissimilarities between the
563 word2vec vector characterizing each word and its immediately preceding context, and word
564 surprisal reflecting long-term predictability of each word given the preceding multi-sentence
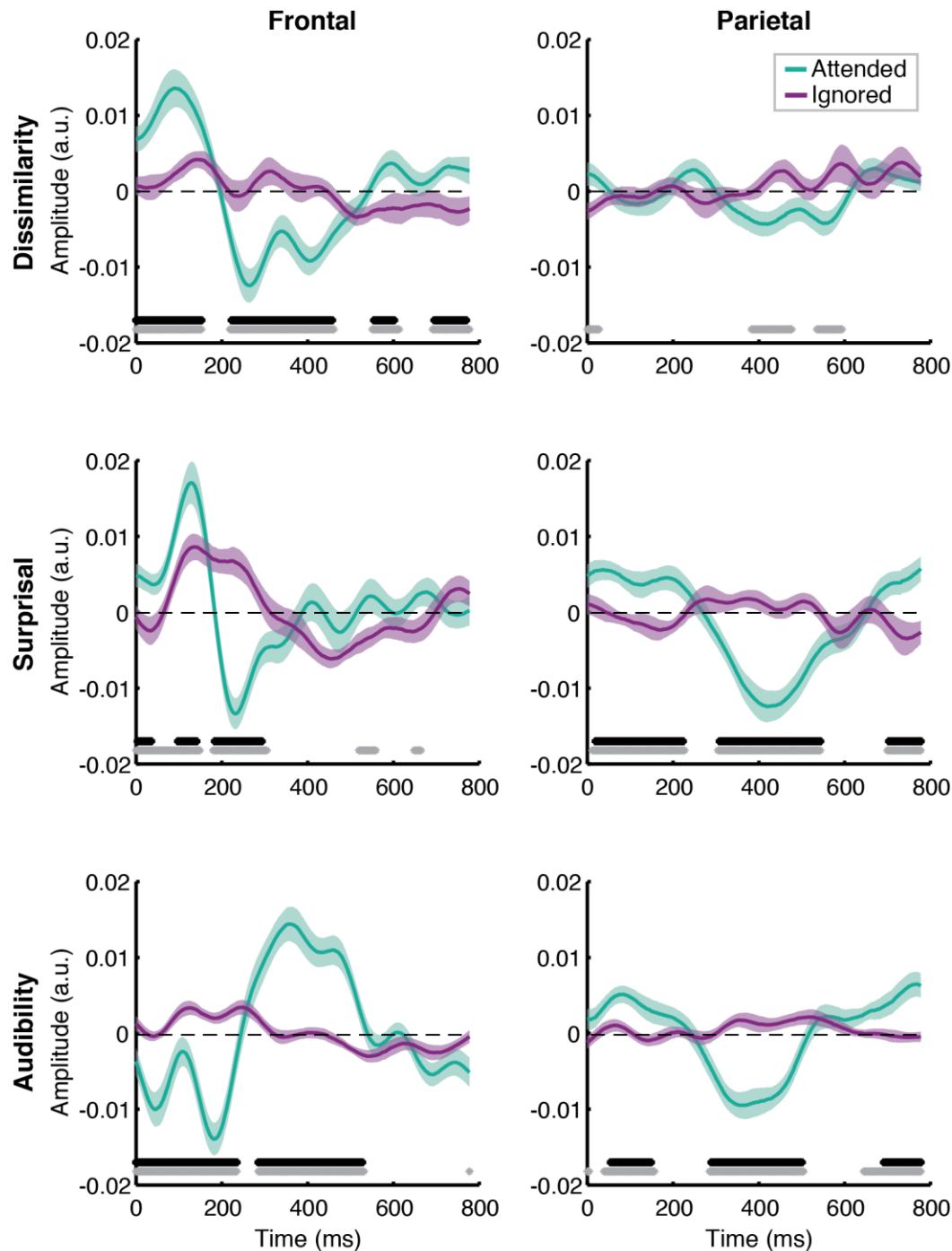565 context.
566

Figure 2. The three-feature model explained a significant amount of variance in responses to both attended and ignored speech. Box plots (top) represent distributions of goodness-of-fit values averaged over electrodes across all participants. The topographic plots (bottom) depict the distribution of goodness-of-fit values for attended and ignored speech across the scalp.

Linear regression of these features against the EEG signal produced responses that explained a significant amount of variance in the data pooled across participant groups and electrodes, as reflected by a significant positive correlation between the full-model EEG prediction and held-out data for both attended [t(38) = 20.87, p < 0.001] and ignored [t(38) = 8.75, p < 0.001] speech, with a significantly stronger fit for the former (t(38) = 10.60, p < 0.001; Fig. 2). The same pattern of results was observed when examining model fits in frontal and parietal ROIs. Figure 3 depicts the average attended (green) and ignored (purple) TRFs in the two ROIs for each of the features included in the model. We observed robust responses to the attended story for each of the features included in the model, with prominent early (~ 100 ms) and late (~ 400 ms) peaks in neural activity. In contrast, the ignored story elicited comparatively flatter responses, with predominantly early peaks in neural activity. Indeed, most features showed extensive periods in the early and late portions of the TRFs where attended and ignored responses differed significantly, as depicted by black horizontal bars at the bottom of each TRF plot (indicating FDR-corrected significant time points).
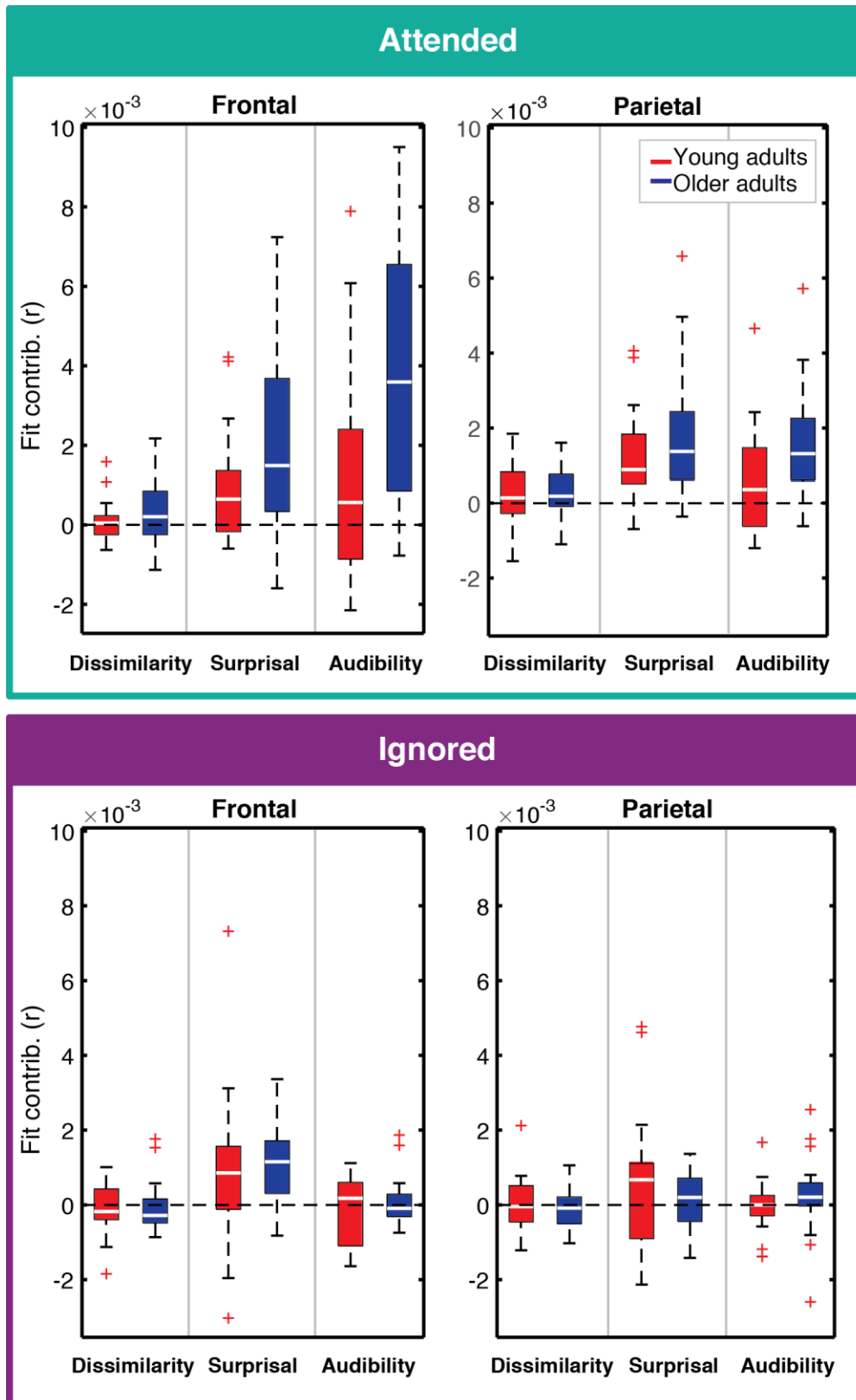
Figure 3. Attentional modulation of feature-specific responses. Each plot depicts the comparison of TRFs averaged across all participants for attended (green) and ignored (purple) speech for each of the features (panel rows) and ROIs (panel columns). The upper and lower bound of each curve represents ± 1 standard error (SE) of the mean. Black and gray horizontal bars at the bottom of the plots indicate time intervals over which attended and ignored TRFs differed significantly at the FDR-corrected and uncorrected level, respectively, with $\alpha = 0.05$.

596        Contributions of each feature to the overall model fit for both age groups are plotted in

597    Fig. 4. Model fit contribution values represent the difference in goodness-of-fit for the held-out

598    EEG data between the full model and null models in which a given feature's regressor was

599    selectively disrupted by shuffling its feature amplitudes (see section 2.9.2). Thus, for a

600    particular feature, a model fit contribution exceeding 0 represents the scenario where the EEG

601    responses scaled, to some degree, with that feature's regressor values. To compare how these

602    model contributions differed in the two age groups, we performed a mixed-design ANOVA with

603    within-subject factors of ROI, model feature, and attention, and a between-subjects factor of

604    age group (Table 1). As expected, we found a main effect of attention [$F(1,37) = 34.28$, $p <$

605    $0.001$, $\eta_p^2 = 0.48$] reflecting generally stronger tracking of high-level features within the

606    attended than ignored speech stream. We also found main effects of ROI [$F(1,37 = 8.89$, $p =$

607    $0.005$, $\eta_p^2 = 0.19$],  feature [$F(2,74) = 18.48$, $p < 0.001$, $\eta_p^2 = 0.33$], and age group [$F(1,37 = 7.92$,

608    $p = 0.008$, $\eta_p^2 = 0.18$].

609

610
611    Figure 4. Feature-specific contributions to the model fit for attended (top) and ignored
612    (bottom) responses. Each panel depicts the box plot of model fit contributions for each of the

613    three features in the younger (red) and older (blue) adult groups. Left and right panels
614    represent results for frontal and parietal ROIs, respectively. Note that some points are depicted
615    with red + signs as outliers in order to better depict where the bulk of the points lie within the
616    fit contribution distributions. However, all data points were utilized in statistical analyses
617    described in the text.

618

619        In addition to these main effects, we detected a number of significant interactions.
620    There was a significant interaction between attention and age group [F(1,37 = 7.64, p = 0.009,
621    $\eta_p^2$ = 0.17], reflecting an overall greater difference between attended and ignored fits in older
622    than younger participants [t(37) = -2.76, p = 0.009]. A significant interaction between ROI and
623    age group [F(1,37 = 7.24, p = 0.011, $\eta_p^2$ = 0.164] was associated with significantly stronger
624    contributions to model fits across features at the frontal compared to the parietal ROI in older
625    adults (p = 0.007; Mann-Whitney U-test). Third, we found a significant interaction between
626    feature and age group [F(2,74 = 4.09, p = 0.021, $\eta_p^2$ = 0.10], and a post hoc analysis revealed
627    this was due to greater difference in contributions to model fit between word audibility and
628    dissimilarity in older than younger participants [t(37) = -3.01, p < 0.005; Bonferroni corrected
629    with α = 0.017].
630        Several interactions did not involve age group, including a significant interaction
631    between attention and feature [F(1.8,66.63) = 8.55, p = 0.001, $\eta_p^2$ = 0.19], a trend towards an
632    interaction between feature and ROI [F(1.68, 62.18] = 3.2, p = 0.056, $\eta_p^2$ = 0.08], and a three-
633    way interaction between attention, feature, and ROI, [F(2,74 = 13.05, p < 0.001, $\eta_p^2$ = 0.21].
634    Because the latter interaction was a combination of factors from the former two, we only
635    pursued post hoc analyses for the three-way interaction. These indicated that in the frontal
636    ROI, the contribution of audibility to the model fit was greater for the attended than the
637    ignored story, and that this differential was greater than that for both dissimilarity and surprisal
638    [t(38) = -3.38, p = 0.002, and t(38) = -3.61, p < 0.001, respectively; Bonferroni corrected with α =
639    0.017]. Comparison of dissimilarity and surprisal showed no difference [t(38) = 1.38, p = 0.18].
640        Although the goodness-of-fit analyses above indicate that there are significant
641    differences in processing of attended and ignored speech between younger and older
642    participants, they do not provide insight into the timing and amplitude of the underlying neural
643    responses. To explore if our data contain evidence of age-related differences in neural
644    responses, we statistically compared TRF amplitudes between the two age groups at each time
645    point in the 0 to 800 ms range.  Because these analyses involved hundreds of point-by-point
646    comparisons between groups, we corrected for false discovery rate (FDF), and focused on
647    comparisons at the level of individual features, rather than utilizing more complex interaction
648    metrics. As such, these analyses were relatively rudimentary, and should be considered as
649    exploratory in their nature.

650
651    Figure 5. Between-group comparison of TRFs for attended speech. Each plot depicts a
652    comparison of TRFs between younger (red curves) and older (blue curves) participants, for
653    different features (panel rows) and ROIs (panel columns). Black and gray horizontal bars at the
654    bottom of the plots indicate time points at which the two age groups differed significantly at
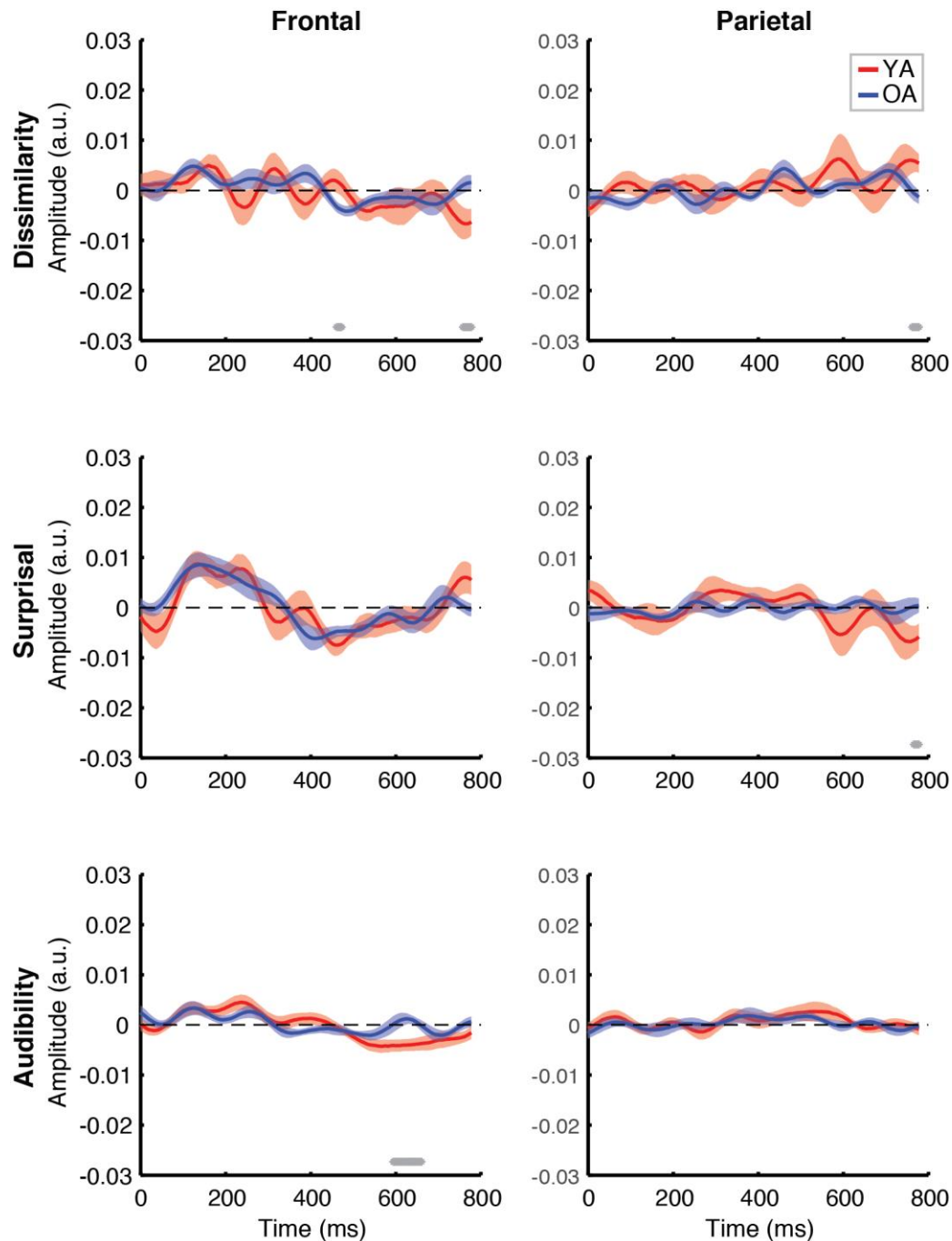655    the FDR-corrected and uncorrected level, respectively, with $\alpha = 0.05$.
656
657            Figure 5 depicts the differences in responses to the attended speech between younger
658    (red lines) and older (blue lines) participants, separately for each feature (plot rows) and ROI

659   (plot columns). Two-tailed statistical outcomes at the p < 0.05 level are depicted at the bottom

660   of each plot in both uncorrected (gray horizontal bars) and FDR-corrected (black horizontal

661   bars) forms. At the FDR-corrected level, we only found two clusters of significant time points in

662   the frontal TRFs for dissimilarity, with older participants showing a significantly more negative

663   response between approximately 260-300 ms, and a significantly more positive response in the

664   620-675 ms time range. While surprisal and audibility showed no robust differences at the FDR-

665   corrected level, several clusters of time points were suggestive of group differences at the level

666   of uncorrected statistics. For surprisal, we found that older adults had a greater negative

667   deflection in the 225-260 ms time range and a pair of positive deflections around 390-430 and

668   515-580 ms that were absent in the TRF of the young adults at the frontal ROI. We also found a

669   single cluster of time points with greater negative deflection for older than younger adults

670   between 395-435 ms in the parietal ROI. For word audibility, we found a prolonged elevated

671   response with portions between 415-480 ms exhibiting larger positive deflection in older than

672   younger participants, at the frontal ROI. Older adults also showed a greater negative deflection

673   in the word-audibility TRF frontally, and a greater positive deflection parietally around 550-600

674   ms.

675         Between-group comparison of TRFs for ignored speech are shown in Figure 6. Unlike

676   responses to attended speech, most features, with the exception of frontal TRFs for surprisal,

677   show largely flat response patterns that do not differ between groups. Several time points

678   showed a difference in uncorrected statistics for each of the features, the most notable of

679   which was a more negative response of younger adults to audibility between 590-660 ms in the

680   frontal ROI. However, given the low amplitude of the TRFs, and long latencies of most of the

681   potential differences, we believe these are likely to simply reflect false discoveries due to

682   hundreds of comparisons. Indeed, fewer than 5% of comparisons for ignored speech were

683   significant at the uncorrected level.

Figure 6. Between-group comparison of TRFs for ignored speech. Subplot arrangement and statistical comparisons are as in Fig. 5.

To complement these exploratory point-by-point analyses, we also conducted between-groups analyses specifically targeted at comparing responses in the time range of the N400 response. To this end, we compared each feature's average TRF amplitudes in the 300-500 ms range. Because previous work found little to no evidence of N400 for ignored speech, these comparisons were only done for attended speech. Although we found both a significantly more

693    negative parietal N400 for the older group to surprisal [t(37) = 2.03, p = 0.05], and a significantly
694    elevated frontal response in the older group for audibility [t(37) = -2.72, p = 0.01], neither of
695    these results remained significant with Bonferroni  correction (α = 0.008, given the total
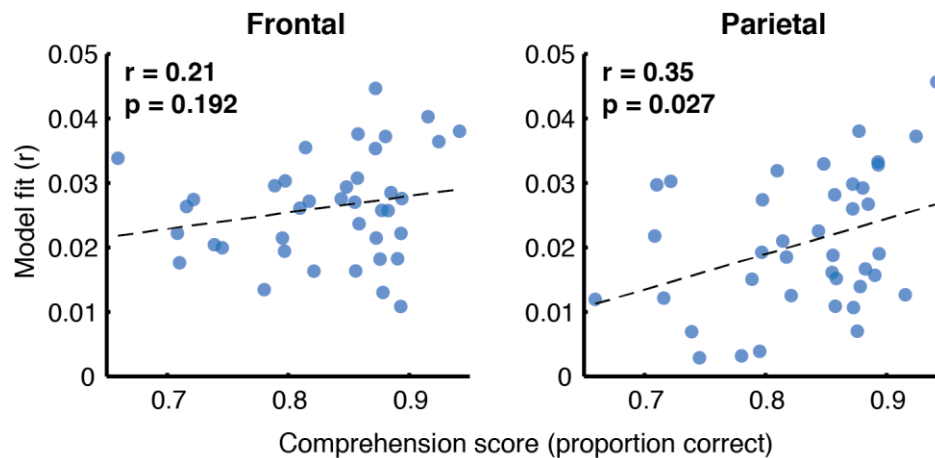696    number of 6 comparisons).

697

698

699    **3.3 Neuro-behavioral correlations**

700

701    We next sought to examine how our electrophysiological measures related to behavioral
702    responses during the experiment, and the SSQ$_m$ scores obtained prior to this experiment. To
703    this end, we conducted a number of exploratory analyses, including correlations between
704    behavioral measures and the overall model goodness-of-fit, feature-specific model
705    contributions, and the average TRF amplitudes in the 300-500 ms time range. Given the number
706    of these analyses, and our limited sample size, we focused our analyses on full participant
707    samples, rather than age group comparisons. Because of the less stringent multiple
708    comparisons correction procedure (only correcting by the number of statistical tests within
709    each analysis), significant effects in this section should be interpreted as trends rather than true
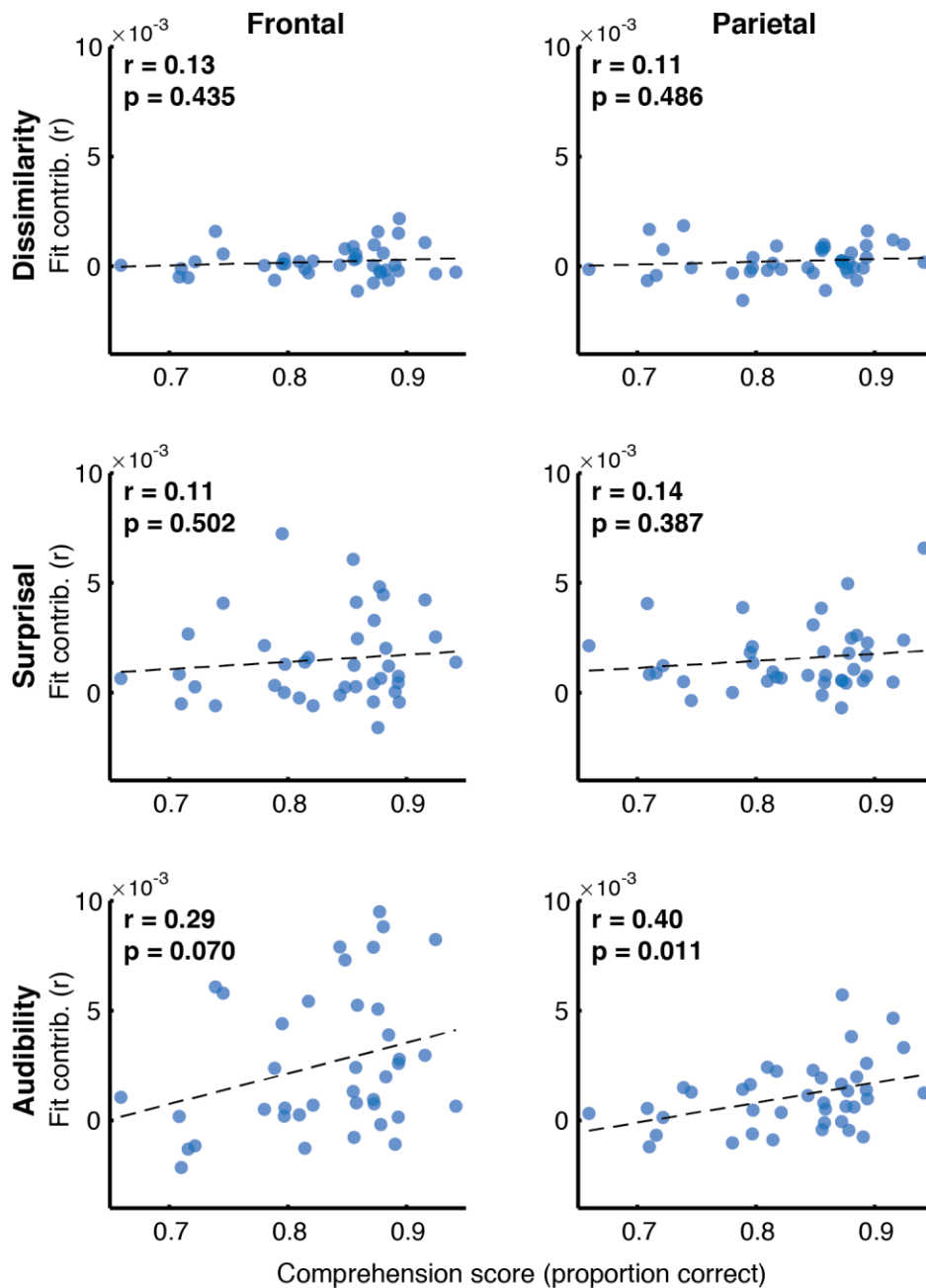710    statistical effects.

711



712
713    Figure 7. Scatterplots showing the relationship between the full model goodness-of-fit and the
714    proportion of correct responses on the comprehension questions. Pearson's correlation
715    coefficients and the corresponding uncorrected p-values are shown for frontal (left plot) and
716    parietal (right plot) ROIs. Symbols represent data from individual participants pooled across the
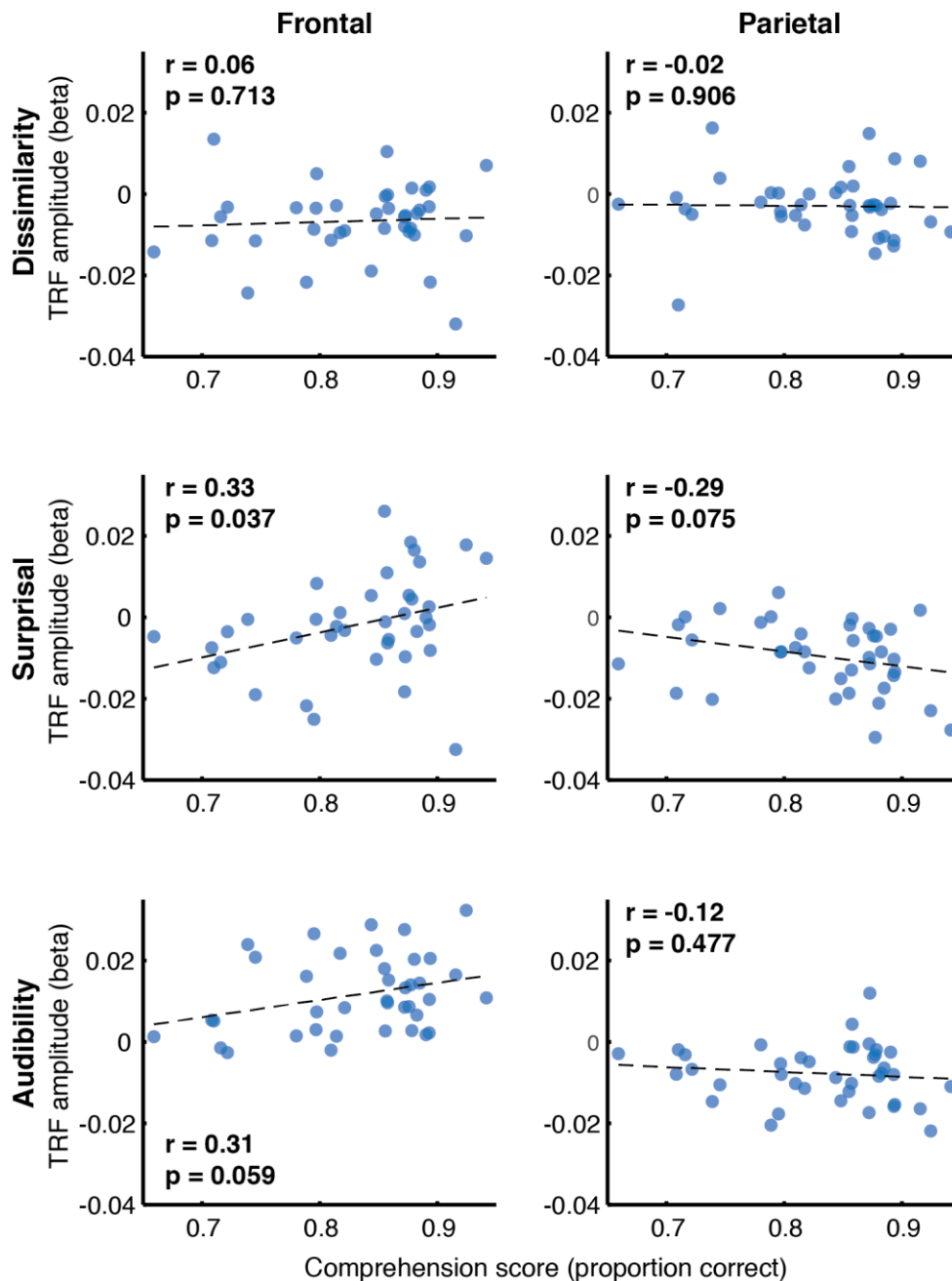717    two age groups, YA and OA.

718

719        Figure 7 depicts the relationship between the proportion of correct responses on
720    comprehension questions during the experiment, and the overall model goodness-of-fit in the
721    frontal (left panel) and parietal (right panel) ROIs. While we observed no relationship in frontal
722    regions (r = 0.21, p = 0.19), there was a marginally significant positive association between the
723    two measures (r = 0.35, p = 0.027, Bonferroni corrected α = 0.025) in the parietal ROI. A similar

724 pattern of results was observed when average confidence ratings for the comprehension
725 questions were used instead of the performance itself. Relationships between the proportion of
726 correct responses and feature-specific contributions to the model fit are depicted in Figure 8.
727 We observed a trend towards a positive association for word audibility in both the frontal (r =
728 0.29, p = 0.07) and parietal ROIs (r = 0.4, p = 0.011), although neither correlation reached
729 significance after correcting for multiple comparisons (α = 0.008). None of the other features
730 showed a significant association with comprehension scores.

731



732
733 Figure 8. Scatterplots of comprehension scores and feature-specific model contributions.
734 Different rows of panels refer to different features and different columns correspond to the
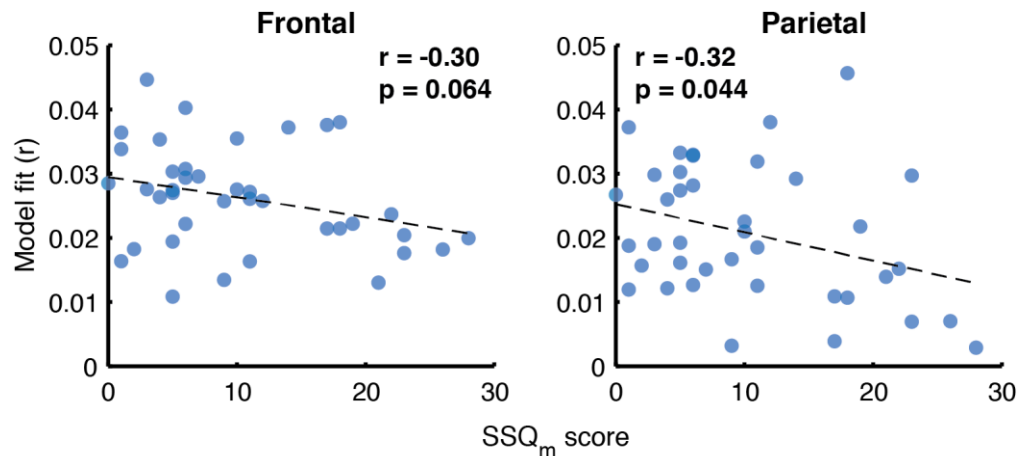
735    two ROIs.  Pearson's correlations and the corresponding uncorrected p-values are shown in the

736    upper portion of each panel.

737



738

739    Figure 9. Scatterplots of comprehension scores and mean TRF amplitudes between 300-500 ms.

740    Figure layout is as in Fig. 8.

741

742         Next, we explored the possible relationship between the comprehension scores

743    (proportion correct) and the average TRF amplitude in the 300-500 ms time range, when N400

744    effects generally appear parietally. These analyses, shown in Figure 9, revealed trends towards
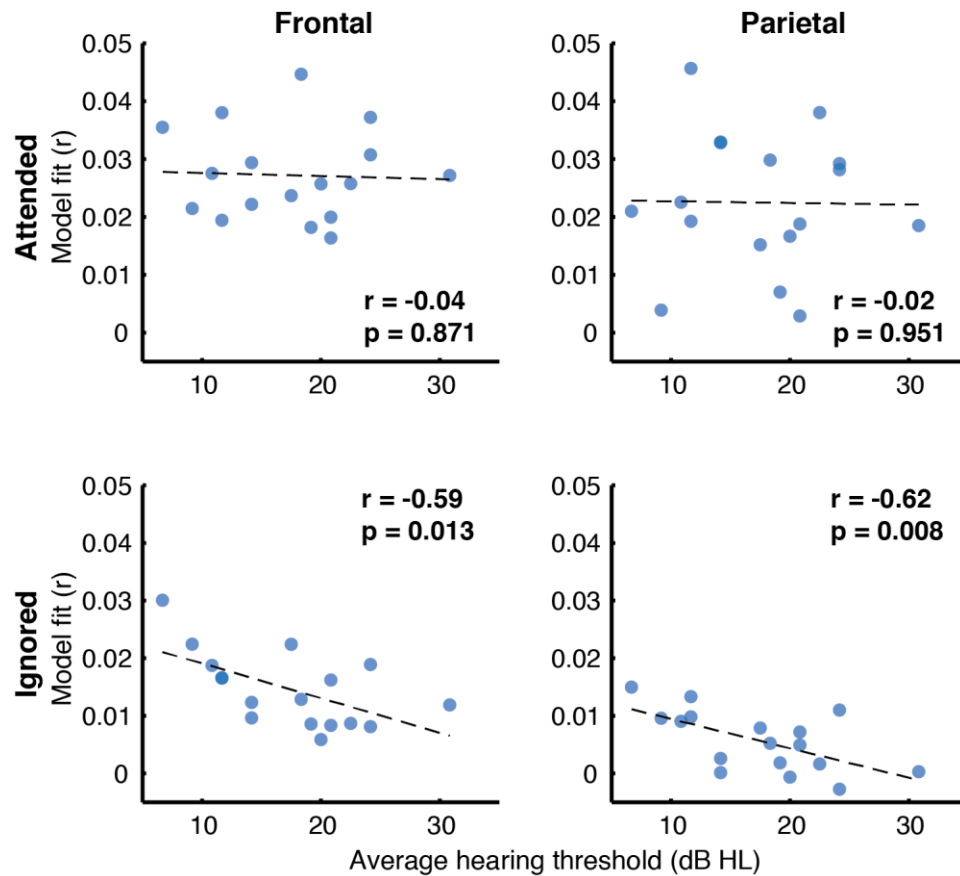
745    a positive relationship in frontal regions for surprisal (r = 0.33, p = 0.037) and audibility (r = 0.31,

746    p = 0.059), as well as a trend towards a negative relationship for surprisal in parietal ROI (r = -

747    0.29, p = 0.075). As before, none of these associations were significant when correcting for

748    multiple comparisons. Although this analysis focused broadly on the time range of N400, two of

749    the frontal trends were associated with positive, rather than negative deflections in the TRF.

750



751

752    Figure 10. Scatterplots of SSQ$_m$ scores and overall model goodness-of-fit for frontal (left panel)

753    and parietal (right panel) ROIs. Note that a higher score on SSQ$_m$ questionnaire reflects a

754    greater difficulty with understanding speech in noise.

755

756        Correlation analyses examining the relationship between subjective SIN perception

757    difficulties, captured by the SSQ$_m$ scores, and the full model goodness-of-fit metric (Fig. 10)

758    revealed trends towards a negative relationship in both the frontal (r = -0.30, p = 0.064) and

759    parietal ROIs (r = -0.32, p = 0.044). However, analyses of relationships with feature-specific TRF

760    amplitudes and model contributions revealed no feature for which these trends were apparent.

761        Finally, because a portion of the participants had mild hearing loss at high frequencies

762    (which was compensated for by amplifying speech in the corresponding frequency ranges; see

763    Methods), we examined if and how high-frequency (2-8 kHz) hearing thresholds related to the

764    overall model fits (Fig. 11). Although we found no relationship between the average hearing

765    thresholds over the 2-8 kHz range and model goodness-of-fit for attended speech (Frontal ROI:

766    r = -0.04, p = 0.87; Parietal ROI: r = -0.02, p = 0.95), there was a significant negative correlation

767    for ignored speech both frontally (r = -0.59, p = 0.013) and parietally (r = -0.62, p = 0.008). At

768    the level of feature-specific contributions to the model fit, there was no indication that this

769    negative correlation was driven by any particular feature, as most features showed low, non-

770    significant negative correlations.

771

Figure 11. Scatterplots of average high-frequency hearing thresholds (2-8 kHz) and overall model goodness-of-fit as a function of attention (panel rows) and ROI (panel columns).

**4. Discussion**

Speech perception is a fundamental capability of the human auditory and language systems, facilitating our abilities to learn and engage in various types of social interaction. However, deficits in SIN perception are commonly experienced by the aging population (e.g., van Rooij and Plomp, 1990; Goossens et al., 2017) and are reported surprisingly frequently even among the younger and nominally normal hearing population (Saunders, 1989; Zhao and Stephens, 2007; Tremblay et al., 2015). Importantly, while subjective SIN perception difficulties may indicate a significant adverse impact on quality of life (Dalton et al., 2003; Chia et al., 2007), existing objective (laboratory and clinical) measures of speech perception have shown surprisingly poor correlations with the self-reported difficulties as measured, for example, by SSQ scores (Phatak et al., 2018; Smith et al., 2019).

In the present study, we measured EEG responses to continuous two-talker speech mixtures in younger (< 40 y.o.) and older (> 40 y.o.) participants. Participants' cortical responses in the 1-8 Hz range were predicted by modeling TRFs for three speech features, short-timescale semantic dissimilarity, long-timescale lexical surprisal, and word-level audibility. We also collected behavioral measures, including participants' subjective ratings of their difficulties with

793  SIN understanding (modified SSQ), and comprehension scores for attended speech during the
794  experiment and the associated confidence ratings.
795       Our three-feature model was able to explain significant variance in the EEG data,
796  especially in responses to attended speech, where each of the features contributed to the
797  neural responses (Fig. 4). The evidence for this was particularly strong for surprisal and
798  audibility, suggesting that these model features captured stimulus characteristics that were
799  actively tracked by our participants' auditory systems. Moreover, we found that participants'
800  performance on the comprehension task (Fig. 7), as well as the associated confidence ratings,
801  showed a trend towards a positive correlation with the goodness of the overall model fit for the
802  attended speech, suggesting that successfully tracking these features is related to speech
803  comprehension. Although our data does not support a strong association between
804  performance and model contributions, or TRF magnitudes, for any one of the model features,
805  we did find trends towards an association between word audibility and performance in both
806  ROIs (for both model fit contributions, and TRF magnitudes), and in the frontal region between
807  the surprisal TRF magnitude and performance. Speculatively, these trends suggest that
808  improved comprehension may be related to at least two cognitive processes. First, the
809  association with audibility suggests that improved performance may stem from more effective
810  weighing of word-level information by word reliability, as reflected by the word SNR. Second,
811  the association with surprisal suggests that high performance may be related to increased
812  sensitivity to lexical and/or semantic associations between different segments of speech.
813       Consistent with previous work on neural representations of two-talker speech (Ding and
814  Simon, 2012; Mesgarani and Chang, 2012; Broderick et al., 2018; O'Sullivan et al., 2019) we
815  found robust differences between responses to attended and ignored speech both in the
816  goodness of model fits and the TRFs. In general, model fits were better for attended than
817  ignored speech (Fig. 4) and the associated TRFs for attended speech showed complex, multi-
818  peaked morphologies, whereas the responses to ignored speech were flatter and contained
819  fewer prominent peaks (Fig. 3). Thus, our results indicate that responses to a speech mixture
820  preferentially reflect attended speech, while representations of distractor speech are largely
821  suppressed.
822       Comparisons of EEG responses between age groups revealed a complex pattern of age-
823  related differences, captured particularly by model fit measures.  Specifically, we found that
824  older participants exhibited on average greater differences in feature-specific model-fit
825  contributions between attended and ignored speech. This age effect was driven primarily by
826  better fits for attended speech in the frontal ROI (see Fig. 4). Although to a weaker degree,
827  these differences were mirrored in attended TRFs, in that older adults showed generally
828  stronger TRF deflections from 0 compared to younger participants (Fig. 5). In most cases,
829  however, these TRF differences did not reach statistical significance when controlling for false
830  discovery rate, possibly due to nuisance factors such as inter-subject variability in cortical
831  geometry, and/or inadequate sample size.
832       With respect to the modelled features, we found that surprisal and audibility both
833  showed stronger frontal contributions in older adults, whereas parietal contributions were

834    relatively similar between the two groups. We speculate that the stronger fits in the frontal

835    region in older adults may be indicative of heightened reliance in this group on both lexical

836    prediction, as reflected by increased accuracy of surprisal fits, and on words with better

837    audibility. Higher word SNRs may have been more important for disambiguation of the masked

838    portions of speech for older compared to younger adults. Although audibility itself reflects a

839    relatively low-level aspect of our stimuli, its frontal TRF profile showed a prolonged positive

840    deflection in the 250-550 ms latency range. Such a long latency is consistent with the

841    possibility that this audibility-related response may reflect engagement of higher-level

842    processes, such as retrospective disambiguation, or prospective prediction.

843          It is notable that participants in the older group exhibited significantly better

844    performance on the comprehension task than younger adults, despite having greater

845    prevalence of hearing loss (15 out of 17 participants with HL were in the older group and the

846    degree of HL was not significantly correlated with performance). This difference in performance

847    difference complicates the interpretation of age-related differences in neural responses. It may

848    be the case that older adults in our participant sample were either more engaged, or exerted

849    greater effort in the task, which in turn led to stronger speech tracking in their EEG data, as well

850    as better performance. This is plausible, since more participants in the older group (12/20 older

851    vs 8/19 younger participants) indicated having a subjective sense of experiencing greater

852    difficulty with SIN understanding compared to their peers. The sense of greater difficulty may

853    have motivated at least some of the older participants to exert greater effort to perform well.

854    However, while the average performance of participants with self-reported SIN difficulties was

855    slightly better than that in participants who did not report such difficulties, these differences

856    were not significant. Despite this, the possibility that differences between the two age groups in

857    effort, attentiveness, or another factor may underlie the neural differences discussed above,

858    deserves further attention in future work.

859

860    **4.1 Relationship to existing work on age-effects on electrophysiological measures of speech**

861    **processing**

862

863    Several studies have examined effects of age (Presacco et al., 2016; Decruy et al., 2019; Zan et

864    al., 2020) and hearing loss (Millman et al., 2017; Decruy et al., 2020) on continuous speech

865    processing in the context of envelope tracking. Generally, these studies have demonstrated

866    that older adults and those with hearing loss exhibit exaggerated cortical tracking of speech

867    envelope both in quiet and in the presence of a competing speaker, as reflected by higher

868    envelope reconstruction accuracies from delta-band EEG or MEG responses in these

869    populations. Our analyses show a similar pattern of amplified feature tracking in the aging

870    population, albeit for word-level features. Responses to the audibility feature, in particular, may

871    reflect similar underlying processes as those involved in envelope processing. However,

872    audibility in our study was defined as the word-by-word ratio between the acoustic energy in

873    the two speech waveforms, rather than the absolute amplitude of each speech signal, making

874    direct comparisons of the two measures difficult. Distinct from envelope TRFs, the audibility

875  TRF in our study contained prolonged deflections from 0 in the 300-500 ms latency range,
876  suggesting that our measure may tap into additional higher-level processes. Although lexical
877  surprisal is seemingly unrelated to speech envelope, it is possible that predictive processes may
878  interact with lower-level stimulus encoding via feedback processes, as has been demonstrated
879  for dissimilarity (Broderick et al., 2019).
880         While measures of envelope tracking have provided important insights into speech
881  processing, they are largely uninformative about the nature of higher-level processes involved
882  in speech perception. In recent years, an increasing number of studies have investigated  the
883  relationship of electrophysiologically-measured cortical responses to both intermediate speech
884  representations such as those evoked by different phoneme categories (Di Liberto, 2015;
885  Lesenfants, 2020; Teoh & Lalor, 2020; but cf. Daube et al. 2019) or phonotactics (Di Liberto,
886  2019), and word-level representations related to lexical (e.g., Brodbeck et al., 2018), as well as
887  syntactic and semantic (Broderick et al., 2018;  Weissbart et al., 2019; Heilbron et al., 2019;
888  Donhauser & Baillet, 2020) processing. Nevertheless, relatively little is known about how these
889  representations change as a function of age, particularly in challenging listening conditions.
890  Recently, Broderick et al. (2020) compared representations of semantic dissimilarity and 5-gram
891  lexical surprisal derived from responses to clean speech in younger and older adults. They
892  showed that although younger adults exhibited robust responses to each feature, older adults
893  only showed strong responses to lexical surprisal (albeit with a delayed peak response), with a
894  nearly absent response to semantic dissimilarity. These results were interpreted as potentially
895  reflecting lesser reliance of older adults on semantic predictive process, thought to be captured
896  by the dissimilarity feature, due to age-related cognitive decline. Consistent with this, older
897  participants with greater semantic verbal fluency, a measure related to the ability to engage in
898  semantic prediction, showed greater contribution of semantic dissimilarity to the model of
899  cortical responses to speech.
900         Because our experimental design involved listening to a more challenging, two-speaker
901  mixture, direct comparisons of our results with those of Broderick et al. (2020) are not possible.
902  Nevertheless, there are marked differences between the patterns of results observed in their
903  study compared to ours. In particular, we observed stronger tracking of both lexical surprisal
904  and word audibility in older than younger adults, and generally weak but otherwise similar
905  tracking of dissimilarity in the two groups. Notably, this was observed predominantly at the
906  frontal ROI, with the posterior ROI showing a smaller difference (albeit in the same direction as
907  the frontal results). In contrast, Broderick et al. focused their analyses on posterior electrode
908  sites, making it unclear how tracking of their features behaved at more frontal sites that are
909  involved in tasks relying on working memory (e.g., Gevins et al., 1997; Onton et al., 2005).
910         In Broderick et al. (2020), the greatest age-related differences were shown for semantic
911  dissimilarity, whereas our goodness-of-fit results showed relatively weak contributions from
912  this feature (compared to surprisal and word audibility) that did not differ significantly between
913  the younger and older age groups. However, we did observe greater frontal TRF deflections in
914  the older group for dissimilarity, with significant group differences around 250 and 600 ms,
915  suggesting an increased gain for this feature in the older population. This underscores the

916    importance of analyzing both model fits and the corresponding TRFs, as morphological
917    differences in the latter may be possible even in the absence of differences in the model
918    goodness-of-fit. The most notable difference in our results with respect to dissimilarity is that
919    we did not observe posterior N400 response in either group, in contrast to the significant
920    parietal N400 in the TRF for dissimilarity in older but not younger adults reported by Broderick
921    et al. Although this discrepancy is puzzling given the use of nearly identical methods for
922    computing dissimilarity, it raises the possibility that the utility of dissimilarity may be limited if
923    other features, which better capture neural responses that would otherwise be attributed to
924    dissimilarity, are included in the model.

925        Another important difference between the two studies pertains to the role of surprisal
926    in the models fitted to the data. Specifically, unlike the relatively simple 5-gram surprisal used
927    by Broderick et al., which was intended to capture responses related to the knowledge of word
928    co-occurrence within 5-word neighborhoods, the surprisal features utilized in our study were
929    computed using an advanced natural language model (GPT-2; Radford et al., 2019) that uses
930    preceding context of up to several hundred words (i.e., dozens of sentences) in order to
931    estimate each upcoming word. As such, surprisal in our study likely captured responses related
932    to higher-level lexical and/or syntactic predictions. Thus, although responses to these two
933    surprisal measures cannot be directly compared, the stronger tracking of surprisal by older
934    adults in our study is consistent with increased reliance on predictive processes in this
935    population. This is in agreement with behavioral results demonstrating greater reliance on
936    semantic context in populations with compromised representations of speech, such as those
937    with hearing loss (Benichov et al., 2012; Lash et al., 2013) and cochlear implants (Amichetti et
938    al., 2018; Dingemanse and Goedegebure, 2019; O'Neill et al., 2019).

939        Importantly, the seemingly conflicting pattern of results between these studies could in
940    fact reflect two distinct contributors to speech perception difficulties in older adults, namely
941    decreases in the fidelity of lower level representations, and cognitive decline. Prevalence of
942    mild high-frequency hearing loss in our sample of older adults was quite high, making it likely
943    that decreased fidelity of peripheral representations had an effect on our results. While
944    Broderick et al. did not report audiogram measures for their sample of older adults, the mean
945    age was considerably greater in their study (mean ± s.d. = 63.9 ± 6.7 years vs 53.5 ± 8.7 years in
946    this study), making it likely that similar or greater hearing difficulties may have impacted their
947    participants. However, because of the age difference in the two samples, the effects of
948    cognitive decline may have contributed more significantly to the results of Broderick et al., and
949    may potentially explain why measures related to predictive processes showed opposite effects
950    in the two studies. This exemplifies the complex combination of etiologies that may underlie
951    speech perception difficulties, and the distinct ways in which they may affect speech
952    processing. Future work should attempt to quantify these factors and use multivariate analyses
953    to better characterize if and how they may relate to different neural measures of speech
954    processing.
955

956    **4.2 Higher-level speech feature tracking as an index of speech in noise perception difficulties**

957

958  A key reason for our choice to study responses to lexical and semantic features is their
959  potentially greater sensitivity to SIN perception difficulties, compared to responses driven by
960  lower-level features such as the speech envelope. Specifically, because dissimilarity and
961  surprisal (but not audibility) depend on preceding lexical and semantic context, in order for
962  language processing mechanisms to accurately track them, each word within the sequence
963  needs to be recognized and integrated with the preceding context. Lower-level SIN processing
964  impairments may thus disproportionately impact tracking of these features. This is because
965  missing a given word may potentially distort neural computations of surprisal and lexical
966  predictions for a large number of subsequent words. This distortion could result in a mismatch
967  between the objectively computed sequences of these features (used in the model) and their
968  internal estimates.

969  Dissimilarity, in particular, depends on local word context (limited to one sentence, in
970  our model). Misperception of individual words may thus greatly distort the internal estimates
971  of the semantic relationships between words within this short-term context, leading to poor
972  correspondence with the objectively computed dissimilarity values. Spectrally degraded speech
973  has previously been shown to elicit weaker N400 responses, and a reduced difference in N400
974  between sentences with high and low cloze probabilities (Aydelott et al., 2006; Obleser and
975  Kotz, 2011; Carey et al., 2014). Similarly, our results showed weak model contributions of
976  dissimilarity with N400 responses essentially absent in the posterior ROI, consistent with the
977  possibility that challenging listening scenarios may indeed disrupt representations related to
978  relationships between words in a local context. Notably, however, we did not observe a reliable
979  association between individual differences in the tracking of this feature, or the magnitude of
980  N400, and performance on the comprehension task, the associated confidence measures, or
981  the SSQ$_m$. As such, the magnitude of dissimilarity tracking, or the associated TRFs, may not
982  actually reflect the degree of SIN perception difficulties, as we hypothesized it would. Thus, it is
983  possible that weak tracking of dissimilarity in our study may reflect that dissimilarity, as
984  computed here, is a relatively unimportant feature for characterizing cortical speech
985  processing. Note that although our results appear to be at odds with Broderick et al. (2018),
986  who demonstrated robust dissimilarity-related N400 responses for both clean and two-talker
987  speech, that study used dissimilarity as the sole feature. It is, therefore, possible that their
988  estimated TRFs may have captured contributions from other features time-locked to word
989  onsets (e.g., ones related to lexical and syntactic processing). Indeed, in a recent reanalysis of
990  cocktail party data from Broderick et al. (2018), Dijkstra et al. (2020) showed that replacing
991  dissimilarity values in a regressor with unit-amplitude impulses leads to estimation of
992  essentially identical TRFs to those obtained with the impulses scaled by dissimilarity features.
993  This insensitivity to impulse scaling calls into question the extent to which said TRFs reflect
994  dissimilarity-related processing. Comparisons of single-feature TRFs derived from our data using
995  word onset and dissimilarity regressors (analyses not shown here) mirrored these observations,
996  suggesting that the utility of dissimilarity in explaining EEG responses to continuous speech may
997  be limited.

33

998    In contrast to dissimilarity, our observation of robust model contributions and posterior
999    N400 responses for surprisal suggests that this feature may be relatively robust to challenging
1000   listening scenarios. This may be the case because surprisal, as defined in the present study,
1001   reflects predictability of each word given a multi-sentence preceding context (vs. single-
1002   sentence context for dissimilarity), potentially making misperception of individual words have
1003   relatively low impact on lexical predictions. In other words, failure to recognize individual words
1004   may have a relatively small impact on the internal predictions, as these may be highly
1005   constrained in natural speech by the successfully identified words within the longer-term
1006   context. Admittedly, the apparent robustness of surprisal to adverse listening conditions may
1007   be specific to longer narratives where long-term semantic dependencies exist, such as
1008   audiobooks used in our study. In contrast to dissimilarity, we did observe weak trends
1009   suggesting an association between the amplitude of the surprisal TRF in the N400 latency
1010   range, and the performance on the comprehension questions. As such, it is possible that
1011   surprisal responses may indeed reflect the extent of SIN perception difficulties. However,
1012   because these trends were not statistically robust to multiple-comparisons correction, and
1013   because similar trends were not observed for $SSQ_m$, it remains unclear if this neuro-behavioral
1014   association is reliable. A replication study with a larger sample size, improved EEG denoising
1015   algorithms, and/or more sensitive behavioral measures may be needed to further explore this
1016   link.
1017        It is notable that the correlations between $SSQ_m$ or task performance and feature-
1018   specific model contributions were overall relatively weak in this study. Although this implies
1019   that none of the features utilized in our study can on their own predict the degree of SIN
1020   perception difficulties, it is possible that such deficits may be better characterized in terms of a
1021   multi-dimensional pattern of feature-specific neural responses. In other words, it may be the
1022   case that in order to predict the extent of SIN perception difficulties, a combination of neural
1023   measures across multiple lower- and higher-level speech features needs to be taken into
1024   account. Along these lines, Lesenfants et al. (2019) showed that speech reception thresholds
1025   can be predicted from EEG responses to speech more accurately using a model that contains
1026   both spectrogram and phonetic features, compared to models containing only one of the
1027   features. Furthermore, because SIN perception difficulties can have different underlying
1028   etiologies, with different relative contributions from peripheral damage and cognitive factors, it
1029   may be the case that distinct patterns of feature-specific responses characterize different
1030   underlying causes of SIN deficits.
1031
1032   **4.3 Behavioral correlates of self-reported SIN difficulties**
1033
1034   Our data revealed a trend towards a negative association between $SSQ_m$ and the overall model
1035   goodness-of-fit for attended speech. This is not surprising, as higher $SSQ_m$ scores reflect greater
1036   subjective difficulty with SIN perception, which would be expected to be related to poorer
1037   tracking of attended speech in the presence of competing speech. However, we found no
1038   correlation between $SSQ_m$ and performance on the comprehension task ($r = -0.17$, $p = 0.29$),

1039    suggesting that even participants with potentially more deteriorated representations of
1040    attended speech had sufficient fidelity of speech representations to achieve high task
1041    performance. The lack of a relationship between subjective SIN perception difficulties and
1042    performance is unintuitive, but mirrors similar results showing only a weak relationship
1043    between subjective and objective measures of SIN difficulties (Phatak et al., 2018; Smith et al.,
1044    2019).
1045         While statistical associations between subjective and objective measures of speech
1046    perception have generally been poor in past work, it is possible that these outcomes are a
1047    result of insufficiently sensitive methods for measuring speech perception. Specifically, typically
1048    used methods for objectively measuring speech perception involve presentations of isolated
1049    sentences, and having participants repeat them back, usually without time constraints (i.e.,
1050    allowing participants to deliberate and piece together their percept). While these measures are
1051    simple and effective in measuring speech perception deficits in populations with moderate and
1052    severe hearing loss (e.g., Phatak et al., 2018), the external validity of these measures may be
1053    limited at best, as they do not reflect real-world listening scenarios. Specifically, real-world
1054    spoken communication generally requires real-time comprehension of complex, multi-sentence
1055    expressions embedded in noisy and reverberant backgrounds, in order to allow for continuous
1056    flow of interaction. Unlike the commonly used speech understanding tasks, these realistic
1057    scenarios allow little time for deliberation about individual words, as new information is
1058    continuous, creating the possibility of falling behind if speech processing is impaired or slowed.
1059    Indeed, Xia et al. (2017) demonstrated marked differences in performance between tasks
1060    involving simple word identification and answering comprehension questions about the
1061    content of narrative stories, with the latter showing a weaker benefit from hearing aids. This
1062    highlights the possibility that traditional speech recognition tasks may indeed be missing
1063    important, behaviorally relevant aspects of speech perception.
1064         In the present study, a continuous multi-talker design with a behavioral task focused on
1065    assessing comprehension was selected in an attempt to mimic some aspects of real-world
1066    speech perception scenarios. Nevertheless, there were important differences that may have
1067    contributed to our failure to detect a relationship between subjective SIN perception difficulty
1068    (reflected in $SSQ_m$) and behavioral performance. First, although we utilized co-located target
1069    and distractor speakers, which are generally more challenging to parse out than spatially-
1070    separated speakers (Marrone et al., 2008; Kidd et al., 2010), their fixed location, predictable
1071    temporal characteristics (e.g., lack of sudden offsets and onsets in speaking), and relatively
1072    monotone speaking styles likely facilitated participants' ability to suppress unwanted
1073    processing of the ignored speaker. In contrast, realistic conversational settings such as
1074    restaurants or bars generally contain distractor signals that vary less predictably in location,
1075    intensity, emotional content, and other characteristics, likely contributing to greater distraction
1076    and informational masking. It is possible that suppression of these types of distractor
1077    information becomes impaired with age due to deterioration of attentional and other cognitive
1078    resources. Second, although we attempted to quantify comprehension, as opposed to mere
1079    word identification, of the content spoken by the target speaker via multiple-choice questions,

1080   it is possible that the implementation of this task lacked sensitivity to detect speech
1081   comprehension deficits. Specifically, the fact that the target story spanned many minutes may
1082   have allowed the participants to utilize much longer semantic context to aid the interpretation
1083   of incoming information, compared to real-world interactions where topics often change more
1084   rapidly. This was compounded by the fact that, for practical purposes, the questions were
1085   framed in a Yes/No format, only requiring participants to identify the more likely of the two
1086   options, rather than to demonstrate their own understanding of the story. While the main
1087   purpose of the comprehension questions was to verify that participants followed the task
1088   instructions, future work should take steps towards optimizing behavioral measures of
1089   comprehension. For example, questions carefully calibrated to require roughly constant reading
1090   time could be used to measure reaction times in addition to mere percent correct measures,
1091   possibly revealing significant response delays in people with self-reported SIN difficulties.
1092
1093   **4.4 Limitations**
1094
1095   Although our work provides evidence of age-related differences in cortical tracking of word-
1096   level features, a notable limitation of our method is that it does not establish the source of this
1097   difference. Specifically, it is unclear from our data if the distinct patterns of feature-tracking
1098   were a result of higher-order linguistic mechanisms receiving inputs with differing fidelities
1099   from lower-level processes, or they reflected age-related changes in the higher-order
1100   mechanisms themselves, or some combination of the two. Furthermore, differential
1101   engagement in cognitive resources (e.g., due to differential effort) may also have contributed to
1102   the observed differences, even in the absence of actual changes in the underlying mechanisms.
1103   Thus, an important goal for future work is to characterize speech representations more
1104   thoroughly at multiple levels of the processing hierarchy in order to elucidate the mechanisms
1105   implicated in the differences in speech processing. Furthermore, the measurement of speech
1106   representations at multiple stages of the language processing hierarchy may be critical for
1107   explaining individual differences in speech perception performance, and subjective measures
1108   such as the $SSQ_m$**.**
1109          The use of artificial neural networks (ANNs) to extract abstract features related to lexical
1110   and semantic content of speech has become increasingly common in studies of language
1111   processing (Huth et al., 2016; Broderick et al., 2018; Weissbart et al., 2019; Donhauser and
1112   Baillet, 2020). While powerful in characterizing brain responses to speech, an important
1113   limitation in the use of these features is that it can be difficult to interpret what aspects of
1114   language they actually capture. Specifically, ANNs are usually trained on a task such as text
1115   prediction on the basis of preceding context, and as such, ANNs may utilize any number of
1116   statistical regularities in the training corpus in order to optimize their performance. Thus,
1117   depending on the ANN architecture, aspects of language including the syntactic structure,
1118   lexical frequency, semantic relationships, and others may all contribute to the performance of
1119   ANNs.  Without knowing the language aspects learned by ANNs, it is difficult, and may be even
1120   impossible, to parse out the relative contributions of the different variables. Consequently,

1121  when cortical responses are found to track these features, as is the case in the present study, it
1122  may remain unclear what linguistic processes underlie this tracking. Thus, improving the
1123  interpretability of neural analyses that utilize complex natural language models remains an
1124  important challenge for future work.
1125
1126  **5 Conclusions**
1127
1128  The present study extends upon the existing body of work demonstrating the plausibility of
1129  measuring cortical tracking of high-level features related to speech meaning and predictability.
1130  The results show evidence of age-related amplification in tracking of these features in
1131  competing speech streams. Moreover, our exploratory analyses showed trends of correlations
1132  between these measures and behavioral measures including comprehension performance and
1133  subjective SIN perception difficulty scores, indicating their potential behavioral relevance.
1134  Taken together, our work demonstrates the utility of modeling cortical responses to multi-
1135  talker speech using complex, word-level features and the potential for their use to study
1136  changes in speech processing due to aging and hearing loss.
1137
1138  **Data availability**
1139
1140  Data is not available publicly, as data sharing was not a part of the informed consent. Requests
1141  to access the dataset should be directed to JM (mesik002@umn.edu).
1142
1143  **Ethics statement**
1144
1145  The Institutional Review Board of the University of Minnesota approved the procedures in this
1146  study. All participants provided written informed consent to participate.
1147
1148  **Author contributions**
1149
1150  JM and MW designed the experiment, analyzed the data, and wrote the manuscript. JM and
1151  LAR implemented experimental procedures and collected the data. All authors commented on
1152  the manuscript and approved the submitted version.
1153

1165

**Conflict of Interest Statement**

1167

1168  The authors declare no conflicts of interest.

1169

**References**

1171

1172  Akeroyd, M. A. (2008). Are individual differences in speech reception related to individual
1173      differences in cognitive ability? A survey of twenty experimental studies with normal and
1174      hearing-impaired adults. *Int. J. Audiol.* 47. doi:10.1080/14992020802301142.
1175  Amichetti, N. M., Atagi, E., Kong, Y.-Y., and Wingfield, A. (2018). Linguistic context versus
1176      semantic competition in word recognition by younger and older adults with cochlear
1177      implants. *Ear Hear.* 39, 101–109. doi:10.1097/AUD.0000000000000469.
1178  Aydelott, J., Dick, F., and Mills, D. L. (2006). Effects of acoustic distortion and semantic context
1179      on event-related potentials to spoken words. *Psychophysiology* 43, 454–464.
1180      doi:10.1111/j.1469-8986.2006.00448.x.
1181  Bell, A. J., and Sejnowski, T. J. (1995). Blind separation and blind deconvolution: an information-
1182      theoretic approach. *ICASSP, IEEE Int. Conf. Acoust. Speech Signal Process. - Proc.* 5, 3415–
1183      3418. doi:10.1109/icassp.1995.479719.
1184  Benichov, J., Cox, L. C., Tun, P. A., and Wingfield, A. (2012). Word recognition within a linguistic
1185      context: effects of age, hearing acuity, verbal ability, and cognitive function. *Ear Hear.* 33,
1186      250–256. doi:10.1097/AUD.0b013e31822f680f.
1187  Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and
1188      powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B* 57, 289–300.
1189      doi:10.1111/j.2517-6161.1995.tb02031.x.
1190  Bigdely-Shamlo, N., Mullen, T., Kothe, C., Su, K.-M., and Robbins, K. A. (2015). The PREP
1191      pipeline: standardized preprocessing for large-scale EEG analysis. *Front. Neuroinform.* 9, 1–
1192      20. doi:10.3389/fninf.2015.00016.
1193  Brainard, D. H. (1997). The psychophysics toolbox. *Spat. Vis.* 10, 433–436.
1194      doi:10.1163/156856897X00357.
1195  Brodbeck, C., Hong, L. E., and Simon, J. Z. (2018). Rapid transformation from auditory to
1196      linguistic representations of continuous speech. *Curr. Biol.* 28, 3976-3983.e5.
1197      doi:10.1016/j.cub.2018.10.042.
1198  Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J., and Lalor, E. C. (2018).
1199      Electrophysiological correlates of semantic dissimilarity reflect the comprehension of
1200      natural, narrative speech. *Curr. Biol.* 28, 803-809.e3. doi:10.1016/j.cub.2018.01.080.
1201  Broderick, M. P., Anderson, A. J., and Lalor, E. C. (2019). Semantic context enhances the early
1202      auditory encoding of natural speech. *J. Neurosci.* 39, 7564–7575.
1203      doi:10.1523/JNEUROSCI.0584-19.2019.

Broderick, M. P., Liberto, G. M. Di, Anderson, A. J., Rofes, A., and Lalor, E. C. (2020). Dissociable electrophysiological measures of natural language processing reveal differences in speech comprehension strategy in healthy ageing. *bioRxiv*, 1–17. Available at: https://www.biorxiv.org/content/10.1101/2020.04.17.046201v1.full.pdf.

Carey, D., Mercure, E., Pizzioli, F., and Aydelott, J. (2014). Auditory semantic processing in dichotic listening: effects of competing speech, ear of presentation, and sentential bias on N400s to spoken words in context. *Neuropsychologia* 65, 102–112. doi:10.1016/j.neuropsychologia.2014.10.016.

Chia, E. M., Wang, J. J., Rochtchina, E., Cumming, R. R., Newall, P., and Mitchell, P. (2007). Hearing impairment and health-related quality of life: the blue mountains hearing study. *Ear Hear.* 28, 187–195. doi:10.1097/AUD.0b013e31803126b6.

Comon, P. (1994). Independent component analysis, a new concept? *Signal Processing* 36, 287–314. doi:10.1016/0165-1684(94)90029-9.

Crosse, M. J., Di Liberto, G. M., Bednar, A., and Lalor, E. C. (2016). The multivariate temporal response function (mTRF) toolbox: A MATLAB toolbox for relating neural signals to continuous stimuli. *Front. Hum. Neurosci.* 10, 1–14. doi:10.3389/fnhum.2016.00604.

Dalton, D. S., Cruickshanks, K. J., Klein, B. E. K., Klein, R., Wiley, T. L., and Nondahl, D. M. (2003). The impact of hearing loss on quality of life in older adults. *Gerontologist* 43, 661–668. doi:10.1093/geront/43.5.661.

Decruy, L., Vanthornhout, J., and Francart, T. (2019). Evidence for enhanced neural tracking of the speech envelope underlying age-related speech-in-noise difficulties. *J. Neurophysiol.* 122, 601–615. doi:10.1152/jn.00687.2018.

Decruy, L., Vanthornhout, J., and Francart, T. (2020). Hearing impairment is associated with enhanced neural tracking of the speech envelope. *Hear. Res.* 393, 107961. doi:10.1016/j.heares.2020.107961.

Delorme, A., and Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *J. Neurosci. Methods* 134, 9–21. doi:10.1016/j.jneumeth.2003.10.009.

Dijkstra, K., Desain, P., and Farquhar, J. (2020). Exploiting electrophysiological measures of semantic processing for auditory attention decoding. *bioRxiv*. doi:10.1101/2020.04.17.046813.

Ding, N., and Simon, J. Z. (2012). Emergence of neural encoding of auditory objects while listening to competing speakers. *Proc. Natl. Acad. Sci.* 109, 11854–11859. doi:10.1073/pnas.1205381109.

Dingemanse, J. G., and Goedegebure, A. (2019). The important role of contextual information in speech perception in cochlear implant users and its consequences in speech tests. *Trends Hear.* 23, 233121651983867. doi:10.1177/2331216519838672.

Donhauser, P. W., and Baillet, S. (2020). Two distinct neural timescales for predictive speech processing. *Neuron* 105, 385-393.e9. doi:10.1016/j.neuron.2019.10.019.

Dryden, A., Allen, H. A., Henshaw, H., and Heinrich, A. (2017). The association between cognitive performance and speech-in-noise perception for adult listeners: a systematic literature review and meta-analysis. *Trends Hear.* 21, 1–21. doi:10.1177/2331216517744675.

Dubno, J. R., Dirks, D. D., and Morgan, D. E. (1984). Effects of age and mild hearing loss on

1248    speech recognition in noise. *J. Acoust. Soc. Am.* 76, 87–96. doi:10.1121/1.391011.

1249    Fiedler, L., Wöstmann, M., Herbst, S. K., and Obleser, J. (2019). Late cortical tracking of ignored
1250    speech facilitates neural selectivity in acoustically challenging conditions. *Neuroimage* 186,
1251    33–42. doi:10.1016/j.neuroimage.2018.10.057.

1252    Fogerty, D., Ahlstrom, J. B., Bologna, W. J., and Dubno, J. R. (2015). Sentence intelligibility
1253    during segmental interruption and masking by speech-modulated noise: effects of age and
1254    hearing loss. *J. Acoust. Soc. Am.* 137, 3487–3501. doi:10.1121/1.4921603.

1255    Fogerty, D., Madorskiy, R., Ahlstrom, J. B., and Dubno, J. R. (2020). Comparing speech
1256    recognition for listeners with normal and impaired hearing: simulations for controlling
1257    differences in speech levels and spectral shape. *J. Speech, Lang. Hear. Res.*, 1–11.
1258    doi:10.1044/2020_JSLHR-20-00246.

1259    Fowler, E. P. (1936). A method for the early detection of otosclerosis: a study of sounds well
1260    above threshold. *Arch. Otolaryngol. - Head Neck Surg.* 24, 731–741.
1261    doi:10.1001/archotol.1936.00640050746005.

1262    Gatehouse, S., and Noble, I. (2004). The speech, spatial and qualities of hearing scale (ssq). *Int.*
1263    *J. Audiol.* 43, 85–99. doi:10.1080/14992020400050014.

1264    Gevins, A., Smith, M. E., McEvoy, L., and Yu, D. (1997). High-resolution EEG mapping of cortical
1265    activation related to working memory: effects of task difficulty, type of processing, and
1266    practice. *Cereb. Cortex* 7, 374–385. doi:10.1093/cercor/7.4.374.

1267    Goossens, T., Vercammen, C., Wouters, J., and van Wieringen, A. (2017). Masked speech
1268    perception across the adult lifespan: impact of age and hearing impairment. *Hear. Res.*
1269    344, 109–124. doi:10.1016/j.heares.2016.11.004.

1270    Heilbron, M., Ehinger, B., Hagoort, P., and de Lange, F. (2019). Tracking naturalistic linguistic
1271    predictions with deep neural language models. in *2019 Conference on Cognitive*
1272    *Computational Neuroscience* (Brentwood, Tennessee, USA: Cognitive Computational
1273    Neuroscience), 424–427. doi:10.32470/CCN.2019.1096-0.

1274    Helfer, K. S., and Wilber, L. A. (1990). Hearing loss, aging, and speech perception in
1275    reverberation and noise. *J. Speech Hear. Res.* 33, 149–155. doi:10.1044/jshr.3301.149.

1276    Huang, Q., and Tang, J. (2010). Age-related hearing loss or presbycusis. *Eur. Arch. Oto-Rhino-*
1277    *Laryngology* 267, 1179–1191. doi:10.1007/s00405-010-1270-7.

1278    Huth, A. G., De Heer, W. A., Griffiths, T. L., Theunissen, F. E., and Gallant, J. L. (2016). Natural
1279    speech reveals the semantic maps that tile human cerebral cortex. *Nature* 532, 453–458.
1280    doi:10.1038/nature17637.

1281    Jutten, C., and Herault, J. (1991). Blind separation of sources, part i: an adaptive algorithm
1282    based on neuromimetic architecture. *Signal Processing* 24, 1–10. doi:10.1016/0165-
1283    1684(91)90079-X.

1284    Kidd, G., Mason, C. R., Best, V., and Marrone, N. (2010). Stimulus factors influencing spatial
1285    release from speech-on-speech masking. *J. Acoust. Soc. Am.* 128, 1965–1978.
1286    doi:10.1121/1.3478781.

1287    Killion, M. C., Niquette, P. A., Gudmundsen, G. I., Revit, L. J., and Banerjee, S. (2004).
1288    Development of a quick speech-in-noise test for measuring signal-to-noise ratio loss in
1289    normal-hearing and hearing-impaired listeners. *J. Acoust. Soc. Am.* 116, 2395–2405.
1290    doi:10.1121/1.1784440.

1291    Kleiner, M., Brainard, D. H., and Pelli, D. G. (2007). What's new in Psychtoolbox-3. in *Perception*

1292    *36 ECVP Abstract Supplement*.

1293    Klem, G. H., Lüders, H. O., Jasper, H. H., and Elger, C. (1999). The ten-twenty electrode system
1294        of the international federation. the international federation of clinical neurophysiology.
1295        *Electroencephalogr. Clin. Neurophysiol. Suppl.* 52, 3–6. Available at:
1296        http://www.ncbi.nlm.nih.gov/pubmed/10590970.

1297    Kong, Y. Y., Mullangi, A., and Ding, N. (2014). Differential modulation of auditory responses to
1298        attended and unattended speech in different listening conditions. *Hear. Res.* 316, 73–81.
1299        doi:10.1016/j.heares.2014.07.009.

1300    Kutas, M., and Hillyard, S. (1980). Reading senseless sentences: brain potentials reflect semantic
1301        incongruity. *Science* 207, 203–205. doi:10.1126/science.7350657.

1302    Kutas, M., and Hillyard, S. A. (1984). Brain potentials during reading reflect word expectancy
1303        and semantic association. *Nature* 307, 161–3. doi:10.1038/307161a0.

1304    Lalor, E. C., and Foxe, J. J. (2010). Neural responses to uninterrupted natural speech can be
1305        extracted with precise temporal resolution. *Eur. J. Neurosci.* 31, 189–193.
1306        doi:10.1111/j.1460-9568.2009.07055.x.

1307    Lash, A., Rogers, C. S., Zoller, A., and Wingfield, A. (2013). Expectation and entropy in spoken
1308        word recognition: effects of age and hearing acuity. *Exp. Aging Res.* 39, 235–253.
1309        doi:10.1080/0361073X.2013.779175.

1310    Lesenfants, D., Vanthornhout, J., Verschueren, E., Decruy, L., and Francart, T. (2019). Predicting
1311        individual speech intelligibility from the neural tracking of acoustic- and phonetic-level
1312        speech representations. *Hear. Res.* 380, accepted. doi:10.1101/471367.

1313    Loughrey, D. G., Kelly, M. E., Kelley, G. A., Brennan, S., and Lawlor, B. A. (2018). Association of
1314        age-related hearing loss with cognitive function, cognitive impairment, and dementia a
1315        systematic review and meta-analysis. *JAMA Otolaryngol. - Head Neck Surg.* 144, 115–126.
1316        doi:10.1001/jamaoto.2017.2513.

1317    Marrone, N., Mason, C. R., and Kidd, G. (2008). The effects of hearing loss and age on the
1318        benefit of spatial separation between multiple talkers in reverberant rooms. *J. Acoust. Soc.*
1319        *Am.* 124, 3064–3075. doi:10.1121/1.2980441.

1320    McAuliffe, M., Socolof, M., Mihuc, S., Wagner, M., and Sonderegger, M. (2017). Montreal
1321        forced aligner: trainable text-speech alignment using kaldi. *Proc. Annu. Conf. Int. Speech*
1322        *Commun. Assoc. INTERSPEECH* 2017-Augus, 498–502. doi:10.21437/Interspeech.2017-
1323        1386.

1324    Mesgarani, N., and Chang, E. F. (2012). Selective cortical representation of attended speaker in
1325        multi-talker speech perception. *Nature* 485, 233–236. doi:10.1038/nature11020.

1326    Mick, P., Kawachi, I., and Lin, F. R. (2014). The association between hearing loss and social
1327        isolation in older adults. *Otolaryngol. Neck Surg.* 150, 378–384.
1328        doi:10.1177/0194599813518021.

1329    Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word
1330        representations in vector space. *arXiv*. Available at: http://arxiv.org/abs/1301.3781.

1331    Mikolov, T., Sutskever, I., Chen, K., Corrado, G., and Dean, J. (2013b). Distributed
1332        representations of words and phrases and their compositionality. in *NIPS*.

1333    Millman, R. E., Mattys, S. L., Gouws, A. D., and Prendergast, G. (2017). Magnified neural
1334        envelope coding predicts deficits in speech perception in noise. *J. Neurosci.* 37, 7727–
1335        7736. doi:10.1523/JNEUROSCI.2722-16.2017.

1336  O'Neill, E. R., Kreft, H. A., and Oxenham, A. J. (2019). Cognitive factors contribute to speech
1337      perception in cochlear-implant users and age-matched normal-hearing listeners under
1338      vocoded conditions. *J. Acoust. Soc. Am.* 146, 195–210. doi:10.1121/1.5116009.
1339  O'Sullivan, J. A., Herrero, J., Smith, E., Sheth, S. A., Mehta, A. D., Mesgarani, N., et al. (2019).
1340      Hierarchical encoding of attended auditory objects in multi-talker speech perception
1341      article hierarchical encoding of attended auditory objects in multi-talker speech
1342      perception. *Neuron*, 1–15. doi:10.1016/j.neuron.2019.09.007.
1343  Obleser, J., and Kotz, S. A. (2011). Multiple brain signatures of integration in the comprehension
1344      of degraded speech. *Neuroimage* 55, 713–723. doi:10.1016/j.neuroimage.2010.12.020.
1345  Onton, J., Delorme, A., and Makeig, S. (2005). Frontal midline EEG dynamics during working
1346      memory. *Neuroimage* 27, 341–356. doi:10.1016/j.neuroimage.2005.04.014.
1347  Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., et al. (2019). PyTorch: an
1348      imperative style, high-performance deep learning library. in *Advances in Neural
1349      Information Processing Systems*, 8026–8037. Available at:
1350      https://papers.nips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf.
1351  Pelli, D. G. (1997). The videotoolbox software for visual psychophysics: transforming numbers
1352      into movies. *Spat. Vis.* 10, 437–442. doi:10.1163/156856897X00366.
1353  Phatak, S. A., Sheffield, B. M., Brungart, D. S., and Grant, K. W. (2018). Development of a test
1354      battery for evaluating speech perception in complex listening environments: effects of
1355      sensorineural hearing loss. *Ear Hear.* 39, 449–456. doi:10.1097/AUD.0000000000000567.
1356  Pickering, M. J., and Gambi, C. (2018). Predicting while comprehending language: a theory and
1357      review. *Psychol. Bull.* 144, 1002–1044. doi:10.1037/bul0000158.
1358  Power, A. J., Foxe, J. J., Forde, E.-J., Reilly, R. B., and Lalor, E. C. (2012). At what time is the
1359      cocktail party? a late locus of selective attention to natural speech. *Eur. J. Neurosci.* 35,
1360      1497–1503. doi:10.1111/j.1460-9568.2012.08060.x.
1361  Presacco, A., Simon, J. Z., and Anderson, S. (2016). Evidence of degraded representation of
1362      speech in noise,in the aging midbrain and cortex. *J. Neurophysiol.* 116, 2346–2355.
1363      doi:10.1152/jn.00372.2016.
1364  Pronk, M., Deeg, D. J. H., Smits, C., Twisk, J. W., Van Tilburg, T. G., Festen, J. M., et al. (2014).
1365      Hearing loss in older persons: does the rate of decline affect psychosocial health? *J. Aging
1366      Health* 26, 703–723. doi:10.1177/0898264314529329.
1367  Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. (2019). Language models
1368      are unsupervised multitask learners. *OpenAI Blog*. Available at:
1369      https://cdn.openai.com/better-language-
1370      models/language_models_are_unsupervised_multitask_learners.pdf.
1371  Ray, J., Popli, G., and Fell, G. (2018). Association of cognition and age-related hearing
1372      impairment in the english longitudinal study of ageing. *JAMA Otolaryngol. - Head Neck
1373      Surg.* 144, 876–882. doi:10.1001/jamaoto.2018.1656.
1374  Rehurek, R., and Sojka, P. (2010). Software framework for topic modelling with large corpora. in
1375      *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, 45–50.
1376      Available at: https://www.fi.muni.cz/usr/sojka/papers/lrec2010-rehurek-sojka.pdf.
1377  Saunders, G. H. (1989). Determinants of objective and subjective auditory disability in patients
1378      with normal hearing. Available at: http://eprints.nottingham.ac.uk/id/eprint/12959.
1379  Smith, S. B., Krizman, J., Liu, C., White-Schwoch, T., Nicol, T., and Kraus, N. (2019). Investigating

1380    peripheral sources of speech-in-noise variability in listeners with normal audiograms. *Hear.*
1381        *Res.* 371, 66–74. doi:10.1016/j.heares.2018.11.008.
1382    Steinberg, J. C., and Gardner, M. B. (1937). The dependence of hearing impairment on sound
1383        intensity. *J. Acoust. Soc. Am.* 9, 11–23. doi:10.1121/1.1915905.
1384    Tremblay, K. L., Pinto, A., Fischer, M. E., Klein, B. E. K., Klein, R., Levy, S., et al. (2015). Self-
1385        reported hearing difficulties among adults with normal audiograms: the beaver dam
1386        offspring study. *Ear Hear.* 36, e290–e299. doi:10.1097/AUD.0000000000000195.
1387    van Rooij, J. C. G. M., and Plomp, R. (1990). Auditive and cognitive factors in speech perception
1388        by elderly listeners. ii: multivariate analyses. *J. Acoust. Soc. Am.* 88, 2611–2624.
1389        doi:10.1121/1.399981.
1390    Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017).
1391        Attention is all you need. in *NIPS* Available at: http://arxiv.org/abs/1706.03762.
1392    Weissbart, H., Kandylaki, K. D., and Reichenbach, T. (2019). Cortical tracking of surprisal during
1393        continuous speech comprehension. *J. Cogn. Neurosci.*, 1–12. doi:10.1162/jocn_a_01467.
1394    Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., et al. (2020). Transformers :
1395        state-of-the-art natural language processing. in *Proceedings of the 2020 Conference on*
1396        *Empirical Methods in Natural Language Processing: System Demonstrations*, 38–45.
1397        Available at: https://www.aclweb.org/anthology/2020.emnlp-demos.6.pdf.
1398    Xia, J., Kalluri, S., Micheyl, C., and Hafter, E. (2017). Continued search for better prediction of
1399        aided speech understanding in multi-talker environments. *J. Acoust. Soc. Am.* 142, 2386–
1400        2399. doi:10.1121/1.5008498.
1401    Yamasoba, T., Lin, F. R., Someya, S., Kashio, A., Sakamoto, T., and Kondo, K. (2013). Current
1402        concepts in age-related hearing loss: epidemiology and mechanistic pathways. *Hear. Res.*
1403        303, 30–38. doi:10.1016/j.heares.2013.01.021.
1404    Zan, P., Presacco, A., Anderson, S., and Simon, J. Z. (2020). Exaggerated cortical representation
1405        of speech in older listeners: mutual information analysis. *J. Neurophysiol.* 124, 1152–1164.
1406        doi:10.1152/jn.00002.2020.
1407    Zhao, F., and Stephens, D. (2007). A critical review of king-kopetzky syndrome: hearing
1408        difficulties, but normal hearing? *Audiol. Med.* 5, 119–124.
1409        doi:10.1080/16513860701296421.

1410

1411

1412    **Tables**

1413

1414    Table 1. Mixed-factors ANOVA results. Significant F-statistic values are bold, with levels of
1415    significance: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

1416

| | df | F | $\eta_p^2$ |
|---|---|---|---|
| Attention | 1, 37 | **34.3*** | 0.48 |
| Feature | 2, 74 | **18.5*** | 0.33 |
| ROI | 1, 37 | **8.9** | 0.19 |

| | | | |
|---|---|---|---|
| Age | 1, 37 | **7.92\*\*** | 0.18 |
| Attention × Age | 1, 37 | **7.6\*\*** | 0.17 |
| Feature × Age | 2, 74 | **4.1\*** | 0.10 |
| ROI × Age | 1, 37 | **7.2\*** | 0.16 |
| Attention × Feature | 1.8, 66.6 | **8.5\*\*** | 0.19 |
| Attention × ROI | 1, 37 | 2 | 0.05 |
| Feature × ROI | 1.7, 62.2 | 3.2 | 0.08 |
| Attention × Feature × Age | 2, 74 | 2 | 0.05 |
| Attention × ROI × Age | 1, 37 | 1.6 | 0.04 |
| Feature × ROI × Age | 2, 74 | 0.7 | 0.02 |
| Attention × Feature × ROI | 2, 74 | **9.7\*\*\*** | 0.21 |
| Attention × Feature × ROI × Age | 2, 74 | 2.4 | 0.06 |

1417

1418