

1     **Electrophysiological indices of hierarchical speech**  
2     **processing differentially reflect the comprehension of**  
3     **speech in noise**

4     Abbreviated title: EEG speech indices differentially reflect comprehension

5     Shyanthony R. Synigal<sup>1</sup>, Andrew J. Anderson<sup>2</sup>, Edmund C. Lalor<sup>1,2</sup>

6  
7     <sup>1</sup>Department of Biomedical Engineering, University of Rochester, 201 Robert B. Goergen Hall,  
8     P.O. Box 270168, Rochester, NY 14627, USA.

9     <sup>2</sup>Department of Neuroscience and Del Monte Institute for Neuroscience, University of Rochester,  
10     201 Robert B. Goergen Hall, P.O. Box 270168, Rochester, NY 14627, USA.

11     Corresponding Author: Edmund C. Lalor PhD, Department of Biomedical Engineering, University  
12     of Rochester, 201 Robert B. Goergen Hall, P.O. Box 270168, Rochester, NY 14627, USA;  
13     [Edmund.Lalor@urmc.rochester.edu](mailto:Edmund.Lalor@urmc.rochester.edu)

14

15     Number of pages: 38

16     Number of figures: 6

17     Number of tables: 3

18     Number of words for abstract: 230

19     Number of words for introduction: 660

20     Number of words for discussion: 1,492

21

22     **Conflict of interest statement**

23     The authors declare no competing financial interests.

24

25     **Acknowledgements and Support**

26     This work was supported by a grant from the Simons Foundation Autism Research Initiative  
27     (SFARI). Additional support was provided by the Del Monte Institute for Neuroscience. The  
28     authors thank Ms. Xueying Wang and Dr. Aaron Nidiffer for some assistance with data collection,  
29     and Dr. Aaron Nidiffer and Dr. Madeline Cappelloni for helpful comments on the manuscript.

30

31 **ABSTRACT**

32           The past few years have seen an increase in the use of encoding models to explain neural  
33 responses to natural speech. The goal of these models is to characterize how the human brain  
34 converts acoustic speech energy into different linguistic representations that enable everyday  
35 speech comprehension. For example, researchers have shown that electroencephalography  
36 (EEG) data can be modeled in terms of acoustic features of speech, such as its amplitude  
37 envelope or spectrogram, linguistic features such as phonemes and phoneme probability, and  
38 higher-level linguistic features like context-based word predictability. However, it is unclear how  
39 reliably EEG indices of these different speech representations reflect speech comprehension in  
40 different listening conditions. To address this, we recorded EEG from neurotypical adults who  
41 listened to segments of an audiobook in different levels of background noise. We modeled how  
42 their EEG responses reflected different acoustic and linguistic speech features and how this  
43 varied with speech comprehension across noise levels. In line with our hypothesis, EEG  
44 signatures of context-based word predictability and phonetic features were more closely  
45 correlated with behavioral measures of speech comprehension and percentage of words heard  
46 than EEG measures based on low-level acoustic features. EEG markers of the influence of top-  
47 down, context-based prediction on bottom-up acoustic processing also correlated with behavior.  
48 These findings help characterize the relationship between brain and behavior by comprehensively  
49 linking hierarchical indices of neural speech processing to language comprehension metrics.

50

51

52

53

54

55

56

57 **SIGNIFICANCE STATEMENT**

58           Acoustic and linguistic features of speech have been shown to be consistently tracked by  
59 neural activity even in noisy conditions. However, it is unclear how signatures of low- and high-  
60 level features covary with one another and relate to behavior across these listening conditions.  
61 Here, we find that categorical phonetic feature processing is more affected by noise than acoustic  
62 and word probability-based speech features. We also find that phonetic features and word  
63 probability-based features better correlate with measures of intelligibility and comprehension.  
64 These results extend our understanding of how various speech features are comparatively  
65 reflected in electrical brain activity and how they relate to perception in challenging listening  
66 conditions.

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

## 83 INTRODUCTION

84           Given the importance of speech communication in human life, tremendous amounts of  
85 research have focused on characterizing the neurophysiology of language comprehension  
86 (Hickok 2015). This research has revealed a network of brain areas that are functionally  
87 specialized for processing different hierarchical levels of speech and language (Hickok and  
88 Poeppel 2007). For example, work has shown that low level acoustic and spectrotemporal  
89 features of speech are chiefly processed in early auditory cortex (de Heer et al. 2017), with various  
90 phonological features being processed in secondary areas like the superior temporal gyrus  
91 (Hamilton, Edwards, and Chang 2018; Mesgarani et al. 2014) and some prefrontal areas (Burton  
92 2009; de Heer et al. 2017), and meaning being represented across large areas of cortex (Huth et  
93 al. 2016; Anderson et al. 2017; Pereira et al. 2018).

94           While much of the above knowledge has been obtained from functional neuroimaging and  
95 invasive recordings in neurosurgical patients, parallel efforts have been made to obtain  
96 noninvasive magneto- and electrophysiological (MEG/EEG) markers reflecting hierarchical  
97 speech features. This includes modeling how EEG and MEG track the amplitude envelope of  
98 natural speech (Lalor and Foxe 2010a) and how neural responses reflect the spectrotemporal (Di  
99 Liberto, O'Sullivan, and Lalor 2015; Daube, Ince, and Gross 2019), phonetic (Di Liberto,  
100 O'Sullivan, and Lalor 2015), phoneme-level probability (Di Liberto et al. 2019; Gwilliams et al.  
101 2020; Brodbeck, Hong, and Simon 2018), lexical (Heilbron et al. 2022), prosodic (Teoh,  
102 Cappelloni, and Lalor 2019), and semantic (Heilbron et al. 2022; Broderick et al. 2018) features  
103 of natural speech.

104           Some advantages of EEG are that it is significantly cheaper and easier to use in applied  
105 research in different cohorts (Peck et al. 2021; Salisbury et al. 2002). Consequently, there has  
106 been considerable interest in exploring how different EEG markers of speech processing reflect  
107 speech intelligibility (Verschueren, Vanthornhout, and Francart 2021) and language  
108 comprehension (Broderick et al. 2022; Ahissar et al. 2001). Many of these studies altered the

109 intelligibility or comprehensibility of speech by adding background noise (Iotzov and Parra 2019)  
110 or by degrading the speech signal itself (Viswanathan et al. 2021). For example, some work has  
111 shown that cortical tracking of the speech envelope decreases as noise levels increase (Etard  
112 and Reichenbach 2019; Lesenfants et al. 2019; Vanthornhout et al. 2018; Zou et al. 2019),  
113 although others have suggested such tracking remains robust until the background noise is more  
114 than twice as loud as the speech it masks (Ding and Simon 2013). Meanwhile, experiments that  
115 have explored EEG indices of semantic processing appear to show a strong correlation with  
116 speech intelligibility and/or understanding (Broderick et al. 2018). Very few studies, however, have  
117 systematically explored how EEG indices of both low- and high-level speech processing covary  
118 across different levels of speech comprehension (Strauß et al. 2022; Yasmin et al. 2023). This is  
119 the goal of the present study.

120         We explore how EEG markers of speech envelope, spectrogram, acoustic onsets,  
121 phonetic features, and lexical surprisal processing vary with subjective measures of the  
122 percentage of words heard (as a proxy for intelligibility) and objective measures of comprehension  
123 across different background noise conditions. We hypothesize that EEG measures of higher-level  
124 processing (e.g., lexical surprisal and phonetic features) will more strongly correlate with behavior  
125 than lower-level measures (e.g., envelope tracking). We also test how listeners exploit linguistic  
126 context to process noisy speech and how that effect might manifest in EEG. Here, we leverage a  
127 recently introduced measure of predictive speech perception that quantifies how the tracking of  
128 low-level speech features varies as a function of the context-based semantic content of that  
129 speech (Broderick, Anderson, and Lalor 2019), but this time using lexical surprisal. We  
130 hypothesize that this measure strengthens for speech in moderate levels of noise (when speech  
131 is still intelligible) relative to speech in quiet, before falling off at high levels of background noise  
132 (when speech is no longer intelligible). With this study we seek to extend our understanding of  
133 the hierarchical processing of continuous speech under a range of realistic listening conditions.

134

## 135 **METHODS**

### 136 ***Participants***

137           28 healthy adults (9 males, 18-35 years old) participated in this study. One subject was  
138 excluded due to an insufficient amount of data and two were excluded due to technical issues,  
139 resulting in a dataset of 25 participants. Each participant provided written informed consent and  
140 reported having normal hearing, normal or corrected-to-normal vision, English as their first and  
141 main language, and no history of neurological disorders. Participants were also compensated for  
142 their participation. All procedures were approved by the University of Rochester Human Subjects  
143 Review Board.

### 144 ***Stimuli and experimental procedure***

145           Participants listened to 70 minutes of *A Wrinkle in Time* by Madeleine L'Engle which was  
146 read by an American female speaker. Each trial was one minute long and was presented at one  
147 of five noise levels: quiet (no noise) and +3 dB, -3 dB, -6 dB, and -9 dB signal-to-noise ratios  
148 (SNRs). The background noise was spectrally matched stationary noise, which was estimated  
149 from the clean speech using a 46<sup>th</sup> order forward linear predictive model. The prediction order  
150 was calculated based on the sampling rate of the audio clips (Crosse, Di Liberto, and Lalor 2016;  
151 Ding and Simon 2013). There were 14 minutes' worth of audio for each of the five noise  
152 conditions. The storyline was preserved from trial to trial, but the conditions were  
153 pseudorandomized such that no noise level occurred consecutively. Participants rated how many  
154 words they heard (on a scale of 0-100%) and answered two multiple choice comprehension  
155 questions after each trial. The comprehension questions used here were the same questions  
156 used from a previous study (Maddox and Lee 2018), except we presented only two out of the four  
157 original questions created for each trial. The stimuli were presented through Sennheiser HD650  
158 headphones at a sampling rate of 44.1 kHz using Psychtoolbox (Kleiner, Brainard, and Pelli 2007)  
159 and custom MATLAB scripts (MATLAB 2019).

160

## 161 ***Data acquisition and preprocessing***

162 EEG data were recorded from 128 scalp electrodes (plus two mastoid channels that were  
163 not analyzed in this work). The data were acquired at a 1024 Hz sampling rate with the BioSemi  
164 Active Two system. The data were preprocessed using the PREP pipeline and its default  
165 parameters (Bigdely-Shamlo et al. 2015). This pipeline first used detrending to high pass filter the  
166 data at 1 Hz followed by 60 Hz line noise removal. Afterwards, robust re-referencing was applied  
167 which allows the data to be referenced to an average of all channels except those contaminated  
168 with noise. This function identifies and interpolates noisy channels in an iterative manner such  
169 that the re-referencing itself is not affected by the noise. The cleaned data was then low pass  
170 filtered at 8 Hz, using a filter with an 8.5 Hz cutoff frequency and 80 dB stopband attenuation.  
171 Next, the data were epoched and independent component analysis (ICA) was applied using  
172 EEGLAB's picard function (Delorme and Makeig 2004; Pion-Tonachini, Kreutz-Delgado, and  
173 Makeig 2019) to remove muscle and eye artifacts. Lastly, the data were downsampled to 128 Hz.

## 174 ***Speech stimulus characterization***

175 Speech is organized in a hierarchical manner where sounds can form syllables, syllables  
176 form words, words form sentences, and so on. To assess how our brains might concurrently  
177 process speech across levels of this hierarchy, we chose to model EEG responses to speech  
178 based on several different representations, all of which were computed on the clean versions of  
179 each trial.

180 *Envelope.* We first calculated the speech envelope, a well-established feature shown to  
181 be robustly tracked by cortical activity (Aiken and Picton 2008; Destoky et al. 2019; Di Liberto,  
182 O'Sullivan, and Lalor 2015; Ding and Simon 2013; Etard and Reichenbach 2019; Lalor and Foxe  
183 2010b; Nourski et al. 2009; Pasley et al. 2012) and to be important for speech recognition and  
184 intelligibility (Ahissar et al. 2001; Drullman, Festen, and Plomp 1994; Shannon et al. 1995). The  
185 speech signal was first lowpass filtered at 20 kHz (22.05 kHz cutoff frequency, 1 dB passband  
186 attenuation, 60 dB stopband attenuation). The broadband speech envelope was calculated using

187 a gammachirp auditory filterbank to mimic the filtering properties of the cochlea (Irino and  
188 Patterson 2006). This filterbank was used to filter the speech into 16 bands from 250 Hz to 8 kHz  
189 with an equal loudness contour (essentially creating a spectrogram). Lastly, the frequency bands  
190 were averaged together.

191 *Acoustic onsets and spectrogram.* We chose to model two additional acoustic features,  
192 acoustic onsets and spectrogram, which were shown to be reflected in cortical activity above and  
193 beyond the speech envelope (Brodbeck et al. 2020; Di Liberto, O'Sullivan, and Lalor 2015;  
194 Sohoglu and Davis 2020). Acoustic onsets were approximated by computing the first derivative  
195 of the speech envelope and then half-wave rectifying the result. A 16-band spectrogram was  
196 calculated using the same filterbank and parameters as the speech envelope, just without the  
197 final averaging step.

198 *Phonetic features.* To calculate phonetic features, the Montreal Forced Aligner (McAuliffe  
199 et al. 2017) was first used to partition and time align each word in the story into phonemes  
200 according to the International Phonetic Alphabet for American English. Then, each phoneme was  
201 linearly mapped onto a set of 19 binary phonetic features based on the University of Iowa's  
202 phonetics project (<http://www.uiowa.edu/~acadtech/phonetics/english/english.html/>).

203 *Lexical surprisal.* Lastly, we calculated the surprisal of each word based on its preceding  
204 context using the Transformer-XL model (Dai et al., 2019). This model contains a recurrence  
205 mechanism that allows it to build and reuse memory from previous segments and learn longer-  
206 term dependencies, while preserving the temporal information of previous word embeddings. This  
207 model was chosen because it can predict the probability of an upcoming word using the context  
208 from all preceding words. The softmax of the values from the output layer of the model were taken  
209 to estimate the probability of each word, and the negative log of a word's probability was computed  
210 to estimate lexical surprisal (Dai et al. 2019).

211

212



## 213 ***Modeling the relationship between speech features and EEG responses***

214 One goal of the present study was to find how the encoding of individual speech features  
215 changes with SNR and find how those changes relate to comprehension and what participants  
216 reported hearing. To index the encoding of individual speech features, we used a forward model  
217 which acted as a filter or kernel that described the transformation from those features to the EEG  
218 responses recorded at each electrode. Here, we modeled acoustic onsets, the spectrogram,  
219 phonetic features, and lexical surprisal. An additional vector with impulses placed at each word's  
220 onset was included to capture any acoustic related onset responses that the acoustic onset  
221 predictor may have missed. All five features were modeled together to control for variance  
222 explained by the competing features and to find which feature best explained certain EEG  
223 responses (Brodbeck, Presacco, and Simon 2018; Gillis et al. 2021; Brodbeck, Hong, and Simon  
224 2018).

225 The data from each of the five experimental (SNR) conditions were modeled separately  
226 using 14-fold leave-one-out cross-validation and ridge regression. Each feature was normalized  
227 between 0-1 and the EEG were z-scored. The features were then concatenated and partitioned  
228 into training and test sets. The stimuli were lagged from -100-700ms to capture both short and  
229 long latency responses to acoustic and linguistic features. Cross-validation was conducted on the  
230 clean condition to select the optimal regularization parameter,  $\lambda$ , which ranged from  $10^{-1}$ – $10^8$ . We  
231 identified the regularization parameter that resulted in the highest prediction accuracy for each  
232 individual test fold. We then selected the parameter that produced the highest reconstruction  
233 accuracy most often (across all test folds) so that we could use one parameter to train the models  
234 for each condition. Using the same parameter for all folds and conditions (within a participant)  
235 allowed for a fairer comparison of model performance since each participant's trials would be on  
236 the same scale and it minimized model overfitting.

237 A temporal response function (TRF),  $w(\tau, n)$ , was trained using the selected regularization  
238 parameter and the training data to predict the neural responses,  $r(t, n)$ , from the set of

239 concatenated speech features,  $s(t - \tau)$ . Then, we separated the model into the segments that  
240 corresponded to acoustic onsets, spectrogram, phonetic features, and surprisal. The forward  
241 modeling procedure can be expressed as follows, including the residual response,  $\varepsilon(t, n)$ , not  
242 explained by the model:

$$243 \quad r(t, n) = \sum_{\tau} w(\tau, n) s(t - \tau) + \varepsilon(t, n)$$

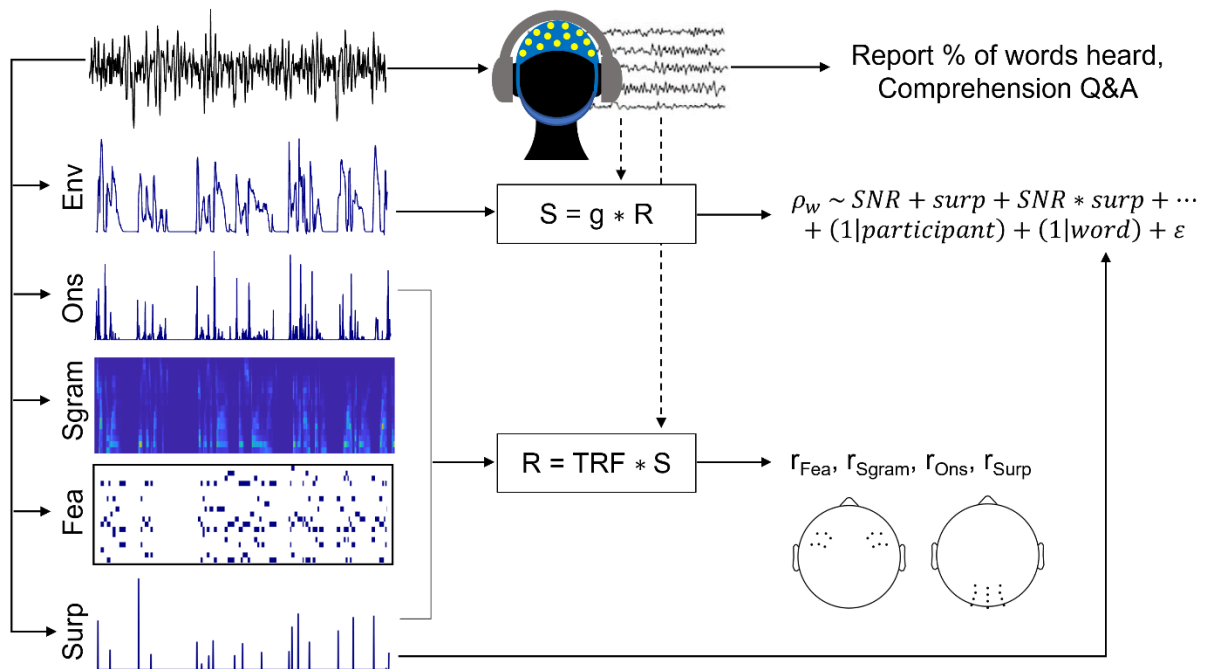
244 Each model segment was tested on the held-out data, and its performance was assessed  
245 by quantifying the correlation between the predicted EEG and the actual EEG. Forward model  
246 performance, or prediction accuracy, was averaged across all folds. The results from the acoustic  
247 onset, spectrogram, and phonetic feature model segments were then averaged across 12 well-  
248 predicted frontotemporal electrodes (6 bilaterally symmetric pairs), and the results from the  
249 surprisal model were averaged across 12 well-predicted parieto-occipital channels. There may be  
250 some positive bias in the prediction accuracies (i.e., higher accuracies) since we have selected  
251 well predicted electrodes based on the current dataset, but we assumed that this bias is present  
252 in all conditions and does not vary systematically across conditions.

253 Backward modeling, where EEG is used to reconstruct an estimate of the speech  
254 envelope was also employed to enable the current results to be directly referenced to other  
255 studies in which backward modeling is more common. The backward model or decoder,  $g(\tau, n)$ ,  
256 describes the transformation from lagged EEG responses at all electrodes,  $r(t + \tau, n)$ , to an  
257 estimate of the speech envelope,  $\hat{s}(t)$ . As detailed elsewhere (Crosse et al. 2016), the modeling  
258 procedure can be expressed as:

$$259 \quad \hat{s}(t) = \sum_n \sum_{\tau} r(t + \tau, n) g(\tau, n)$$

260 Similar to the forward models, this analysis was conducted separately for each  
261 experimental (SNR) condition using 14-fold leave-one-out cross-validation. First, the stimuli and  
262 responses were normalized and then partitioned into train and test sets. Here the EEG data were

263 lagged,  $\tau$ , from -100–300ms since we were interested in short latency responses to the speech  
264 envelope,  $s(t)$ . We employed the same method of cross-validation and regularization parameter  
265 ( $10^{-1}$ – $10^8$ ) selection as before. A decoder was trained using the selected regularization parameter  
266 and the training data to reconstruct an estimate of the speech envelope. After testing the decoder  
267 on held-out data, we assessed model performance by computing the correlation between the  
268 actual speech envelope and the reconstructed speech envelope. This model performance, also  
269 known as reconstruction accuracy, was averaged across folds and participants (**Figure 1**).



270  
271 **Figure 1.** Methods. EEG data were recorded while participants listened to an audiobook in different levels  
272 of noise. Forward modeling was used to estimate EEG responses ( $R$ ) from the clean representations of the  
273 acoustic onsets, spectrogram, phonetic features, and word surprisal ( $S$  being a concatenation of the  
274 features). Model performance ( $r$ ) was assessed by calculating the correlation between the predicted EEG  
275 and the actual EEG and then averaged across the selected channels. Backward modeling was also used  
276 to reconstruct an estimate of the clean speech envelope. Model performance (reconstruction accuracy,  $\rho_w$ )  
277 was assessed by calculating the correlation between the original speech envelope and the reconstructed

278 speech envelope. A linear mixed-effects model (LME) was then used to determine the influence of word  
279 surprisal (*surp*) on envelope reconstruction accuracy.

### 280 ***Assessing the role of context on acoustic encoding in different levels of background noise***

281 Another major goal of the present study was to test how listeners might rely more on  
282 context to encode speech that is masked by moderate levels of background noise. To do this, we  
283 used a variant of a recently introduced approach that involves quantifying how the tracking of low-  
284 level speech features varies as a function of the context-based semantic content of that speech  
285 (Broderick, Anderson, and Lalor 2019). Specifically, we used a linear mixed-effects (LME) model  
286 to explore the extent to which word predictability in the form of word surprisal influences how the  
287 envelope of that word was reflected in EEG and how that influence changes across SNRs. Using  
288 an LME in this way, one can measure the relationship between the main variable(s) of interest  
289 while controlling for variability caused by random factors. We used the *lmerTest* (version 3.1-3)  
290 and *lme4* (version 1.1-30) packages in R to model the following equation:

$$291 \quad \rho_w \sim 1 + SNR + surp + envStd + f_{rel} + res + f_{rel} * envStd + f_{rel} * res + surp * envStd \\ 292 \quad \quad \quad + surp * res + surp * SNR + (1|participant) + (1|word)$$

293 The dependent variable is word reconstruction accuracy which was calculated as the  
294 Spearman's correlation between the actual word envelope and the predicted word envelope for  
295 the first 100 ms of each word. The independent variables are SNR, lexical surprisal (*surp*),  
296 envelope variability (*envStd*), relative frequency (*f<sub>rel</sub>*), resolvability (*res*), and various interactions  
297 such as the interaction between surprisal and SNR.

298 Envelope variability, relative pitch, and resolvability were selected as nuisance  
299 regressors. This was because these measures can correlate with surprisal, with one another, and  
300 with envelope tracking; so, they are included here to ensure that they aren't inherently  
301 contaminating the lexical surprisal effects. Relative pitch is pitch normalized to the vocal range of  
302 the speaker (Tang, Hamilton, and Chang 2017) and resolvability measures whether a sound's  
303 harmonics are processed between distinct (resolved) or within the same (unresolved) filters of the

304 cochlea (Shackleton and Carlyon 1994). Prosodic cues such as relative pitch and resolvability  
305 can uniquely predict EEG activity even after accounting for other acoustic and phonetic features  
306 (Teoh, Cappelloni, and Lalor 2019).

307         Relative pitch was extracted using Praat (Boersma and Weenink 2013). Once the software  
308 estimates absolute pitch, this result is then z-scored, resulting in relative pitch. Resolvability was  
309 extracted using custom scripts based on a model of the human auditory periphery (McDermott  
310 and Simoncelli 2011; Teoh, Cappelloni, and Lalor 2019). Envelope variability was shown to  
311 correlate with relative pitch and resolvability and all three features have been shown to influence  
312 envelope reconstruction accuracy, so we too included these features to control for acoustic  
313 related changes in the speaker's voice (Broderick, Anderson, and Lalor 2019). Envelope  
314 variability is represented as the standard deviation of the speech envelope.

315         The LME model also included by-word and by-participant random intercepts as some  
316 words may be easier to reconstruct than others and some participants may, on average, have  
317 higher reconstructions than other participants. No random slopes were included, as they caused  
318 the model to not converge even with the addition of an optimizer. Like Broderick et al. 2019, word  
319 reconstruction accuracy and the nuisance regressors were measured in the first 100ms following  
320 each word's onset (Broderick, Anderson, and Lalor 2019).

### 321 ***Statistical analyses***

322         All statistical analyses were performed in R (version 4.2.0) and in MATLAB R2021b  
323 (MATLAB 2021). Due to the skewed distribution of the behavioral results, comparisons were  
324 calculated using a nonparametric Friedman's test, followed by a Wilcoxon Rank Sum test. All  
325 corrections for multiple comparisons were performed using false discovery rate (FDR), specifically  
326 Benjamini & Yekutieli (BY) correction, unless otherwise stated. FDR (BY) corrected pairwise t-  
327 tests were used to determine differences in envelope reconstruction accuracy and EEG prediction  
328 accuracy between conditions. Permutation testing was performed to test the significance of the  
329 EEG-speech model predictions. A null model was created for each SNR by shuffling the stimulus

330 of interest between trials (except for the surprisal vectors which were shuffled within trial) and  
331 calculating a new model with each shuffle. This procedure, including regularization, was repeated  
332 30 times. The null prediction accuracies were averaged across folds, permutations, and  
333 electrodes to result in one value per person. Pairwise t-tests were performed between the actual  
334 and null prediction accuracies across participants.

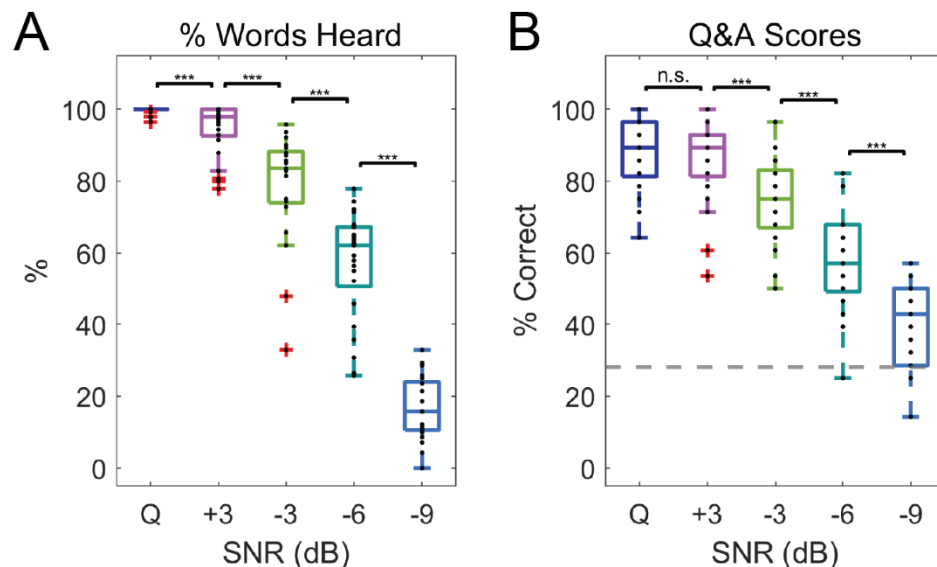
335 Unaggregated prediction accuracies (each accuracy for each trial per person) were used  
336 in an LME model (lmerTest version 3.1-3 and lme4 version 1.1-30 in R) to test how prediction  
337 accuracies changed across SNRs. This model included SNR, speech feature type, and the  
338 interaction between the two as fixed effects, and a by-participant random intercept. We then  
339 computed estimated marginal means (emmeans version 1.8.2 in R) with Tukey adjustment to  
340 perform multiple comparisons tests between each feature. Each LME model in this study was  
341 calculated using the default parameters which included fitting the models with restricted maximum  
342 likelihood (REML) and using Satterthwaite's method for the t-tests. LME models were also used  
343 to model the relationship between behavior and reconstruction/prediction accuracy. Permutation  
344 tests were used to test the significance of the final LME analysis where we tested the relationship  
345 between lexical surprisal and envelope tracking. Surprisal values were shuffled 5000 times while  
346 all other variables remained fixed, and an LME model was calculated for each shuffle. We  
347 calculated the proportion of coefficients that were greater than the observed values. All data and  
348 scripts are available upon request.

## 349 **RESULTS**

### 350 ***Speech perception decreases as SNR decreases***

351 Behavioral scores were collected to show how participants' perception of the story  
352 changed as listening conditions became more challenging. After each one-minute-long trial,  
353 participants rated an estimate of the number of words they were able to hear on a scale from 0-  
354 100% and answered two multiple-choice comprehension questions, each of which had four  
355 possible answers. As expected, there was a significant reduction in percentage of words heard

356 ( $X^2(4) = 98.894$ ,  $p = 2.20 \times 10^{-16}$ ) and comprehension scores ( $X^2(4) = 83.44$ ,  $p = 2.20 \times 10^{-16}$ ,  
357 Friedman test followed by Wilcoxon rank sum test with FDR correction (BY), **Figure 2, Table 1**)  
358 with the addition of noise. After the experiment, participants reported hearing very few words in  
359 the -9 dB SNR condition but that they had attempted to use context clues from previous trials to  
360 answer the questions in the -9 dB SNR trials. To control for this potential confound in our measure  
361 of comprehension, we used a procedure similar to Orf and colleagues, where nine additional  
362 volunteers answered the comprehension questions for each trial without listening to the story (Orf  
363 et al. 2022). Based on the performance of those new participants, we determined a new empirical  
364 chance level of 28% rather than 25%. All participants performed above chance in the quiet, +3  
365 dB SNR, and -3 dB SNR conditions. Comprehension scores were not above chance for five  
366 participants in the -6 dB SNR condition ( $p = 0.066$  and above) and 13 participants in the -9 dB  
367 SNR condition ( $p = 0.066$  and above). Although participants reported hearing fewer words in the  
368 +3 dB SNR condition compared to speech in quiet, they performed similarly in terms of  
369 comprehension ( $p = 1.000$ ).



370

371 **Figure 2.** Behavioral results. **A.** Average percentage of words participants reported hearing in each

372 condition. **B.** Average percentage of correctly answered comprehension questions for each condition. The

373 dotted gray line is the chance level at 28% which was calculated using a separate set of participants who  
374 answered the comprehension questions without listening to the audiobook. For both plots, significance is  
375 indicated by \* if  $p < 0.05$ , \*\* if  $p < 0.01$ , and \*\*\* if  $p < 0.001$  using pairwise Wilcoxon rank sum tests. The  
376 black points in each are individual participants.

377 **Table 1.** Wilcoxon rank sum test results comparing behavioral scores between all conditions.

% Words Heard					Q&A Scores				
	quiet	+3 dB	-3 dB	-6 dB		quiet	+3 dB	-3 dB	-6 dB
+3 dB	0.00093	-	-	-	+3 dB	1.00000	-	-	-
-3 dB	4.2e-05	4.2e-05	-	-	-3 dB	0.00020	0.00080	-	-
-6 dB	4.2e-05	4.2e-05	4.2e-05	-	-6 dB	8.5e-05	8.5e-05	0.00014	-
-9 dB	4.2e-05	4.2e-05	4.2e-05	4.2e-05	-9 dB	8.5e-05	8.5e-05	8.5e-05	0.00027

378

### 379 ***Hierarchical feature encoding declines across SNRs***

380 Studies have typically modeled brain responses to isolated speech features. Recent work,  
381 however, is increasingly beginning to model acoustic and linguistic features simultaneously  
382 (Brodbeck, Hong, and Simon 2018; Brodbeck, Presacco, and Simon 2018; de Heer et al. 2017;  
383 Gillis et al. 2021; Heilbron et al. 2022) to find how the brain uniquely encodes a feature of interest  
384 when accounting for others, as some features may be correlated and explain similar neural activity  
385 (Daube, Ince, and Gross 2019). Since few studies have modeled the encoding of simultaneous  
386 features in challenging listening conditions (Brodbeck et al. 2020), we were interested in how a  
387 range of acoustic and linguistic features were encoded in noise.

388 To test this, we first fit one complete model comprised of acoustic onsets, the speech  
389 spectrogram, phonetic features, and word surprisal. Word onset (not pictured in the subsequent  
390 figures) was also included to ensure the surprisal measure was not reflecting variance that would  
391 be better explained by onset responses to individual words. The full model was separated into its  
392 constituent model pieces (where, for example, constituent models may have one feature in the  
393 case of surprisal, or many features in the case of phonemes) and tested on left out data. The  
394 acoustic onset, spectrogram, and phonetic feature prediction accuracies were averaged across



395 12 well-predicted frontotemporal electrodes and the surprisal prediction accuracies were  
396 averaged across 12 well-predicted parieto-occipital electrodes (electrodes shown in **Figure 1**).

397 Each model performed above chance for each SNR ( $p < 0.01$ , permutations followed by  
398 one-tailed, paired t-tests). As expected, the prediction accuracies for each constituent model  
399 decreased with increasing levels of noise (**Figure 3A**, paired t-tests with FDR [BY] correction).  
400 We then elected to run an LME model analysis exploring how the encoding of different speech  
401 features falls off with decreasing SNR. This helps account for the fact that EEG prediction  
402 accuracies can vary greatly between subjects (based on, for example, cortical folding or  
403 skull/scalp geometry). We fit the following LME model:

$$404 \quad r \sim SNR + Feature + SNR * Feature + (1|participant)$$

405 Since we were interested in trends across noise levels, SNR was treated as a continuous  
406 variable rather than a categorical variable. EEG predictions based on the speech spectrogram  
407 had the shallowest slope across noise conditions ( $\beta_{\text{sgram}} = -0.009$  which is the difference between  
408 SNR and SNR:Sgram, **Table 2**)—suggesting that the spectrogram was relatively well reflected in  
409 the EEG even as noise levels increased. This trend in prediction accuracy, however, was similar  
410 to both the acoustic onset ( $p = 0.708$ ) and surprisal trends ( $p = 0.395$ ). The spectrogram trend  
411 was significantly greater than the phonetic feature trend ( $p = 4 \times 10^{-4}$ ), suggesting that phonetic  
412 feature representations are more sensitive to noise. Altogether, high- and low-level features, with  
413 the exception of phonetic features in some cases, declined similarly across SNRs.

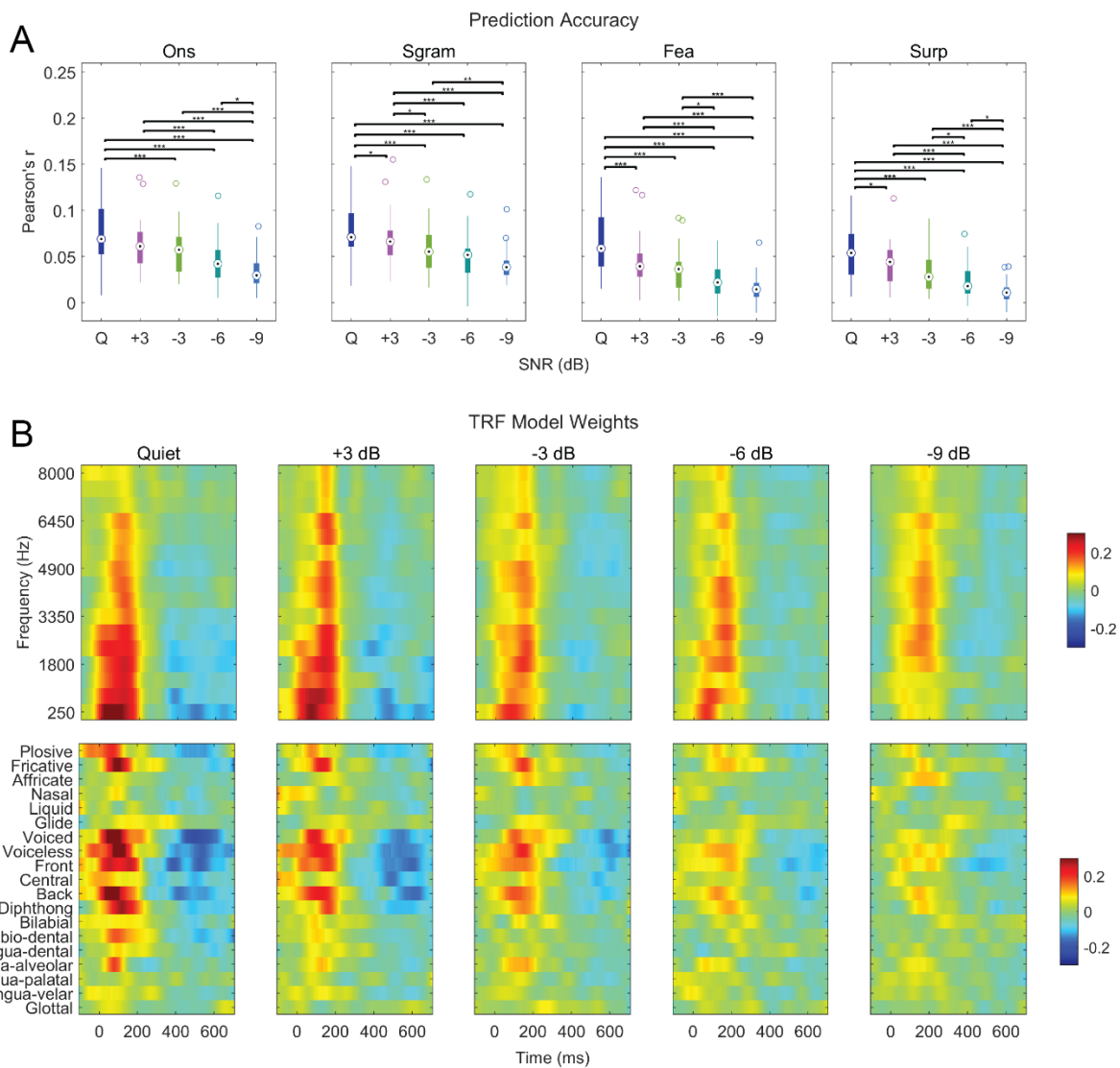
414 Another way to examine the effect of SNR on speech feature encoding is to visualize the  
415 TRF models themselves. In fact, recent work has shown that the amplitude and latency of acoustic  
416 onset and semantic TRFs are influenced by speech SNR (Yasmin et al. 2023). Given that the  
417 encoding of specific frequency bands or phonetic properties may decrease differentially with  
418 noise, we examined the change in the spectrogram and phonetic feature TRF weights at each  
419 SNR. In quiet, we saw the strongest TRF weights in the lower frequency bands, suggesting the  
420 importance of low frequency spectral tracking when no noise is present. Overall, we found a

421 decrease and narrowing in time of the TRF weights across a broad range of frequencies as SNR  
 422 decreased. Notably, model weights below 1,800 Hz seem to completely diminish in the -9 dB SNR  
 423 condition (**Figure 3B top**). Altogether, these results support the importance of low frequency  
 424 speech signals in speech encoding which is reasonable given that others have shown that this  
 425 range is important for speech intelligibility in noise (Chang, Bai, and Zeng 2006; Turner et al.  
 426 2004) and given that the type of noise used in the present study matches the speech spectrum.  
 427 **Table 2.** LME model results showing how prediction accuracies changed across SNRs and post-hoc  
 428 contrasts between each feature

Marginal R <sup>2</sup> = 0.204, Conditional R <sup>2</sup> = 0.457					
Fixed effects	Estimate	Std. Error	t value	p-value	
Onsets (intercept)	8.488 x 10 <sup>-2</sup>	4.513 x 10 <sup>-3</sup>	18.809	< 2 x 10 <sup>-16</sup>	***
SNR	-1.013 x 10 <sup>-2</sup>	5.171 x 10 <sup>-4</sup>	-19.581	< 2 x 10 <sup>-16</sup>	***
Sgram	2.875 x 10 <sup>-3</sup>	2.426 x 10 <sup>-3</sup>	1.185	0.236	
Fea	-1.112 x 10 <sup>-2</sup>	2.426 x 10 <sup>-3</sup>	-4.585	4.62 x 10 <sup>-6</sup>	***
Surp	-2.035 x 10 <sup>-2</sup>	2.426 x 10 <sup>-3</sup>	-8.391	< 2 x 10 <sup>-16</sup>	***
SNR:Sgram	7.829 x 10 <sup>-4</sup>	7.313 x 10 <sup>-4</sup>	1.070	0.284	
SNR:Fea	-2.113 x 10 <sup>-3</sup>	7.313 x 10 <sup>-4</sup>	-2.889	0.004	***
SNR:Surp	-3.668 x 10 <sup>-4</sup>	7.313 x 10 <sup>-4</sup>	-0.501	0.616	
Contrasts	Estimate	Std. Error	z ratio	p-value	
Ons – Sgram	-7.830 x 10 <sup>-4</sup>	7.310 x 10 <sup>-4</sup>	-1.070	0.708	
Ons – Fea	2.113 x 10 <sup>-3</sup>	7.310 x 10 <sup>-4</sup>	2.889	0.020	*
Ons – Surp	3.670 x 10 <sup>-4</sup>	7.310 x 10 <sup>-4</sup>	0.501	0.958	
Sgram – Fea	2.896 x 10 <sup>-3</sup>	7.310 x 10 <sup>-4</sup>	3.960	4 x 10 <sup>-4</sup>	***
Sgram – Surp	1.150 x 10 <sup>-3</sup>	7.310 x 10 <sup>-4</sup>	1.572	0.395	
Fea – Surp	-1.746 x 10 <sup>-3</sup>	7.310 x 10 <sup>-4</sup>	-2.388	0.079	

429  
 430 In the phonetic features case, although there is a decrease across practically all features  
 431 around 100 ms in the +3 dB SNR condition, voicing and some manner of articulation (plosive and  
 432 fricative) and vowel backness features (front, back, and diphthong) remained largely intact in this  
 433 condition (**Figure 3B bottom**). As noise levels continued to increase, we saw even more of a  
 434 decrease in each feature. In the -6 dB SNR condition, only some plosive, fricative, and vowel  
 435 backness feature weights remained whereas all other features here and in the -9 dB SNR  
 436 condition diminished. Although TRF weights for individual features decreased in the +3 dB SNR,  
 437 we began to see larger reductions in TRF weight amplitudes in the -6 dB SNR similar to

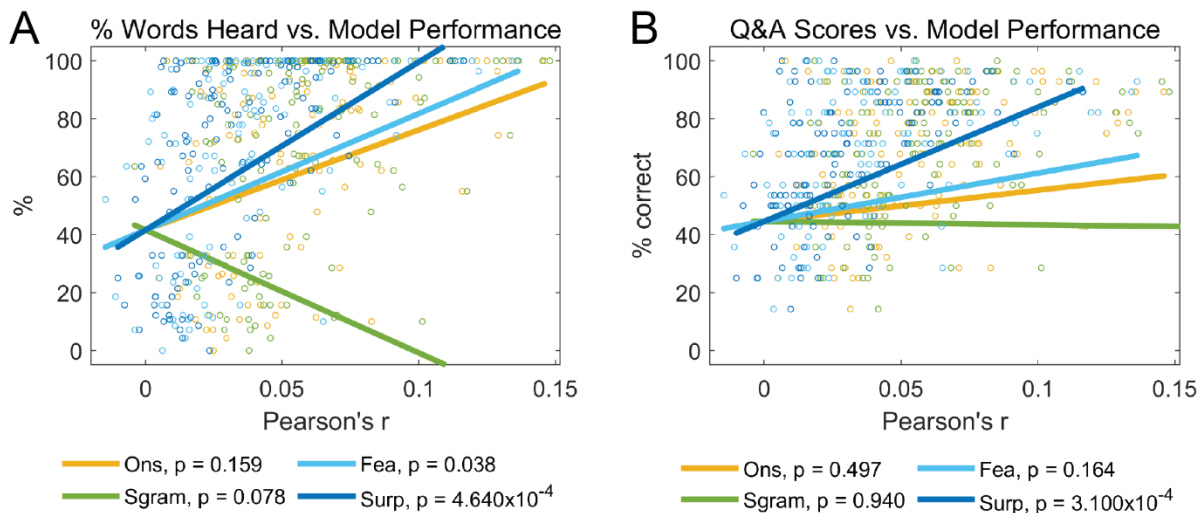
438 Swaminathan and Heinz who found their greatest lapse in phonetic feature reception around -5  
439 dB SNR (Swaminathan and Heinz 2012).



440  
441 **Figure 3.** Forward modeling results. **A.** Prediction accuracies which were calculated using acoustic onsets,  
442 spectrogram, phonetic feature, and surprisal from left to right. Significance is indicated by \* if  $p < 0.05$ , \*\* if  
443  $p < 0.01$ , and \*\*\* if  $p < 0.001$  using pairwise t-tests with FDR (BY) correction. Spectrogram (**B top**) and  
444 phonetic feature (**B bottom**) TRF model weights across the experimental conditions. The model weights  
445 were averaged across 12 frontotemporal electrodes and originated from a model that included acoustic  
446 onsets, spectrogram, phonetic features, and surprisal.  
447

## 448 **Linguistic features are highly predictive of behavior**

449 One of the main hypotheses we had in this study was that, across SNRs, EEG signatures  
450 of linguistic processing would be more closely correlated with behavior than EEG measures of  
451 low-level acoustic processing. To test this, we modeled the relationship between each speech  
452 feature's model prediction accuracy with the percentage of words heard and comprehension  
453 scores using LME models. We used this method so that we could pool data across SNRs which  
454 would result in 5 data points per participant. As such, we included fixed effects that corresponded  
455 to the model accuracies for each feature and a random effect for participant. We modeled both  
456 behavioral metrics separately.



457  
458 **Figure 4.** The relationship between forward model performance (prediction accuracy) and behavior. LME  
459 models were used to relate the model performance for each speech feature to the percentage of words  
460 heard (**A**) and comprehension scores (**B**). The marginal and conditional  $R^2$ s are 0.444 and 0.583 for A and  
461 0.429 and 0.701 for B.

462 In line with our hypothesis, these analyses show that lexical surprisal is highly predictive  
463 of both the percentage of words heard ( $\beta = 580.101$ ,  $p = 4.640 \times 10^{-4}$ , **Figure 4A**) and  
464 comprehension scores ( $\beta = 394.314$ ,  $p = 3.100 \times 10^{-4}$ , **Figure 4B**). Phonetic feature performance  
465 was significantly predictive of percentage of words heard ( $\beta = 401.892$ ,  $p = 0.038$ ), but not  
466 comprehension ( $\beta = 166.587$ ,  $p = 0.164$ ). The model did not produce any significant effects for

467 the acoustic onset ( $\beta = 345.738$ ,  $p = 0.159$  for % words heard;  $\beta = 107.556$ ,  $p = 0.497$  for  
468 comprehension) or spectrogram ( $\beta = -424.404$ ,  $p = 0.078$  for % words heard;  $\beta = -11.802$ ,  $p =$   
469  $0.940$  for comprehension) fixed effects. The spectrogram fixed effect may have negative trends  
470 because the spectrogram and acoustic onset model performances are highly (yet negatively)  
471 correlated; the percentage of words heard LME model reported a correlation of  $-0.623$  and the  
472 comprehension LME model reported a correlation of  $-0.602$ . The spectrogram fixed effect may be  
473 fitting to the noise that remains after the acoustic onsets have explained all it can behavior-wise.  
474 Nevertheless, these results support that linguistic features are more predictive of subjective and  
475 objective measures of speech perception than acoustic features.

#### 476 ***Envelope reconstruction accuracy maps well to behavior***

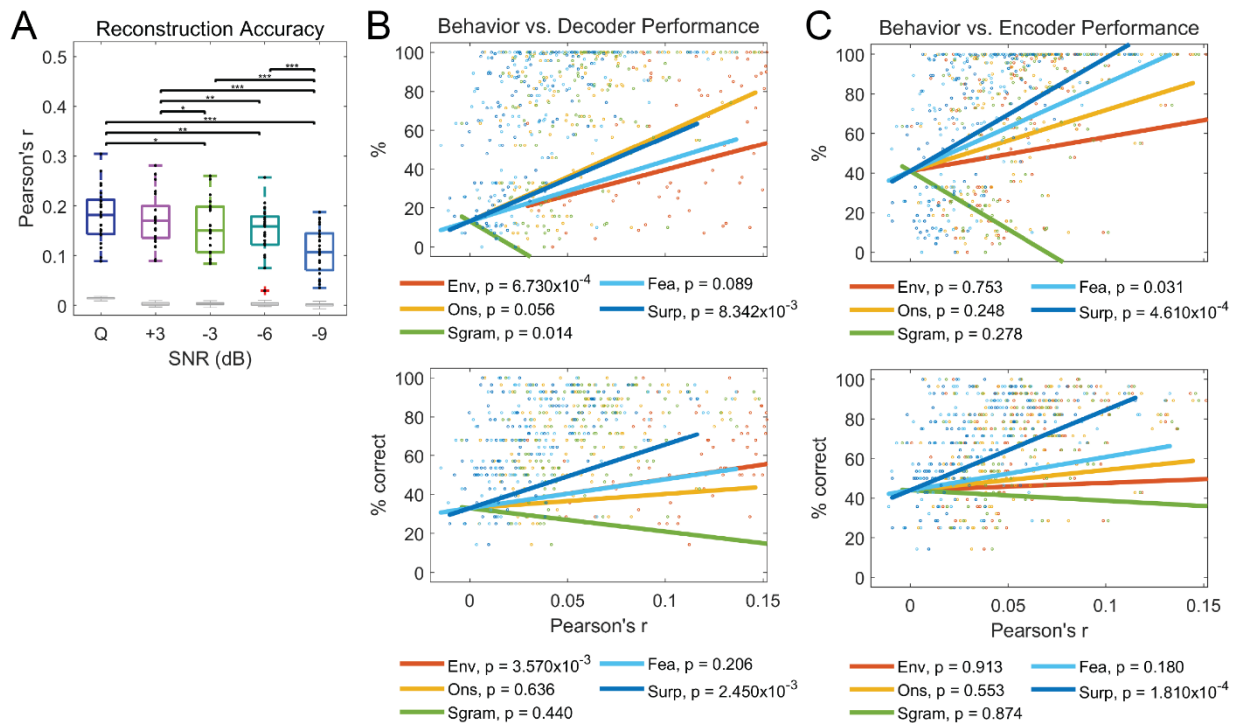
477 Numerous studies have shown that cortical activity tracks the temporal modulations of the  
478 speech envelope (Aiken and Picton 2008; Destoky et al. 2019; Di Liberto, O'Sullivan, and Lalor  
479 2015; Ding and Simon 2013; Etard and Reichenbach 2019; Lalor and Foxe 2010b; Nourski et al.  
480 2009; Pasley et al. 2012). Speech envelope reconstruction is a powerful method of indexing  
481 speech tracking given its overall advantage of using all scalp data to provide a better SNR. It is  
482 unknown, however, exactly what information envelope reconstruction indexes because the  
483 envelope has been shown to capture syllabic boundaries (Hertrich et al. 2012; Oganian and  
484 Chang 2019), phonetic feature information (Rosen 1992), and prosodic cues (Myers, Lense, and  
485 Gordon 2019). Nevertheless, due to its common usage in speech research and its improved SNR,  
486 we were interested in how speech envelope tracking would compare to our acoustic onset,  
487 spectrogram, phonetic feature, and surprisal encoding results and how it would relate to behavior.

488 The first step in this analysis was to use a backward modeling procedure to find how  
489 speech envelope tracking is affected by different levels of background noise. Speech in quiet and  
490 the +3 dB SNR condition shared similar reconstruction accuracies ( $p = 0.497$ , paired t-test with  
491 FDR [BY] correction) and were reconstructed more reliably than all other SNRs ( $p < 0.05$ ). The -  
492 3 dB SNR and -6 dB SNR accuracies were not significantly different from each other ( $p = 0.547$ )

493 but were both higher than the -9 dB SNR condition ( $p = 2.2 \times 10^{-7}$  and  $p = 2.5 \times 10^{-6}$ , **Figure 5A**).  
494 We compared the changes in reconstruction accuracy across SNRs to the other features by  
495 including it in the original LME model. Interestingly, the speech envelope feature had the steepest  
496 slope in model performance across SNRs compared to all other features (marginal/conditional  $R^2$   
497 = 0.521/0.644,  $\beta = -0.017$ ,  $p < 0.0001$ ). In other words, the speech envelope caused the greatest  
498 change in prediction accuracy across SNRs when all other features were fixed.

499         Next, we tested if the fidelity of envelope tracking was indicative of how well participants  
500 could hear and understand the story in comparison to the forward modeled features. This was  
501 tested using an LME model that contained fixed effects for envelope reconstruction accuracy and  
502 prediction accuracies for models trained on acoustic onsets, spectrograms, phonetic features,  
503 and lexical surprisal (and word onset). As expected, envelope reconstruction scores were highly  
504 predictive of percent words heard (marginal/conditional  $R^2 = 0.512/0.695$ ,  $\beta = 264.240$ ,  $p = 6.730$   
505  $\times 10^{-4}$ ) and comprehension (marginal/conditional  $R^2 = 0.476/0.705$ ,  $\beta = 149.070$ ,  $p = 3.570 \times 10^{-$   
506  $^3$ , **Figure 5B**). Surprisal was still highly predictive of both measures ( $\beta = 431.332$ ,  $p = 8.342 \times 10^{-$   
507  $^3$  for % words heard;  $\beta = 326.367$ ,  $p = 2.450 \times 10^{-3}$  for comprehension). Even though the phonetic  
508 feature predictors had greater or similar slopes to the envelope predictor, the LME models did not  
509 report these results as significant ( $\beta = 130.305$ ,  $p = 0.089$  for % words heard;  $\beta = 147.956$ ,  $p =$   
510  $0.206$  for comprehension). To our surprise, the spectrogram model performance now significantly  
511 (negatively) predicted percentage of words heard in this model as well ( $\beta = -587.122$ ,  $p = 0.014$ ).  
512 Furthermore, we calculated EEG prediction accuracy based on a full model trained on the  
513 concatenation of acoustic onsets, spectrogram, phonetic features, surprisal, and word onset. We  
514 then related the full model accuracy and envelope reconstruction accuracies to behavior. This full  
515 model was more predictive of behavior (% words heard, marginal/conditional  $R^2 = 0.497/0.731$ ,  $\beta$   
516 = 486.279,  $p = 5.900 \times 10^{-5}$ ; comprehension, marginal/conditional  $R^2 = 0.470/0.729$ ,  $\beta = 311.894$ ,  
517  $p = 3.180 \times 10^{-5}$ ) than envelope reconstruction accuracy (% words heard,  $\beta = 296.395$ ,  $p = 2.500$   
518  $\times 10^{-4}$ ; comprehension,  $\beta = 167.173$ ,  $p = 1.190 \times 10^{-3}$ ).

519 The main caveat of the results above is that we're comparing two different modeling  
 520 methods. Backward models have the ability to take advantage of all EEG channels and up-weight  
 521 more informative channels, resulting in higher model performance compared to forward models.  
 522 As such, we reperformed the analysis with a forward model that included acoustic onsets,  
 523 spectrograms, phonetic features, word surprisal, and the speech envelope (and word onsets). In  
 524 this case, we found that the change in the envelope encoding across SNRs is similar to all other  
 525 features modeled except the phonetic features (marginal/conditional  $R^2 = 0.199/0.491$ ,  $p =$   
 526  $0.0485$ , LME model).



527

528 **Figure 5.** Envelope modeling results. **A.** The colored boxplots are the envelope reconstruction accuracies  
 529 for each condition, and the black points represent each participant. The gray boxplots represent mean  
 530 reconstruction accuracies (per person) based on shuffled permutations. Significance is indicated by \* if  $p <$   
 531  $0.05$ , \*\* if  $p < 0.01$ , and \*\*\* if  $p < 0.001$  using pairwise t-tests with FDR (BY) correction. **B.** We then used  
 532 LME models to determine the relationship between envelope decoder performance and behavior  
 533 (percentage of words heard on top and comprehension scores on the bottom). Each circle represents a

534 participant, so there are 125 circles (5 conditions x 25 participants). **C.** Same as B, except using results  
535 from a forward model trained on all features including the envelope.

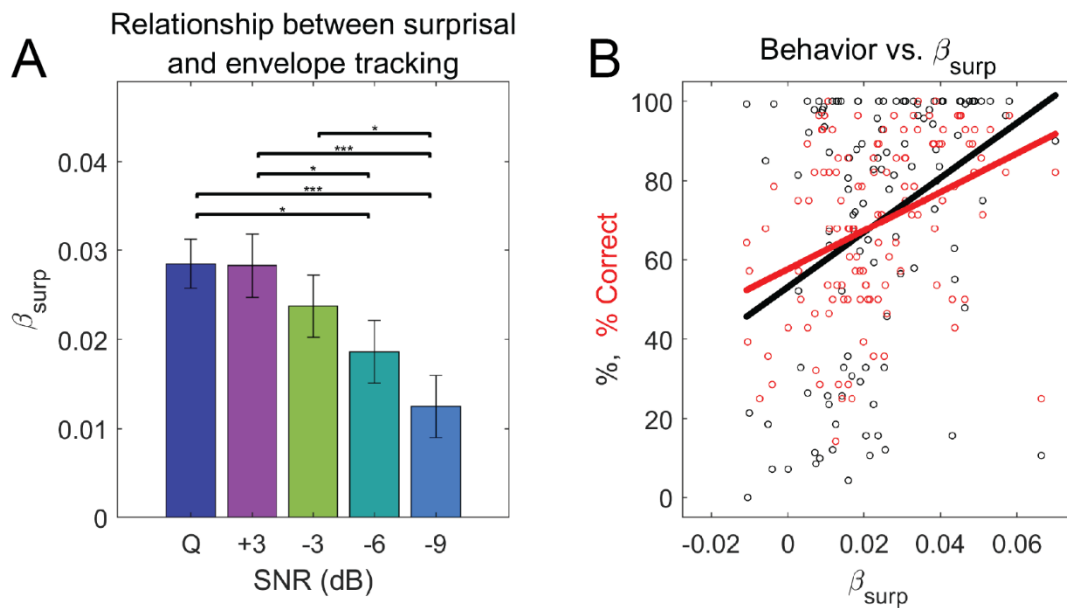
536 This reanalysis also showed that higher-level linguistic features are most correlated with  
537 behavior (surprisal vs. % words heard,  $\beta = 570.767$ ,  $p = 4.610 \times 10^{-4}$ ; surprisal vs. comprehension,  
538  $\beta = 405.383$ ,  $p = 1.810 \times 10^{-4}$ ; phonetic features vs. % words heard,  $\beta = 442.387$ ,  $p = 0.031$ ). The  
539 envelope (% words heard,  $\beta = 171.309$ ,  $p = 0.753$ ; comprehension,  $\beta = 37.104$ ,  $p = 0.913$ ),  
540 acoustic onset (% words heard,  $\beta = 307.682$ ,  $p = 0.248$ ; comprehension,  $\beta = 102.556$ ,  $p = 0.553$ ),  
541 and spectrogram (% words heard,  $\beta = -587.997$ ,  $p = 0.278$ ; comprehension,  $\beta = -53.038$ ,  $p =$   
542  $0.874$ ) model effects were not significant in addition to phonetic features in the comprehension  
543 model ( $\beta = 167.955$ ,  $p = 0.180$ , **Figure 5C**). The marginal and conditional  $R^2$  are 0.445 and 0.568  
544 for the % words heard model and 0.434 and 0.697 for the comprehension model. In short, when  
545 modeling the envelope in the forward direction, its encoding decreases at a similar rate to other  
546 low-level acoustic features and it predicts behavior similar to them as well. Conversely, when  
547 using stimulus reconstruction, envelope models perform best across SNRs and predict  
548 comprehension scores similar to higher level features.

#### 549 ***Lexical surprisal's influence on early auditory encoding decreases at high noise levels***

550 Recent work by Broderick and colleagues has shown that the cortical tracking of an  
551 individual word's acoustic and phonetic representations was enhanced the more semantically  
552 similar it was to its preceding context (Broderick, Anderson, and Lalor 2019). Since higher-level  
553 representations can bias perception when incoming stimuli are noisy (de Lange, Heilbron, and  
554 Kok 2018), perhaps in the present study, participants relied more on a higher-level feature (i.e.,  
555 next word probability/surprisal) when noise levels slightly increased—thereby strengthening the  
556 relationship between word surprisal and lower-level feature encoding. This influence of lexical  
557 surprisal would then decrease the noisier the speech became, i.e., as people begin to fail to  
558 understand the speech. We aimed to test these hypotheses using a two-stage regression analysis  
559 that consists of stimulus reconstruction followed by LME modeling.



560 Stimulus reconstruction was used (due to its increased SNR) to reconstruct an estimate  
561 of the speech envelope for each individual word. We then compared the predicted and actual  
562 envelopes using Spearman's correlation. All words in the story were also scored on their lexical  
563 surprisal, which was calculated as the negative logarithm of a word's probability given the words  
564 that came before it. Envelope variability, relative pitch, and resolvability (the nuisance regressors)  
565 were also calculated to control for acoustic related properties of the stimuli, many of which could  
566 correlate with word surprisal and word reconstruction accuracy. Surprisal, SNR, and the nuisance  
567 regressors were included in an LME model to predict word reconstruction accuracy for the first  
568 100 ms of each word.



569  
570 **Figure 6.** The relationship between lexical surprisal and envelope tracking in different levels of noise. **A.**  
571 Surprisal coefficient in quiet and the interaction between surprisal and SNR using a 100 ms word window.  
572 The error bars are the standard errors of each measure, calculated using the LME model. Significance is  
573 indicated by \* if  $p < 0.05$ , \*\* if  $p < 0.01$ , and \*\*\* if  $p < 0.001$  using R's "emtrends" function to compare  
574 estimated marginal means of linear trends. **B.** The relationship between the surprisal coefficients across  
575 SNRs and behavior. Percentage of words heard are in black, and the comprehension scores are in red.  
576 There are 125 circles (5 SNRs x 25 participants).

577 There was a significantly positive relationship between lexical surprisal and word  
 578 reconstruction accuracy in quiet (marginal/conditional  $R^2 = 0.017/0.030$ ,  $\beta = 2.850 \times 10^{-2}$ ,  $t =$   
 579  $10.519$ ,  $p < 2.00 \times 10^{-16}$ , **Table 3**). In other words, the more surprising a word (or the less probable  
 580 that word was given its preceding context), the greater that word's envelope was reflected in the  
 581 EEG. The interaction between surprisal and SNR shows how this coefficient changes in the other  
 582 noise levels (**Figure 6A**). The influence of lexical surprisal on word reconstruction accuracy in  
 583 quiet was similar to the +3 dB ( $\beta = -2.172 \times 10^{-4}$ ,  $t = -0.062$ ,  $p = 0.951$ ) and -3 dB ( $\beta = -4.767 \times 10^{-$   
 584  $3$ ,  $t = -1.363$ ,  $p = 0.173$ ) SNR conditions and greater than the -6 dB ( $\beta = -9.900 \times 10^{-3}$ ,  $t = -2.817$ ,  
 585  $p = 0.005$ ) and -9 dB ( $\beta = -1.600 \times 10^{-2}$ ,  $t = -4.570$ ,  $p = 4.890 \times 10^{-6}$ ) SNR conditions. This shows  
 586 that lexical surprisal impacts envelope tracking in low levels of background noise (+3 dB and -3  
 587 dB SNRs) similar to when speech is in quiet, but this influence significantly decreases with  
 588 moderate to high levels of noise. It is also important to note that the nuisance regressors and  
 589 additional interactions we've included also significantly influenced word reconstruction accuracy,  
 590 reinforcing the value of controlling for those variables and ensuring our results were not  
 591 confounded by other acoustic related measures.

592 **Table 3.** LME model showing the relationship between lexical surprisal and envelope tracking in each  
 593 SNR.

	Estimate	Std. Error	t value	Pr(> t )	
Quiet (Intercept)	$8.992 \times 10^{-2}$	$5.451 \times 10^{-3}$	16.496	$< 2 \times 10^{-16}$	***
SNR +3	$-4.703 \times 10^{-3}$	$3.514 \times 10^{-3}$	-1.338	0.18077	
SNR -3	$-1.671 \times 10^{-2}$	$3.524 \times 10^{-3}$	-4.724	$2.12 \times 10^{-6}$	***
SNR -6	$-3.196 \times 10^{-2}$	$3.496 \times 10^{-3}$	-9.141	$< 2 \times 10^{-16}$	***
SNR -9	$-5.822 \times 10^{-2}$	$3.511 \times 10^{-3}$	-16.586	$< 2 \times 10^{-16}$	***
surp	$2.850 \times 10^{-2}$	$2.709 \times 10^{-3}$	10.519	$< 2 \times 10^{-16}$	***
surp:SNR +3	$-2.172 \times 10^{-4}$	$3.528 \times 10^{-3}$	-0.062	0.95091	
surp:SNR -3	$-4.767 \times 10^{-3}$	$3.497 \times 10^{-3}$	-1.363	0.17287	
surp:SNR -6	$-9.900 \times 10^{-3}$	$3.514 \times 10^{-3}$	-2.817	0.00484	**
surp:SNR -9	$-1.600 \times 10^{-2}$	$3.501 \times 10^{-3}$	-4.570	$4.89 \times 10^{-6}$	***
envStd	$6.042 \times 10^{-2}$	$1.430 \times 10^{-3}$	42.251	$< 2 \times 10^{-16}$	***
$f_{rel}$	$8.697 \times 10^{-3}$	$1.576 \times 10^{-3}$	5.519	$3.42 \times 10^{-8}$	***
res	$-1.605 \times 10^{-2}$	$1.453 \times 10^{-3}$	-11.051	$< 2 \times 10^{-16}$	***
envStd: $f_{rel}$	$6.692 \times 10^{-3}$	$1.148 \times 10^{-3}$	5.832	$5.49 \times 10^{-9}$	***
$f_{rel}$ :res	$-6.327 \times 10^{-3}$	$1.450 \times 10^{-3}$	-4.363	$1.28 \times 10^{-5}$	***
surp:envStd	$1.624 \times 10^{-2}$	$1.305 \times 10^{-3}$	12.448	$< 2 \times 10^{-16}$	***
surp:res	$-6.872 \times 10^{-3}$	$1.325 \times 10^{-3}$	-5.186	$2.17 \times 10^{-7}$	***

595           Permutation tests were conducted to test if each surprisal estimate was significantly  
596 different from chance. The surprisal values were shuffled 5000 times while keeping all other  
597 variables intact, and a new LME model was computed with each shuffle. This procedure resulted  
598 in none of the null coefficients being greater than the observed values. Despite studies showing  
599 that low SNRs benefit from the utilization of context (Mayo, Florentine, and Buus 1997), we found  
600 the permutation results for the -6 dB and -9 dB conditions to be surprising. Participants reported  
601 hearing either no words or very few words in, for instance, the -9 dB SNR condition, yet we still  
602 see a significant influence of lexical surprisal on envelope tracking. We did not control for any  
603 other features (including surprisal) in our envelope reconstruction model, so these unaccounted-  
604 for features may have contributed to the significantly positive interaction coefficients in the -6 dB  
605 and -9 dB SNR conditions. Or the fact that participants were able to hear any words in both  
606 conditions may have been enough to cause this significant effect. Lastly, we found that single  
607 subject surprisal coefficients also predicted their self-reported percentage of words heard  
608 (marginal/conditional  $R^2 = 0.124/NA$ ,  $\beta = 691.434$ ,  $p = 5.290 \times 10^{-5}$ ) and comprehension scores  
609 (marginal/conditional  $R^2 = 0.138, 0.177$ ,  $\beta = 488.401$ ,  $p = 2.140 \times 10^{-5}$ ) across SNRs (**Figure 6B**).  
610 That is to say, the stronger the influence of surprisal on envelope tracking, the better the  
611 participants were able to hear and comprehend the story.

## 612 **DISCUSSION**

613           This study sought to establish how well indices of hierarchical neural speech processing  
614 reflect language comprehension—advancing on prior work that has typically tested specific  
615 hierarchical levels without controlling for the others. We first characterized how the encoding of a  
616 range of hierarchical speech features diminished in noise and if those changes in encoding were  
617 predictive of behavior. We found that the encoding of acoustic and surprisal features declined  
618 similarly as noise levels increased, and that phonetic feature encoding was more affected by noise  
619 than the acoustic features. In addition, lexical surprisal and phonetic feature encoding were the  
620 most predictive of participants' behavioral scores across SNRs. Speech envelope models were

621 predictive of behavior, but only when employing decoding models. Lastly, we investigated how  
622 lexical surprisal influenced the neural tracking of the speech envelope. In general, we found that  
623 the envelopes of more unexpected words were better reflected in the EEG. This was true in quiet  
624 and in low levels of background noise, but this relationship weakened as noise levels increased.

625 We hypothesized that acoustic features would be the most invariant to noise, but this was  
626 only partially supported by our results. The degree to which envelope and spectrogram features  
627 were reflected in EEG decreased at a slower rate only in comparison to phonetic features.  
628 However, when we analyzed envelope reconstruction accuracies, rather than EEG predictions  
629 based on the speech envelope, decoding accuracies decreased at a faster rate than all other  
630 features. This was surprising, first, due to a previous finding that the synchronization between  
631 neural activity and the speech envelope remained unaffected until the speech signal had an SNR  
632 of -9 dB (Ding and Simon 2013). These stark differences may have been due to a combination of  
633 factors: neural recording modality, data preprocessing, model training and testing procedures  
634 between conditions, or the regularization method used (e.g., boosting versus ridge regression).  
635 Instead, our results show a gradual decrease in envelope tracking across SNRs similar to  
636 Vanthornhout and colleagues (Vanthornhout et al. 2018).

637 Although the rate at which our acoustic and linguistic model accuracies declined did not  
638 completely support our hypothesis, these results may not be surprising given recent work. Kell  
639 and McDermott measured primary and non-primary auditory cortices' invariance to background  
640 noise using fMRI. Invariance was measured by correlating voxel responses to natural sounds in  
641 quiet with the voxel response to those same sounds in noise. They found that primary and non-  
642 primary auditory cortices were similarly invariant to natural sounds in spectrally matched  
643 background noise tested at a 0 dB SNR. However, non-primary areas became more robust to  
644 noise than primary areas when sounds were presented in real-world noise (Kell and McDermott  
645 2019). So, our model performances may result from how the brain represents speech in the type  
646 of synthetic noise we used. Models in the present study could have also been affected by

647 attention. Participants may have allocated less attention to the -9 dB SNR trials due to the large  
648 amount of noise, in turn skewing the encoding/decoding of the different speech features.

649 Our findings also showed that phonetic features and lexical surprisal were most predictive  
650 of subjective behavioral metrics (percentage of words heard) and lexical surprisal was most  
651 predictive of objective metrics (comprehension). Previously published work has shown that neural  
652 measures of lexical surprisal is highly predictive of behavior (Mesik, Ray, and Wojtczak 2021).  
653 Many studies have also shown that the speech envelope (using stimulus reconstruction or cross-  
654 correlation) contributes and relates to speech intelligibility and comprehension (Ahissar et al.  
655 2001; Decruy et al. 2020; Itzov and Parra 2019; Lesenfants et al. 2019; Muncke, Kuruvila, and  
656 Hoppe 2022; Vanthornhout et al. 2018). Once we included envelope reconstructions in our  
657 analysis, it also proved to be an accurate predictor of behavior. However, backward modeling  
658 greatly improves overall model performance due to its ability to utilize all recorded neural  
659 channels, thereby increasing neural signal-to-noise ratio. Spectrogram and phonetic features  
660 have previously been shown to better predict EEG than the speech envelope (Di Liberto,  
661 O'Sullivan, and Lalor 2015), so we believe that our behavior-prediction accuracy correlations were  
662 due to how the envelope was modeled, rather than the information the speech envelope itself  
663 carries or how well it is reflected in the brain.

664 Interestingly, we found that phonetic features uniquely predicted neural activity even when  
665 controlling for the speech spectrogram and acoustic onsets. This is in line with previous studies  
666 showing that the addition of phonetic features to spectrotemporal representations improve EEG  
667 prediction (Di Liberto, O'Sullivan, and Lalor 2015; Sohoglu and Davis 2020) and its correlation  
668 with speech intelligibility (Lesenfants et al. 2019) and that phonetic features uniquely predict EEG  
669 responses even when attending to a specific talker (Teoh, Ahmed, and Lalor 2022). However, the  
670 present phonetic feature results contrast with previous work which suggested that responses to  
671 articulations could be explained by simpler acoustic features (Daube, Ince, and Gross 2019).  
672 Nevertheless, our findings that phonetic feature encoding declines at a different rate and better

673 predicts behavior compared to the spectrogram, provides further evidence that the two features  
674 are dissociable.

675 Another one of our key hypotheses was that participants would use lexical context to  
676 predict and encode the acoustic features of each word. This was found to be true: our LME  
677 analysis (stage two of the two-stage regression) revealed that the more unexpected a word, the  
678 better we were able to reconstruct that word's envelope. However, we had also hypothesized that  
679 participants would rely more on these predictions for speech in moderate levels of noise (when  
680 speech is still intelligible) relative to speech in quiet, before falling off at high levels of background  
681 noise (when speech is no longer intelligible). This result was only partially borne out. Specifically,  
682 the use of lexical context in processing the speech acoustics did decrease as the speech became  
683 noisier, but there was no evidence to support a stronger reliance on context in moderate levels of  
684 noise. In particular, while there was no difference in comprehension scores between the quiet and  
685 +3 dB SNR conditions, there was no increase in the influence of surprisal on envelope tracking  
686 for the latter condition compared to speech in quiet. We did notice a larger spread of the  
687 percentage of words heard scores across subjects in the +3 dB SNR condition. So, we explored  
688 the possibility that the subjects who were starting to struggle might put forth more effort to  
689 understand and process the story by relying more on context (and thus might have a higher  
690 surprisal weight in **Figure 6**) than those who remained at ceiling. But we found no significant  
691 difference. This was a little surprising given that context has been known to affect behavior  
692 (Golestani et al. 2013) and neural activity (Boulenger et al. 2011; Koskinen et al. 2020; Strauß et  
693 al. 2022) in challenging listening conditions. Future work with larger subject numbers and perhaps  
694 even lower levels of background noise (e.g., + 6 dB SNR) might reveal such an effect.

695 Our LME model analysis based on word surprisal seems on face value to be at odds with  
696 Broderick and colleagues who found that the envelopes of words that were more semantically  
697 *similar* to their context were better reflected in the EEG. That is to say, envelope tracking is  
698 enhanced for words that share a similar meaning with their context (Broderick, Anderson, and

699 Lalor 2019). However, semantic similarity and lexical surprisal tend to share a moderate, and  
700 sometimes weak negative correlation (Frank and Willems 2017). Indeed, a re-analysis of  
701 Broderick et al.'s original EEG data has revealed that both semantic similarity and lexical surprisal  
702 play complementary (positive) roles in estimating when envelope tracking is enhanced (Broderick  
703 and Lalor 2020). Nevertheless, the nature of this duality remains mysterious, and we hope it will  
704 provide the grounds for an exciting body of future work. We anticipate it will take a substantial  
705 battery of future experiments to shape a unifying explanation, with stimuli that can disentangle  
706 correlations between semantic similarity, lexical surprisal, and other linguistic factors that could  
707 come into play (e.g., semantic content, next-word entropy, phonetic surprisal, next-phoneme  
708 entropy). In any case, what seems clear in the present results is that lexical context influences  
709 the neural tracking of speech acoustics on a word-by-word basis, and this influence drops as  
710 speech becomes unintelligible.

711 In summary, the current results show that phonetic features are more susceptible to noise  
712 than acoustic speech features. While linguistic features are more predictive of behavior than  
713 acoustic features, envelope decoding models can be used to improve this relationship. We have  
714 also found that the encoding of certain phonetic features decreases in even low levels of noise,  
715 and that the encoding of frequencies below 1.3k essentially disappears in high noise levels.  
716 Lastly, we show support that context influences a word's acoustic encoding. This influence  
717 lessens in high background noise levels. Future work will aim to further characterize how people  
718 might rely more or less on top-down context to process bottom-up speech input as a function of  
719 stimulus type, task, and listening conditions.

## 720 REFERENCES

721 Ahissar, Ehud, Srikantan Nagarajan, Merav Ahissar, Athanassios Protopapas, Henry Mahncke,  
722 and Michael M Merzenich. 2001. 'Speech comprehension is correlated with temporal  
723 response patterns recorded from auditory cortex', *Proceedings of the National Academy  
724 of Sciences*, 98: 13367-72.

- 725 Aiken, S. J., and T. W. Picton. 2008. 'Human cortical responses to the speech envelope', *Ear*  
726 *Hear*, 29: 139-57.
- 727 Anderson, A. J., J. R. Binder, L. Fernandino, C. J. Humphries, L. L. Conant, M. Aguilar, X. Wang,  
728 D. Doko, and R. D. S. Raizada. 2017. 'Predicting Neural Activity Patterns Associated with  
729 Sentences Using a Neurobiologically Motivated Model of Semantic Representation',  
730 *Cereb Cortex*, 27: 4379-95.
- 731 Bigdely-Shamlo, N., T. Mullen, C. Kothe, K. M. Su, and K. A. Robbins. 2015. 'The PREP pipeline:  
732 standardized preprocessing for large-scale EEG analysis', *Front Neuroinform*, 9: 16.
- 733 Boersma, Paul, and David Weenink. 2013. 'Praat: doing phonetics by computer [Computer  
734 program]. Version 5.3. 51', Online: <http://www.praat.org/retrieved>, last viewed on, 12.
- 735 Boulenger, Véronique, Michel Hoen, Caroline Jacquier, and Fanny Meunier. 2011. 'Interplay  
736 between acoustic/phonetic and semantic processes during spoken sentence  
737 comprehension: An ERP study', *Brain and language*, 116: 51-63.
- 738 Brodbeck, C., L. E. Hong, and J. Z. Simon. 2018. 'Rapid Transformation from Auditory to Linguistic  
739 Representations of Continuous Speech', *Curr Biol*, 28: 3976-83 e5.
- 740 Brodbeck, C., A. Jiao, L. E. Hong, and J. Z. Simon. 2020. 'Neural speech restoration at the cocktail  
741 party: Auditory cortex recovers masked speech of both attended and ignored speakers',  
742 *PLoS Biol*, 18: e3000883.
- 743 Brodbeck, C., A. Presacco, and J. Z. Simon. 2018. 'Neural source dynamics of brain responses  
744 to continuous stimuli: Speech processing from acoustics to comprehension', *Neuroimage*,  
745 172: 162-74.
- 746 Broderick, M. P., A. J. Anderson, G. M. Di Liberto, M. J. Crosse, and E. C. Lalor. 2018.  
747 'Electrophysiological Correlates of Semantic Dissimilarity Reflect the Comprehension of  
748 Natural, Narrative Speech', *Curr Biol*, 28: 803-09 e3.
- 749 Broderick, M. P., A. J. Anderson, and E. C. Lalor. 2019. 'Semantic Context Enhances the Early  
750 Auditory Encoding of Natural Speech', *J Neurosci*, 39: 7564-75.
- 751 Broderick, Michael P, and Edmund C Lalor. 2020. 'Co-existence of prediction and error signals in  
752 electrophysiological responses to natural speech', *bioRxiv*: 2020.11. 20.391227.



- 753 Broderick, Michael P, Nathaniel J Zuk, Andrew J Anderson, and Edmund C Lalor. 2022. 'More  
754 than words: Neurophysiological correlates of semantic dissimilarity depend on  
755 comprehension of the speech narrative', *European Journal of Neuroscience*, 56: 5201-14.
- 756 Burton, M. W. 2009. 'Understanding the role of the prefrontal cortex in phonological processing',  
757 *Clin Linguist Phon*, 23: 180-95.
- 758 Chang, J. E., J. Y. Bai, and F. G. Zeng. 2006. 'Unintelligible low-frequency sound enhances  
759 simulated cochlear-implant speech recognition in noise', *IEEE Trans Biomed Eng*, 53:  
760 2598-601.
- 761 Crosse, M. J., G. M. Di Liberto, A. Bednar, and E. C. Lalor. 2016. 'The Multivariate Temporal  
762 Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to  
763 Continuous Stimuli', *Front Hum Neurosci*, 10: 604.
- 764 Crosse, M. J., G. M. Di Liberto, and E. C. Lalor. 2016. 'Eye Can Hear Clearly Now: Inverse  
765 Effectiveness in Natural Audiovisual Speech Processing Relies on Long-Term  
766 Crossmodal Temporal Integration', *J Neurosci*, 36: 9888-95.
- 767 Dai, Zihang, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov.  
768 2019. 'Transformer-xl: Attentive language models beyond a fixed-length context', *arXiv*  
769 *preprint arXiv:1901.02860*.
- 770 Daube, Christoph, Robin AA Ince, and Joachim Gross. 2019. 'Simple acoustic features can  
771 explain phoneme-based predictions of cortical responses to speech', *Current Biology*, 29:  
772 1924-37. e9.
- 773 de Heer, W. A., A. G. Huth, T. L. Griffiths, J. L. Gallant, and F. E. Theunissen. 2017. 'The  
774 Hierarchical Cortical Organization of Human Speech Processing', *J Neurosci*, 37: 6539-  
775 57.
- 776 de Lange, F. P., M. Heilbron, and P. Kok. 2018. 'How Do Expectations Shape Perception?',  
777 *Trends Cogn Sci*, 22: 764-79.
- 778 Decruy, L., D. Lesenfants, J. Vanthornhout, and T. Francart. 2020. 'Top-down modulation of  
779 neural envelope tracking: The interplay with behavioral, self-report and neural measures  
780 of listening effort', *Eur J Neurosci*, 52: 3375-93.
- 781 Delorme, A., and S. Makeig. 2004. 'EEGLAB: an open source toolbox for analysis of single-trial  
782 EEG dynamics including independent component analysis', *J Neurosci Methods*, 134: 9-  
783 21.

- 784 Destoky, F., M. Philippe, J. Bertels, M. Verhasselt, N. Coquelet, M. Vander Ghinst, V. Wens, X.  
785 De Tieghe, and M. Bourguignon. 2019. 'Comparing the potential of MEG and EEG to  
786 uncover brain tracking of speech temporal envelope', *Neuroimage*, 184: 201-13.
- 787 Di Liberto, G. M., J. A. O'Sullivan, and E. C. Lalor. 2015. 'Low-Frequency Cortical Entrainment to  
788 Speech Reflects Phoneme-Level Processing', *Curr Biol*, 25: 2457-65.
- 789 Di Liberto, Giovanni M, James A O'Sullivan, and Edmund C Lalor. 2015. 'Low-Frequency Cortical  
790 Entrainment to Speech Reflects Phoneme-Level Processing', *Current Biology*, 25: 2457-  
791 65.
- 792 Di Liberto, Giovanni M, Daniel Wong, Gerda Ana Melnik, and Alain de Cheveigné. 2019. 'Low-  
793 frequency cortical responses to natural speech reflect probabilistic phonotactics',  
794 *Neuroimage*, 196: 237-47.
- 795 Ding, N., and J. Z. Simon. 2013. 'Adaptive temporal encoding leads to a background-insensitive  
796 cortical representation of speech', *J Neurosci*, 33: 5728-35.
- 797 Drullman, R., J. M. Festen, and R. Plomp. 1994. 'Effect of reducing slow temporal modulations  
798 on speech reception', *J Acoust Soc Am*, 95: 2670-80.
- 799 Etard, O., and T. Reichenbach. 2019. 'Neural Speech Tracking in the Theta and in the Delta  
800 Frequency Band Differentially Encode Clarity and Comprehension of Speech in Noise', *J  
801 Neurosci*, 39: 5750-59.
- 802 Frank, Stefan L, and Roel M Willems. 2017. 'Word predictability and semantic similarity show  
803 distinct patterns of brain activity during language comprehension', *Language, Cognition  
804 and Neuroscience*, 32: 1192-203.
- 805 Gillis, M., J. Vanthornhout, J. Z. Simon, T. Francart, and C. Brodbeck. 2021. 'Neural Markers of  
806 Speech Comprehension: Measuring EEG Tracking of Linguistic Speech Representations,  
807 Controlling the Speech Acoustics', *J Neurosci*, 41: 10316-29.
- 808 Golestani, N., A. Hervais-Adelman, J. Obleser, and S. K. Scott. 2013. 'Semantic versus perceptual  
809 interactions in neural processing of speech-in-noise', *Neuroimage*, 79: 52-61.
- 810 Gwilliams, Laura, Jean-Remi King, Alec Marantz, and David Poeppel. 2020. 'Neural dynamics of  
811 phoneme sequencing in real speech jointly encode order and invariant content', *bioRxiv*.
- 812 Hamilton, L. S., E. Edwards, and E. F. Chang. 2018. 'A Spatial Map of Onset and Sustained  
813 Responses to Speech in the Human Superior Temporal Gyrus', *Curr Biol*, 28: 1860-71 e4.

- 814 Heilbron, Micha, Kristijan Armeni, Jan-Mathijs Schoffelen, Peter Hagoort, and Floris P De Lange.  
815 2022. 'A hierarchy of linguistic predictions during natural language comprehension',  
816 *Proceedings of the National Academy of Sciences*, 119: e2201968119.
- 817 Hertrich, I., S. Dietrich, J. Trouvain, A. Moos, and H. Ackermann. 2012. 'Magnetic brain activity  
818 phase-locked to the envelope, the syllable onsets, and the fundamental frequency of a  
819 perceived speech signal', *Psychophysiology*, 49: 322-34.
- 820 Hickok, Greg. 2015. *Neurobiology of language* (Elsevier: Boston, MA).
- 821 Hickok, Gregory, and David Poeppel. 2007. 'The cortical organization of speech processing',  
822 *Nature Reviews Neuroscience*, 8: 393-402.
- 823 Huth, Alexander G, Wendy A De Heer, Thomas L Griffiths, Frédéric E Theunissen, and Jack L  
824 Gallant. 2016. 'Natural speech reveals the semantic maps that tile human cerebral cortex',  
825 *Nature*, 532: 453-58.
- 826 Iotzov, I., and L. C. Parra. 2019. 'EEG can predict speech intelligibility', *J Neural Eng*, 16: 036008.
- 827 Irino, T., and R. D. Patterson. 2006. 'A Dynamic Compressive Gammachirp Auditory Filterbank',  
828 *IEEE Trans Audio Speech Lang Process*, 14: 2222-32.
- 829 Kell, A. J. E., and J. H. McDermott. 2019. 'Invariance to background noise as a signature of non-  
830 primary auditory cortex', *Nat Commun*, 10: 3958.
- 831 Kleiner, Mario, David Brainard, and Denis Pelli. 2007. 'What's new in Psychtoolbox-3?'
- 832 Koskinen, M., M. Kurimo, J. Gross, A. Hyvarinen, and R. Hari. 2020. 'Brain activity reflects the  
833 predictability of word sequences in listened continuous speech', *Neuroimage*, 219:  
834 116936.
- 835 Lalor, E. C., and J. J. Foxe. 2010a. 'Neural responses to uninterrupted natural speech can be  
836 extracted with precise temporal resolution', *European Journal of Neuroscience*, 31: 189–  
837 93.
- 838 Lalor, Edmund C, and John J Foxe. 2010b. 'Neural responses to uninterrupted natural speech  
839 can be extracted with precise temporal resolution', *European journal of neuroscience*, 31:  
840 189-93.

- 841 Lesenfants, D., J. Vanthornhout, E. Verschueren, L. Decruy, and T. Francart. 2019. 'Predicting  
842 individual speech intelligibility from the cortical tracking of acoustic- and phonetic-level  
843 speech representations', *Hear Res*, 380: 1-9.
- 844 Maddox, R. K., and A. K. C. Lee. 2018. 'Auditory Brainstem Responses to Continuous Natural  
845 Speech in Human Listeners', *eNeuro*, 5.
- 846 MATLAB. 2019. Natick, Massachusetts: The MathWorks Inc.
- 847 ———. 2021. Natick, Massachusetts: The MathWorks Inc.
- 848 Mayo, L. H., M. Florentine, and S. Buus. 1997. 'Age of second-language acquisition and  
849 perception of speech in noise', *J Speech Lang Hear Res*, 40: 686-93.
- 850 McAuliffe, Michael, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger.  
851 2017. "Montreal Forced Aligner: Trainable Text-Speech Alignment Using Kaldi." In  
852 *Interspeech*, 498-502.
- 853 McDermott, J. H., and E. P. Simoncelli. 2011. 'Sound texture perception via statistics of the  
854 auditory periphery: evidence from sound synthesis', *Neuron*, 71: 926-40.
- 855 Mesgarani, N., C. Cheung, K. Johnson, and E. F. Chang. 2014. 'Phonetic feature encoding in  
856 human superior temporal gyrus', *Science*, 343: 1006-10.
- 857 Mesik, J., L. Ray, and M. Wojtczak. 2021. 'Effects of Age on Cortical Tracking of Word-Level  
858 Features of Continuous Competing Speech', *Front Neurosci*, 15: 635126.
- 859 Muncke, J., I. Kuruvila, and U. Hoppe. 2022. 'Prediction of Speech Intelligibility by Means of EEG  
860 Responses to Sentences in Noise', *Front Neurosci*, 16: 876421.
- 861 Myers, B. R., M. D. Lense, and R. L. Gordon. 2019. 'Pushing the Envelope: Developments in  
862 Neural Entrainment to Speech and the Biological Underpinnings of Prosody Perception',  
863 *Brain Sci*, 9.
- 864 Nourski, K. V., R. A. Reale, H. Oya, H. Kawasaki, C. K. Kovach, H. Chen, M. A. Howard, 3rd, and  
865 J. F. Brugge. 2009. 'Temporal envelope of time-compressed speech represented in the  
866 human auditory cortex', *J Neurosci*, 29: 15564-74.
- 867 Oganian, Y., and E. F. Chang. 2019. 'A speech envelope landmark for syllable encoding in human  
868 superior temporal gyrus', *Sci Adv*, 5: eaay6279.

- 869 Orf, Martin, Malte Wöstmann, Ronny Hannemann, and Jonas Obleser. 2022. 'Auditory neural  
870 tracking reflects target enhancement but not distractor suppression in a psychophysically  
871 augmented continuous-speech paradigm', *bioRxiv*: 2022.06.18.496558.
- 872 Pasley, B. N., S. V. David, N. Mesgarani, A. Flinker, S. A. Shamma, N. E. Crone, R. T. Knight,  
873 and E. F. Chang. 2012. 'Reconstructing speech from human auditory cortex', *PLoS Biol*,  
874 10: e1001251.
- 875 Peck, Fleming C, Laurel J Gabard-Durnam, Carol L Wilkinson, William Bosl, Helen Tager-  
876 Flusberg, and Charles A Nelson. 2021. 'Prediction of autism spectrum disorder diagnosis  
877 using nonlinear measures of language-related EEG at 6 and 12 months', *Journal of*  
878 *neurodevelopmental disorders*, 13: 1-13.
- 879 Pereira, F., B. Lou, B. Pritchett, S. Ritter, S. J. Gershman, N. Kanwisher, M. Botvinick, and E.  
880 Fedorenko. 2018. 'Toward a universal decoder of linguistic meaning from brain activation',  
881 *Nat Commun*, 9: 963.
- 882 Pion-Tonachini, L., K. Kreutz-Delgado, and S. Makeig. 2019. 'ICLabel: An automated  
883 electroencephalographic independent component classifier, dataset, and website',  
884 *Neuroimage*, 198: 181-97.
- 885 Rosen, S. 1992. 'Temporal information in speech: acoustic, auditory and linguistic aspects', *Philos*  
886 *Trans R Soc Lond B Biol Sci*, 336: 367-73.
- 887 Salisbury, D. F., M. E. Shenton, C. B. Griggs, A. Bonner-Jackson, and R. W. McCarley. 2002.  
888 'Mismatch negativity in chronic schizophrenia and first-episode schizophrenia', *Arch Gen*  
889 *Psychiatry*, 59: 686-94.
- 890 Shackleton, T. M., and R. P. Carlyon. 1994. 'The role of resolved and unresolved harmonics in  
891 pitch perception and frequency modulation discrimination', *J Acoust Soc Am*, 95: 3529-  
892 40.
- 893 Shannon, R. V., F. G. Zeng, V. Kamath, J. Wygonski, and M. Ekelid. 1995. 'Speech recognition  
894 with primarily temporal cues', *Science*, 270: 303-4.
- 895 Sohoglu, E., and M. H. Davis. 2020. 'Rapid computations of spectrotemporal prediction error  
896 support perception of degraded speech', *Elife*, 9.
- 897 Strauß, A., T. Wu, J. M. McQueen, O. Scharenborg, and F. Hintz. 2022. 'The differential roles of  
898 lexical and sublexical processing during spoken-word recognition in clear and in noise',  
899 *Cortex*, 151: 70-88.

- 900 Swaminathan, J., and M. G. Heinz. 2012. 'Psychophysiological analyses demonstrate the  
901 importance of neural envelope coding for speech perception in noise', *J Neurosci*, 32:  
902 1747-56.
- 903 Tang, C., L. S. Hamilton, and E. F. Chang. 2017. 'Intonational speech prosody encoding in the  
904 human auditory cortex', *Science*, 357: 797-801.
- 905 Teoh, E. S., F. Ahmed, and E. C. Lalor. 2022. 'Attention Differentially Affects Acoustic and  
906 Phonetic Feature Encoding in a Multispeaker Environment', *J Neurosci*, 42: 682-91.
- 907 Teoh, E. S., M. S. Cappelloni, and E. C. Lalor. 2019. 'Prosodic pitch processing is represented in  
908 delta-band EEG and is dissociable from the cortical tracking of other acoustic and phonetic  
909 features', *Eur J Neurosci*, 50: 3831-42.
- 910 Turner, C. W., B. J. Gantz, C. Vidal, A. Behrens, and B. A. Henry. 2004. 'Speech recognition in  
911 noise for cochlear implant listeners: benefits of residual acoustic hearing', *J Acoust Soc  
912 Am*, 115: 1729-35.
- 913 Vanthornhout, J., L. Decruy, J. Wouters, J. Z. Simon, and T. Francart. 2018. 'Speech Intelligibility  
914 Predicted from Neural Entrainment of the Speech Envelope', *J Assoc Res Otolaryngol*,  
915 19: 181-91.
- 916 Verschueren, E., J. Vanthornhout, and T. Francart. 2021. 'The effect of stimulus intensity on  
917 neural envelope tracking', *Hear Res*, 403: 108175.
- 918 Viswanathan, V., H. M. Bharadwaj, B. G. Shinn-Cunningham, and M. G. Heinz. 2021. 'Modulation  
919 masking and fine structure shape neural envelope coding to predict speech intelligibility  
920 across diverse listening conditions', *J Acoust Soc Am*, 150: 2230.
- 921 Yasmin, Sonia, Vanessa Irsik, Ingrid S Johnsrude, and Björn Herrmann. 2023. 'The Effects of  
922 Speech Masking on Neural Tracking of Acoustic and Semantic Features of Natural  
923 Speech', *bioRxiv*: 2023.02. 10.527537.
- 924 Zou, J., J. Feng, T. Xu, P. Jin, C. Luo, J. Zhang, X. Pan, F. Chen, J. Zheng, and N. Ding. 2019.  
925 'Auditory and language contributions to neural encoding of speech features in noisy  
926 environments', *Neuroimage*, 192: 66-75.
- 927
- 928