
Research Article: Methods/New Tools | Novel Tools and Methods

Cortical tracking of complex sound envelopes: modeling the changes in response with intensity

Denis P. Drennan¹ and Edmund C. Lalor^{1,2}

¹*School of Engineering, Trinity Centre for Bioengineering and Trinity College Institute of Neuroscience, Trinity College Dublin, Dublin 2, Ireland*

²*Department of Biomedical Engineering, Department of Neuroscience, and Del Monte Institute for Neuroscience, University of Rochester, 201 Robert B. Goergen Hall, Rochester, NY, 14627, USA*

<https://doi.org/10.1523/ENEURO.0082-19.2019>

Received: 7 March 2019

Revised: 2 May 2019

Accepted: 3 May 2019

Published: 6 June 2019

D.D. and E.L. designed research; D.D. performed research; D.D. analyzed data; D.D. and E.L. wrote the paper.

Funding: Irish Research Council
EPSPG/2014/54

;

Funding: Science Foundation Ireland (SFI)
CDA/15/3316

.

Conflict of Interest: The authors declare no competing financial interests.

Irish Research Council [EPSPG/2014/54]; Science Foundation Ireland (SFI) [CDA/15/3316]

Correspondence should be addressed to Edmund C. Lalor at edmund_lalor@urmc.rochester.edu

Cite as: eNeuro 2019; 10.1523/ENEURO.0082-19.2019

Alerts: Sign up at www.eneuro.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

Accepted manuscripts are peer-reviewed but have not been through the copyediting, formatting, or proofreading process.

Copyright © 2019 Drennan and Lalor

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license, which permits unrestricted use, distribution and reproduction in any medium provided that the original work is properly attributed.

1 Cortical tracking of complex sound envelopes:
2 modeling the changes in response with intensity

3
4 *Abbreviated Title:* Improved modeling of responses to complex
5 sounds

6
7 Denis P. Drennan¹, Edmund C. Lalor^{1,2}

8
9 ¹School of Engineering, Trinity Centre for Bioengineering and Trinity College Institute of
10 Neuroscience, Trinity College Dublin, Dublin 2, Ireland,

11 ²Department of Biomedical Engineering, Department of Neuroscience, and Del Monte
12 Institute for Neuroscience, University of Rochester, 201 Robert B. Goergen Hall, P.O. Box
13 270168, Rochester, NY 14627, USA.

14
15 Corresponding Author: Edmund C. Lalor PhD, Department of Biomedical Engineering,
16 Department of Neuroscience, and Del Monte Institute for Neuroscience, University of
17 Rochester, 201 Robert B. Goergen Hall, P.O. Box 270168, Rochester, NY 14627, USA
18 edmund_lalor@urmc.rochester.edu

19
20 Number of Figures: 7, Number of Tables: 0; Number of Multimedia: 0

21 Number of Words: *Abstract* – 241, *Significance Statement* – 116, *Introduction* – 732,
22 *Discussion* – 1,143

23
24 Conflict of Interest: The authors declare no competing financial interests.

25
26 Acknowledgements: This work was supported by the Irish Research Council through an
27 Enterprise Partnership Postgraduate Scholarship and by a Career Development Award from
28 Science Foundation Ireland (CDA/15/3316). The authors thank Dr. Kevin Prinsloo for
29 assistance with data collection, and Dr. Nathaniel Zuk for helpful comments on this
30 manuscript.

31 **Abstract**

32 Characterizing how the brain responds to stimuli has been a goal of sensory neuroscience for
33 decades. One key approach has been to fit linear models to describe the relationship between
34 sensory inputs and neural responses. This has included models aimed at predicting spike
35 trains, local field potentials, BOLD responses and EEG/MEG. In the case of EEG/MEG, one
36 explicit use of this linear modeling approach has been the fitting of so-called temporal
37 response functions (TRFs). TRFs have been used to study how auditory cortex tracks the
38 amplitude envelope of acoustic stimuli, including continuous speech. However, such linear
39 models typically assume that variations in the amplitude of the stimulus feature (i.e., the
40 envelope) produce variations in the magnitude but not the latency or morphology of the
41 resulting neural response. Here we show that by amplitude binning the stimulus envelope,
42 and then using it to fit a multivariate TRF, we can better account for these amplitude-
43 dependent changes, and that this leads to a significant improvement in model performance for
44 both amplitude modulated noise and continuous speech in humans. We also show that this
45 performance can be further improved through the inclusion of an additional envelope
46 representation that emphasizes onsets and positive changes in the stimulus, consistent with
47 the idea that while some neurons track the entire envelope, others respond preferentially to
48 onsets in the stimulus. We contend that these results have practical implications for
49 researchers interested in modeling brain responses to amplitude modulated sounds.

50 **Significance Statement**

51 A key approach in sensory neuroscience has been to fit linear models to describe the
52 relationship between stimulus features and neural responses. However, these linear models
53 often assume that the response to a stimulus feature will be consistent across its time course,
54 but just scaled linearly as a function of the stimulus feature's intensity. Here, using EEG in

55 humans, we show that allowing a linear model to vary as a function of the stimulus feature's
56 intensity leads to improved prediction of unseen neural data. We do so using both amplitude
57 modulated noise stimuli as well as continuous natural speech. This approach provides more
58 robust measures of envelope tracking and facilitates the study of its underlying mechanisms.

59 **Introduction**

60 Characterizing how the brain responds to stimuli has been a major goal of sensory
61 neuroscience for decades (Hubel and Wiesel, 1962). One key approach to this problem has
62 been to fit models to describe the relationship between sensory inputs and neural responses
63 (Wu et al., 2006). Our understanding of the sensory system can then be assessed by
64 quantifying how well such models predict neural responses to novel stimuli (Carandini et al.,
65 2005).

66 A central feature of such models has been a 'linear receptive field' stage that seeks to
67 account for some of the neural response as a linear weighted sum (i.e., a linear filter) of
68 particular features of the sensory input (e.g., the contrast of a visual stimulus across space or
69 the amplitude of an acoustic stimulus across time and frequency). In neural spiking models,
70 this linear filtering stage is typically just one of several stages (e.g., linear, nonlinear, and
71 Poisson) that seek to capture how stimulus variations are reflected in spike trains
72 (Chichilnisky, 2001). However, with macroscopic data like fMRI (Boynton et al., 1996) or
73 EEG/MEG (Crosse et al., 2016), this linear filtering stage often represents the entirety of the
74 model. Implicitly (or explicitly), these linear models assume that responses to a stimulus
75 feature will be temporally and morphologically consistent across its time-course, but just
76 scaled linearly as a function of the stimulus feature's intensity. In other words, they assume
77 that responses to a particular stimulus feature can be modeled by a linear impulse response
78 function.

79 In auditory neuroscience these filters have been viewed as representing the
80 ‘spectrotemporal receptive fields’ (STRFs) of auditory cortical neurons (Aertsen and
81 Johannesma, 1981). They are often fit using audio stimuli with broad spectrotemporal
82 statistics so as to characterize how neurons might respond to any sound (Depireux et al.,
83 2001), although there has been increasing interest in the use of more naturalistic stimuli such
84 as animal vocalizations (Theunissen et al., 2001; Machens et al., 2004) and human speech.
85 The latter has included efforts to fit linear response functions between various speech features
86 (e.g., envelope, spectrogram, phonemes, or phonetic features) and population responses in
87 animals (David et al., 2007; Mesgarani et al., 2008), as well as macroscopic measures in
88 humans (Lalor and Foxe, 2010; Ding and Simon, 2012; Di Liberto et al., 2015).

89 One explicit use of this linear modeling approach has been the fitting of so-called
90 temporal response functions (TRFs) to describe how EEG is affected by variations in visual
91 (Gonçalves et al., 2014) or auditory (Lalor et al., 2009) stimuli. This includes univariate
92 TRFs that model how EEG changes based on a single stimulus feature (e.g., an envelope),
93 and multivariate TRFs that simultaneously model responses to multiple features (e.g., a
94 spectrogram; Ding and Simon, 2012; Di Liberto et al., 2015). In both cases, however, it is
95 typically assumed that changes in the intensity of the stimulus feature produce variations in
96 the magnitude but not the latency or morphology of the responses. While this assumption
97 may be reasonable for certain brain responses in certain brain areas to certain stimulus
98 features (Boynton et al., 1996), there is definitive evidence that it is imperfect for EEG-based
99 TRFs.

100 One such piece of evidence is the long-known relationship between auditory stimulus
101 amplitude and response latency (Beagley and Knight, 1967). Specifically, while there is a
102 monotonic (although not necessarily linear) relationship between auditory stimulus amplitude
103 and response magnitude, there is also an inverse relationship between stimulus amplitude and

104 response latency. Therefore, to model neural responses to an ongoing auditory stimulus using
105 a linear univariate TRF is likely to be suboptimal given that it ignores the dependence of
106 response latency (and morphology) on stimulus amplitude.

107 Here, we aim to demonstrate that by allowing the stimulus-response model to vary as
108 a function of the stimulus amplitude, we can improve the modeling of responses to
109 continuous auditory stimuli. To do so, we propose a simple extension to the standard linear
110 TRF estimation approach that involves amplitude binning a *single* feature, namely the
111 envelope, and then using it to fit a *multivariate* TRF. This should allow the TRF to vary
112 across the different amplitude ranges, thus enabling it to account for associated changes in
113 response magnitude, latency, and morphology. We aim to validate that this represents an
114 improved model by comparing how well it predicts EEG data relative to more standard
115 univariate models and discuss other methods that can be used to improve model performance.

116 **Materials and Methods**

117 EEG data from two experiments were used in this study: one acquired in response to
118 amplitude modulated broadband noise (AM BBN), the other in response to continuous
119 natural speech (Natural Speech Dataset from <https://doi.org/10.5061/dryad.070jc>, including
120 amplitude envelopes; Broderick et al., 2018).

121 *Subjects*

122 13 subjects participated in the AM BBN experiment; 5 male, aged 23-35 years. 19 subjects
123 participated in the speech experiment; 13 male, aged 19-38 years, although data from 2
124 subjects were later excluded because of uncertainties in response timing due to differences in
125 their data acquisition setup. All subjects had self-reported normal hearing. The protocol for

126 both studies was approved by the Ethics Committee of the Health Sciences Faculty at Trinity
127 College Dublin, Ireland, and all subjects gave written informed consent.

128 *Stimuli*

129 As mentioned, this study involved experiments using two different types of stimuli, AM BBN
130 and continuous natural speech.

131 The carrier signal for the AM BBN stimulus was uniform broadband noise with
132 energy limited to a bandwidth of 0–24000 Hz. Its modulating signal (envelope) had a log-
133 uniform amplitude distribution (by design, although less so after envelope extraction, please
134 see below) and a bottom-heavy (right-skewed) frequency (modulation rate) distribution
135 (Figure 1), so chosen as it has been shown that auditory cortical areas tend to be most
136 sensitive to AM BBN presented at lower modulation frequencies (Liégeois-Chauvel et al.,
137 2004). The envelope was created by first generating a signal with discrete amplitude values
138 with the desired statistical properties, and then interpolating between those discrete points to
139 provide a smooth transition from one modulation amplitude to the next.

140 The speech stimulus had a bottom-heavy (right-skewed) frequency distribution with
141 energy limited to a bandwidth of 0–22050 Hz. Its envelope had a log-top-heavy (left-skewed)
142 amplitude distribution, and a bottom-heavy (right-skewed) frequency (modulation rate)
143 distribution, similar to that of the AM BBN stimulus (Figure 1). It comprised extracts from a
144 professional audio-book version of a popular mid-20th century American work of fiction (i.e.,
145 ‘The Old Man and the Sea’ by Ernest Hemingway) written in an economical and understated
146 style and read by a single male American speaker.

147 *Experimental Procedure*

148 In the AM BBN experiment, subjects were presented with 80 repetitions of the same 60 s
149 long AM BBN stimulus as they reclined in a comfortable chair, in a quiet, darkened room,
150 and watched a silent animated cartoon presented on a tablet computer. They were asked *not*
151 to attend to the auditory stimuli, which were presented monaurally to their right ear at a peak
152 level equivalent to that of a 1 kHz pure-tone at 80 dB SPL, using a Sound Blaster X-Fi
153 Surround 5.1 Pro external sound card, a TPA3118D2EVM amplifier, and electromagnetically
154 shielded Etymotic Research ER-2 earphones, via VLC Media Player from VideoLan
155 (<http://www.videolan.org>). Compensation for the 1 ms sound-tube delay introduced by the
156 ER-2 earphones was applied post-hoc.

157 In the speech experiment, subjects were presented with 28 trials of ~155s long
158 audiobook extracts. The trials preserved the storyline, with neither repetitions nor
159 discontinuities. Subjects sat in a comfortable chair, in a quiet, darkened room, and were
160 instructed to maintain visual fixation on a crosshair centered on a computer monitor, and to
161 minimize eye blinking and all other motor activities for the duration of each trial. They were
162 asked to attend to the auditory stimuli, which were presented diotically at a comfortable
163 listening level, using Sennheiser HD 650 headphones, via Presentation software from
164 Neurobehavioral Systems (<http://www.neurobs.com>). For the purposes of analysis, all trials
165 were truncated to 150 s, and a peak level of 80 dB SPL was estimated (as the original
166 presentation level was not available).

167 *EEG Acquisition*

168 In the AM BBN experiment, 40 channels of EEG data were recorded at 16384 Hz (analog -3
169 dB point of 3276.8 Hz), using a BioSemi ActiveTwo system (<http://www.biosemi.com>). 32

170 cephalic electrodes were positioned according to the standard 10-20 system. A further eight
171 non-cephalic electrodes were also collected although only two – those over the left and right
172 mastoids – were used in the analysis. Triggers indicating the start of each 60 s trial were
173 encoded in a separate channel in the stimulus WAV file as three cycles of a 16 kHz tone
174 burst. These triggers were interpreted by custom hardware before being fed into the
175 acquisition laptop for synchronous recording along with the EEG.

176 In the speech experiment, 130 channels of EEG data were recorded at 512 Hz (analog
177 -3 dB point of 409.6 Hz), using a BioSemi ActiveTwo system. 128 cephalic electrodes were
178 positioned according to the BioSemi Equiradial system, with another 2 electrodes located
179 over the left and right mastoids. Triggers indicating the start of each ~155 s trial were
180 presented using Neurobehavioral Systems Presentation software for synchronous recording
181 along with the EEG.

182 *EEG Preprocessing*

183 The EEG data were first resampled to 128 Hz using the *decimate* function in MATLAB
184 (<http://www.mathworks.com>). The *decimate* function incorporates an 8th order low-pass
185 Chebyshev Type I infinite impulse response (IIR) anti-aliasing filter. This filter was applied
186 with a cutoff frequency of 64 Hz and was implemented using the *filtfilt* function, ensuring
187 zero phase distortion and in effect doubling the order of the filter. A 1st order high-pass
188 Butterworth filter was then applied with a cutoff frequency of 1 Hz, also using the *filtfilt*
189 function. Bad channels were determined as those whose variance was either less than half or
190 greater than twice that of the surrounding 2–4 channels for the AM BBN dataset, and 3–7
191 channels for the speech dataset (depending on location). These were then replaced through
192 spherical spline interpolation using EEGLAB (Delorme and Makeig, 2004). Finally, the data

193 were rereferenced to the average of the mastoids, separated into trials based on the triggers
194 provided, and z-scored.

195 *Temporal Response Function Estimation*

196 The models were fit using TRF estimation implemented via the mTRF Toolbox (Crosse et al.,
197 2016). With TRF estimation, the assumption is that the output EEG, $y(t)$, consists of the
198 convolution of a particular input stimulus feature, $x(t)$, with an unknown system response
199 $w(\tau)$ (i.e., the TRF), plus noise (Lalor et al., 2009), i.e.,

$$y(t) = w(\tau) * x(t) + noise$$

200 where τ represents the range of time-lags over which the TRF is estimated. Given the known
201 stimulus feature and the measured EEG, the TRF can be derived (in this case) by performing
202 regularized linear (ridge) regression (see Crosse et al., 2016 for details). Baseline correction
203 was performed on each subject's average TRF (by subtracting the mean value between -20
204 and 0 ms) before being combined to form the grand average.

205 *Amplitude Binned Envelope*

206 The choice of stimulus feature can have a significant influence on the resulting model. Such
207 features could include the envelope (Lalor et al., 2009) or spectrogram (Ding and Simon,
208 2012; Di Liberto et al., 2015), or in the case of speech, phonemes, phonetic features
209 (Di Liberto et al., 2015), or its semantic content (Broderick et al., 2018). The envelope (time
210 x amplitude) however is probably the most commonly used stimulus feature and is the one
211 chosen for use in this study. For both the AM BBN and speech stimuli the envelopes were
212 calculated by taking the absolute value of their Hilbert transforms, and then resampling them
213 to 128 Hz using the *decimate* function in MATLAB.

214 As mentioned, it has long been known that the magnitude and latency of auditory
215 system responses vary directly and inversely with stimulus amplitude, respectively, i.e., as
216 the stimulus amplitude increases, the response magnitude increases, and the response latency
217 decreases (and vice versa). Univariate TRFs, like those modeled using envelopes, cannot
218 account for all these amplitude-dependent changes. In fact, univariate TRFs can only account
219 for linear changes in magnitude and cannot account for any changes in latency or
220 morphology. However, by simply amplitude binning the envelope (time x [amplitude] x
221 amplitude), i.e., by dividing the envelope up into multiple sub-envelopes comprising the
222 different amplitude ranges of the full envelope, normalizing the values in each bin to be
223 between 0 and 1, and then using it to fit a multivariate TRF, should allow the TRF to vary
224 across the different amplitude ranges, potentially enabling it to account for more of these
225 amplitude-dependent changes than its univariate counterpart.

226 The amplitude binned (AB) envelope was created by logarithmically binning the
227 envelope into 8 dB bins using the *histcounts* function in MATLAB, and then normalizing the
228 values in each bin to between 0 and 1 (an important step in ensuring the stability of the
229 resulting TRF). This bin size was chosen empirically after comparing the prediction
230 accuracies attained across a range of bin sizes, with broader bins perhaps being less able to
231 capture changes in the response with amplitude, and narrower bins perhaps suffering from the
232 limited amount of data available for training. The logarithmic bin edges were determined by
233 taking 10 to the power of the desired bin edges in dB (i.e., 8, 16, 24, etc.) divided by 20, and
234 then normalizing the resulting range to between 0 and 1 (Figure 2B).

235 *Other Stimulus Representations*

236 A number of other approaches have already been put forward that attempt to modify the
237 stimulus representation in order to account for certain properties of the auditory system. So,

238 rather than just comparing the AB envelope model with the standard envelope model, we also
239 chose to compare it with two others, i.e., the SPL envelope and onset envelope models. The
240 SPL envelope model was fit using an envelope that was transformed into its equivalent
241 logarithmic (sound pressure level; SPL) representation, and the onset envelope model was fit
242 using an envelope that was modified to place a greater emphasis on onsets and positive
243 changes in amplitude.

244 The motivation for using the SPL envelope model derives from the well-known fact
245 that electrophysiological responses generally vary in proportion to the log of the stimulus
246 amplitude (Aiken and Picton, 2008). The SPL envelope was generated by taking 20 times the
247 base 10 logarithm of the envelope (Aiken and Picton, 2008; Figure 2A), and it was hoped that
248 this would help linearize the amplitude to magnitude mapping between the stimulus
249 representation and the EEG. It was presumed that the AB envelope model might outperform
250 the SPL envelope model however, given that they both attempt to account for nonlinearities
251 in the relationship between stimulus amplitude and response magnitude, but only the former
252 accounts for changes in response latency and morphology.

253 The motivation for using the onset envelope model comes from the idea that many
254 auditory neurons are particularly sensitive to onsets, offsets, and changes in the stimulus
255 (Bieser and Müller-Preuss, 1996), and that this approach has been used effectively in the past
256 (Aiken and Picton, 2008; Hertrich et al., 2012; Fiedler et al., 2017). The onset envelope was
257 explicitly designed to reflect onsets and positive changes in the stimulus, and was created by
258 half-wave rectifying the first-derivative of the envelope (Hertrich et al., 2012; Figure 2A).

259 *Experimental Design and Statistical Analyses*

260 In order to compare the different models tested as part of this study, a nested ‘leave-one-out’
261 cross-validation approach was employed. Specifically, for each stimulus representation, a

262 separate TRF (univariate for the envelope, SPL envelope, and onset envelope, and
263 multivariate for the AB envelope) was fit for each of M trials across several ridge parameters
264 (usually denoted λ) used to regularize the models. One trial was then chosen to be ‘left
265 out’, i.e., to be used as a ‘test set’, with the remaining $M-1$ trials to be used for the inner
266 cross-validation. Of these inner $M-1$ trials, one trial was again chosen to be ‘left out’, i.e., to
267 be used as a ‘validation set’, with the remaining $M-2$ trials to be used as a ‘training set’.

268 For each λ value, an average model was obtained by averaging over the single-
269 trial models in the ‘training set’. These were then convolved with the stimulus representation
270 associated with the ‘validation set’ to predict its neural response. Model performance was
271 assessed by quantifying how accurately these predicted responses matched the actual
272 recorded response from the ‘validation set’, using Pearson’s correlation coefficient. This
273 process was then repeated $M-2$ times such that each trial was ‘left out’ of the ‘training set’
274 once. The overall model performance was then determined by averaging over the individual
275 model performances for each trial, and the optimal λ value was chosen.

276 Using this optimal λ value, another average model was then obtained by
277 averaging over the single-trial models in both the ‘training’ and ‘validation’ sets. This was
278 then convolved with the stimulus representation associated with the ‘test set’ to predict its
279 neural response. Model performance was then assessed by quantifying how accurately the
280 predicted response matched the actual recorded response from the ‘test set’. This entire
281 procedure was then repeated $M-1$ times such that each trial was ‘left out’ of the inner cross-
282 validation procedure once. The overall model performance was then finally determined by
283 averaging over the individual model performances for each trial. Importantly, the parameter
284 optimization was done separately for each stimulus representation and subject, so that we
285 were left comparing each model based on its respective optimal performance.

286 Again, the performance of each model was assessed by quantifying how accurately
287 the predicted response matched the actual recorded response, using Pearson's correlation
288 coefficient. The normality of these performance measures for each model was confirmed
289 using the Anderson Darling test, and model comparisons were carried out using paired-
290 sample t-tests and Cohen's d effect size for paired-sample t-tests. Cohen's d effect size was
291 calculated by dividing each t-value by the square-root of the sample size. One potential
292 concern when comparing models with different numbers of parameters is that models with
293 more parameters may perform better simply due to their greater complexity. To account for
294 this, supplementary comparisons were also carried out using the Akaike Information
295 Criterion (AIC) which penalizes models based on their complexity. As the results of these
296 analyses were not normal, model comparisons were carried out using Wilcoxon signed-rank
297 tests.

298 Permutation tests were also used to assess the null distributions of the envelope
299 models. For the AM BBN dataset, as the stimulus was the same for each trial, a pool of 80
300 circularly-shifted envelopes (i.e., the original envelope plus 79 circularly-shifted envelopes,
301 each iteratively shifted by $1/80$ times the length of the envelope with respect to the previously
302 shifted envelope) were first created. 80 envelopes from this pool were then chosen at random
303 with replacement for use in the cross-validation procedure. This selection and cross-
304 validation procedure was repeated 100 times to determine the null-distribution of the
305 envelope model for each subject. For the speech dataset, as the stimuli were different for each
306 trial, envelopes were simply chosen at random with replacement from the original set of
307 envelopes, for use in the cross-validation procedure. This selection and cross-validation
308 procedure was also repeated 100 times to determine the null-distribution of the envelope
309 model for each subject.

310 **Results**311 *Channel Selection*

312 EEG prediction accuracies will vary across channels depending on how related the data on
313 those channels are to the stimulus representation. For the AM BBN analyses, the seven
314 channels (of 32) with the highest prediction accuracies for the envelope model were used
315 (Figure 3A). For the speech analyses, the 42 channels (of 128) with the highest prediction
316 accuracies for the envelope model, plus three other channels (to ensure symmetry) were used
317 (Figure 5A). In both cases, these channels tended to reside over fronto-central to temporal
318 scalp (see Di Liberto et al., 2015). The overall prediction accuracy for each model was
319 calculated by averaging the prediction accuracies over these electrodes.

320 *Individual Model Comparisons*321 AM BBN

322 For the AM BBN dataset, prediction accuracies were determined for each model, and each
323 subject (Figure 3B). All four stimulus representations (i.e., envelope, SPL envelope, onset
324 envelope, and AB envelope) and their associated models were able to predict EEG responses
325 with an accuracy that was significantly above 0.0012, i.e., the null hypothesis obtained using
326 the permutation tests, for all subjects ($t(12)$, all $p < 0.001$), and greater than all values
327 obtained using the permutation tests. However, the AB envelope model significantly
328 outperformed all three of the other models, in each case with a large to very large positive
329 effect size ($t(12) = 5.471$, $p < 0.001$, $d = 1.518$ vs. the envelope model; $t(12) = 4.070$, $p <$
330 0.01 , $d = 1.129$ vs. the SPL envelope model; $t(12) = 4.800$, $p < 0.001$, $d = 1.331$ vs. the onset
331 envelope model). These results were also seen when comparing the models using AIC ($p <$

332 0.001; Wilcoxon signed rank test). Neither the SPL envelope nor onset envelope models
333 managed to outperform the standard envelope model ($t(12)$, both $p > 0.05$).

334 Exactly how the TRF changes as a function of stimulus amplitude becomes more
335 apparent on closer inspection of the AB envelope TRF (Figure 4A). As the stimulus
336 amplitude decreases, the TRF magnitude decreases, latency increases, and morphology
337 changes in accordance with our hypothesis. The influence of stimulus amplitude on TRF
338 latency is perhaps better emphasized in Figure 4B. For example, the “N1”, which is quite
339 large in magnitude in the uppermost amplitude bin, decreases in magnitude and increases in
340 latency, with decreasing stimulus amplitude. To quantify this relationship, the N1 peak in
341 each bin was determined as being the largest negative peak in the TRF at lags between 70 and
342 210 ms (the corresponding latencies can be seen in Figure 4C). A line was then fit to the data
343 (R -squared = 0.9143, $p < 0.001$), which showed that the N1 peak latency increases by ~11 ms
344 with every unit decrease in amplitude bin. A similar effort was made to quantify the
345 relationship between stimulus amplitude and TRF magnitude, i.e., the “P1” peak in each bin
346 was determined as being the largest positive peak in the TRF at lags between 0 to 130 ms (the
347 corresponding P1-N1 peak-peak amplitudes can be seen in Figure 4D). However, while there
348 does seem to be some relationship between stimulus amplitude and TRF magnitude, it was
349 not well fit by a line (R -squared = 0.597, $p < 0.01$).

350 Speech

351 For the speech dataset, prediction accuracies were again determined for each model, and each
352 subject, with very similar results to before (Figure 5B). All four stimulus representations and
353 their associated models were able to predict EEG responses with an accuracy that was
354 significantly above 0.0015, i.e., the null-hypothesis obtained using the permutation tests, for
355 all subjects ($t(16)$, all $p < 0.001$), and greater than all values obtained using the permutation

356 tests. The AB envelope model significantly outperformed all three of the other models, in
357 each case with a large to very large positive effect size ($t(16) = 5.472$, $p < 0.001$, $d = 1.327$
358 vs. the envelope model; $t(16) = 7.649$, $p < 0.001$, $d = 1.855$ vs. the SPL envelope model; $t(16)$
359 $= 4.666$, $p < 0.001$, $d = 1.132$ vs. the onset envelope model). These results were also seen
360 when comparing the models using AIC ($p < 0.001$; Wilcoxon signed rank test). Neither the
361 SPL envelope nor onset envelope models managed to outperform the standard envelope
362 model ($t(16)$, both $p > 0.05$).

363 Again, exactly how the TRF changes as a function of stimulus amplitude becomes
364 more apparent on closer inspection of the AB envelope TRF (Figure 6A and B). While the
365 overall relationship between stimulus amplitude and TRF magnitude, latency, and
366 morphology appears similar to before, in this case, the magnitude of the TRF for some of the
367 lower amplitude bins seems unexpectedly high. It is not entirely clear why this would have
368 been the case. The “P1” and “N1” peaks were also determined in the same manner as before
369 (and the corresponding “N1” latencies and “P1-N1” peak-peak amplitudes can be seen in
370 Figure 6C and D respectively). To quantify the relationship between stimulus amplitude and
371 TRF latency, a line was fit to the N1 latency data ($R\text{-squared} = 0.5008$, $p < 0.05$), which
372 again showed that the N1 peak latency increases by ~ 11 ms with every unit decrease in
373 amplitude bin. However, there was no simple relationship between stimulus amplitude and
374 TRF magnitude.

375 *Combined Model Comparisons*

376 The comparison between the AB envelope and onset envelope models is not necessarily as
377 straightforward as one might expect. This is because each model is likely reflecting different
378 envelope tracking mechanisms in the cortex (Bieser and Müller-Preuss, 1996). Specifically,
379 the onset envelope model likely reflects contributions from neurons that track onsets and

380 positive changes in amplitude while the AB envelope model likely reflects contributions from
381 neurons that track along with all of the amplitude fluctuations (Bieser and Müller-Preuss,
382 1996).

383 To test the idea that these two models are capturing complementary information on
384 envelope tracking, we investigated whether there would be any advantage in combining these
385 two models (i.e., by combining the two stimulus representations and then using that to fit a
386 multivariate TRF). Indeed, the combined AB envelope plus onset envelope model
387 significantly outperformed the individual onset envelope and AB envelope models, for both
388 the AM BBN ($t(12) = 5.717$, $p < 0.001$, $d = 1.586$ vs. the onset envelope; $t(12) = 4.184$, $p <$
389 0.01 , $d = 1.161$ vs. the AB envelope) and speech datasets ($t(16) = 6.139$, $p < 0.001$, $d = 1.489$
390 vs. the onset envelope; $t(16) = 3.312$, $p < 0.01$, $d = 0.8032$ vs. the AB envelope), suggesting
391 that they are capturing complementary information on envelope tracking in the cortex (Figure
392 7A and B). These results were also seen when comparing the models using AIC (all p
393 < 0.001 ; Wilcoxon signed rank test).

394 One obvious extension to this approach then might be to also amplitude bin the onset
395 envelope representation, producing an ‘AB onset envelope’ model. However, while this AB
396 onset envelope TRF exhibits a similar dependence on stimulus amplitude to the AB envelope
397 TRF (Figure 7C and D) and significantly outperformed the onset envelope model alone for
398 the AM BBN dataset ($t(12) = 3.887$, $p < 0.05$, $d = 1.078$) although not for the speech dataset
399 ($t(16)$, $p > 0.05$), the combined AB envelope plus AB onset envelope model failed to
400 outperform the combined AB envelope plus onset envelope model for either the AM BBN or
401 speech datasets.

402 **Discussion**

403 Despite it long being known that the latency and morphology (and not just the magnitude) of
404 auditory system responses are dependent on the stimulus amplitude, this has been overlooked
405 in previous efforts at linearly modeling the auditory system. Here we have shown that by
406 allowing the stimulus-response model to vary as a function of the stimulus amplitude, we can
407 improve the modeling of responses to continuous auditory stimuli.

408 Specifically, we saw that by amplitude binning the envelope and then using that to fit
409 a multivariate TRF, we could improve the prediction accuracy over the standard envelope
410 model with a very large effect size for both the AM BBN and speech datasets. This was not
411 the case for the SPL envelope or onset envelope models however, which both failed to
412 outperform the standard envelope model. We also evaluated the offset envelope (created by
413 half-wave rectifying the negative portion of the first-derivative of the envelope and then
414 using that to fit a univariate TRF) and derivative envelope models but again, neither managed
415 to outperform the envelope model for either dataset and indeed mostly performed worse.
416 Finally, we saw that by combining the AB envelope and onset envelope models, we could
417 further improve the prediction accuracy over the AB envelope model with a large effect size
418 for both the AM BBN and speech datasets.

419 Interestingly, despite having lower prediction accuracies overall, the improvement in
420 prediction accuracy was greater for the AM BBN dataset. This is likely due to the differences
421 in amplitude distribution seen between the two types of stimuli (Figure 1E and F). While the
422 speech stimuli predominantly vary within a narrow amplitude range, the wider ‘active’
423 amplitude range of the AM BBN stimulus may allow it to benefit more from taking
424 amplitude-dependent variations into account. The reason that the prediction accuracies were
425 higher for the speech dataset overall is likely due to attention effects, e.g., as seen in

426 O’Sullivan et al., 2014 (albeit it in that case with two competing speech streams and
427 reconstruction accuracy).

428 Previous work has shown that the use of other stimulus representations can also
429 improve modeling performance. For example, for speech it has been shown that models
430 based on spectrograms, phonemes, and phonetic features, outperform those based on the
431 standard envelope (Di Liberto et al., 2015). However, for each of these stimulus
432 representations, the same assumption of unchanging TRF morphology applies. For
433 categorical representations such as those reflecting the phonemic/phonetic content of the
434 speech, this could be considered a strength, but for lower-level representations such as the
435 spectrogram, this could be considered a weakness. In the same way as we have done for the
436 envelope in this study, an amplitude binning approach could also be applied to the
437 spectrogram representation of speech by binning the stimulus in each frequency band (time x
438 [frequency x [amplitude]] x amplitude), and then using it to fit a multivariate TRF. This could
439 then potentially account for both amplitude- and frequency- dependent changes in the
440 response, which could lead to improved model performance. Furthermore, as before, the
441 onset envelope could also be included to explain even more of the variance. That said, it
442 should be noted that this representation would be high-dimensional and so would come with
443 increased computational requirements as well as an increased chance of overfitting.

444 Researchers interested in improving the performance of their envelope tracking
445 measures could benefit from using the AB envelope approach and/or including the onset
446 envelope as part of their stimulus representation. The sensitivity and robustness of such
447 measures could be further improved however, if this work was adapted into a “decoding”
448 framework. Such approaches have become quite popular in recent years and often involve
449 mapping backwards from the multivariate neural data to reconstruct an estimate of the
450 univariate speech envelope that caused those data (O’Sullivan et al., 2014). This approach

451 takes advantage of the large increase in modeling performance that comes with incorporating
452 all of the neural data simultaneously in one multivariate mapping. This stands in contrast with
453 the forward channel-by-channel modeling approach we have used in the present study. As
454 such, it would be practically valuable to incorporate the AB envelope approach into a
455 multivariate-to-multivariate decoding framework. While we have not done that here, such
456 frameworks have been implemented before for multivariate auditory stimuli (e.g., Mesgarani
457 et al., 2009) and there are several flexible methods available that would be well suited to such
458 a task (e.g., de Cheveigné et al., 2018).

459 Finally, we suggest that the results of our study should factor into theories on the
460 generative mechanisms underlying the cortical tracking of acoustic envelopes. There are at
461 least two such theories. One proposes that intrinsic, ongoing oscillatory brain rhythms
462 “entrain” to the rhythms of the speech signal by aligning their phase with the stimulus in an
463 anticipatory, behaviorally effective manner (Giraud and Poeppel, 2012; Rimmele et al.,
464 2018). An alternative idea is that cortical tracking of speech (or any auditory stimulus for that
465 matter) occurs as a result of the stimulus providing a driving input to auditory cortex that
466 evokes transient responses in neuronal populations that are tuned to the features of that
467 stimulus and that scale with the strength of those features. It is well known that sensory
468 neurons are tuned to certain features of the stimuli that they encounter – including features
469 such as frequency and intensity in the auditory domain (Phillips and Irvine, 1981). As such,
470 researchers have explicitly modeled cortical tracking of the speech envelope as a series of
471 transient responses to changes in that speech envelope (Aiken and Picton, 2008). Indeed, this
472 assumption is at the core of the TRF analysis used in this paper. Moreover, in other work, we
473 have shown that EEG responses to continuous speech are well modeled as a series of
474 transient responses to changes in frequency and phonetic features within the speech
475 (Di Liberto et al., 2015). While the present study cannot definitively adjudicate between

476 oscillatory entrainment and transient evoked responses as the underlying mechanism – and
477 indeed maybe both are at play – we do suggest that the relationship we have shown between
478 stimulus amplitude and response latency needs to be considered when positing one or other of
479 these mechanisms. Do smaller amplitude changes evoke later and slower transient responses?
480 Or do they entrain slower oscillations? Or some combination of the two? Work in our lab has
481 begun to look directly at this issue (Lalor, 2019) and will continue to do so.

482 In summary, here we have shown that by allowing the stimulus-response model to vary as
483 a function of the stimulus amplitude, we can improve the modeling of responses to
484 continuous auditory stimuli, and that the inclusion of an onset stimulus representation can
485 improve this performance even further. This obviously has implications for how people
486 model auditory processing in humans, but, more generally, points to the importance of
487 incorporating stimulus-dependencies when modeling the activity of sensory systems.

488 **References**

- 489 Aertsen AMHJ, Johannesma PIM (1981) The Spectro-Temporal Receptive Field. *Biol*
490 *Cybern* 42:133–143.
- 491 Aiken SJ, Picton TW (2008) Human Cortical Responses to the Speech Envelope. *Ear Hear*
492 29:139–157.
- 493 Beagley HA, Knight JJ (1967) Changes in Auditory Evoked Response with Intensity. *J*
494 *Laryngol Otol* 81:861–873.
- 495 Bieser A, Müller-Preuss P (1996) Auditory responsive cortex in the squirrel monkey: neural
496 responses to amplitude-modulated sounds. *Exp Brain Res* 108:273–284.
- 497 Boynton GM, Engel SA, Glover GH, Heeger DJ (1996) Linear systems analysis of functional
498 magnetic resonance imaging in human V1. *J Neurosci* 16:4207–4221.
- 499 Broderick MP, Anderson AJ, Liberto GMD, Crosse MJ, Lalor EC (2018)
500 Electrophysiological Correlates of Semantic Dissimilarity Reflect the Comprehension
501 of Natural, Narrative Speech. *Curr Biol* 28:803-809.e3.
- 502 Carandini M, Demb JB, Mante V, Tolhurst DJ, Dan Y, Olshausen BA, Gallant JL, Rust NC
503 (2005) Do we know what the early visual system does? *J Neurosci* 25:10577–10597.

- 504 Chichilnisky EJ (2001) A simple white noise analysis of neuronal light responses. *Netw*
505 *Comput Neural Syst* 12:199–213.
- 506 Crosse MJ, Di Liberto GM, Bednar A, Lalor EC (2016) The Multivariate Temporal Response
507 Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to
508 Continuous Stimuli. *Front Hum Neurosci* 10.
- 509 David SV, Mesgarani N, Shamma SA (2007) Estimating sparse spectro-temporal receptive
510 fields with natural stimuli. *Netw Comput Neural Syst* 18:191–212.
- 511 de Cheveigné A, Wong DDE, Di Liberto GM, Hjørtkjær J, Slaney M, Lalor E (2018)
512 Decoding the auditory brain with canonical component analysis. *NeuroImage*
513 172:206–216.
- 514 Delorme A, Makeig S (2004) EEGLAB: an open source toolbox for analysis of single-trial
515 EEG dynamics including independent component analysis. *J Neurosci Methods*
516 134:9–21.
- 517 Depireux DA, Simon JZ, Klein DJ, Shamma SA (2001) Spectro-Temporal Response Field
518 Characterization With Dynamic Ripples in Ferret Primary Auditory Cortex. *J*
519 *Neurophysiol* 85:1220–1234.
- 520 Di Liberto GM, O’Sullivan JA, Lalor EC (2015) Low-Frequency Cortical Entrainment to
521 Speech Reflects Phoneme-Level Processing. *Curr Biol* 25:2457–2465.
- 522 Ding N, Simon JZ (2012) Neural coding of continuous speech in auditory cortex during
523 monaural and dichotic listening. *J Neurophysiol* 107:78–89.
- 524 Fiedler L, Wöstmann M, Graversen C, Brandmeyer A, Lunner T, Obleser J (2017) Single-
525 channel in-ear-EEG detects the focus of auditory attention to concurrent tone streams
526 and mixed speech. *J Neural Eng* 14:036020.
- 527 Giraud A-L, Poeppel D (2012) Cortical oscillations and speech processing: emerging
528 computational principles and operations. *Nat Neurosci* 15:511.
- 529 Gonçalves NR, Whelan R, Foxe JJ, Lalor EC (2014) Towards obtaining spatiotemporally
530 precise responses to continuous sensory stimuli in humans: a general linear modeling
531 approach to EEG. *Neuroimage* 97:196–205.
- 532 Hertrich I, Dietrich S, Trouvain J, Moos A, Ackermann H (2012) Magnetic brain activity
533 phase-locked to the envelope, the syllable onsets, and the fundamental frequency of a
534 perceived speech signal. *Psychophysiology* 49:322–334.
- 535 Hubel DH, Wiesel TN (1962) Receptive fields, binocular interaction and functional
536 architecture in the cat’s visual cortex. *J Physiol* 160:106–154.
- 537 Lalor EC. Evoked activity plays a very substantial role in the cortical tracking of natural
538 speech. In: Annual Meeting of the Cognitive Neuroscience Society; 2019 Mar 23-26;
539 San Francisco, California (CA) Available at:
540 [https://www.cogneurosociety.org/2019/wordpress/wp-content/uploads/2019/03/CNS-](https://www.cogneurosociety.org/2019/wordpress/wp-content/uploads/2019/03/CNS-2019-Abstract-Book.pdf)
541 [2019-Abstract-Book.pdf](https://www.cogneurosociety.org/2019/wordpress/wp-content/uploads/2019/03/CNS-2019-Abstract-Book.pdf). Abstract nr C111.

- 542 Lalor EC, Foxe JJ (2010) Neural responses to uninterrupted natural speech can be extracted
543 with precise temporal resolution. *Eur J Neurosci* 31:189–193.
- 544 Lalor EC, Power AJ, Reilly RB, Foxe JJ (2009) Resolving Precise Temporal Processing
545 Properties of the Auditory System Using Continuous Stimuli. *J Neurophysiol*
546 102:349–359.
- 547 Liégeois-Chauvel C, Lorenzi C, Trébuchon A, Régis J, Chauvel P (2004) Temporal Envelope
548 Processing in the Human Left and Right Auditory Cortices. *Cereb Cortex* 14:731–
549 740.
- 550 Machens CK, Wehr MS, Zador AM (2004) Linearity of Cortical Receptive Fields Measured
551 with Natural Sounds. *J Neurosci* 24:1089–1100.
- 552 Mesgarani N, David SV, Fritz JB, Shamma SA (2008) Phoneme representation and
553 classification in primary auditory cortex. *J Acoust Soc Am* 123:899–909.
- 554 Mesgarani N, David SV, Fritz JB, Shamma SA (2009) Influence of Context and Behavior on
555 Stimulus Reconstruction From Neural Activity in Primary Auditory Cortex. *J*
556 *Neurophysiol* 102:3329–3339.
- 557 O’Sullivan JA, Power AJ, Mesgarani N, Rajaram S, Foxe JJ, Shinn-Cunningham BG, Slaney
558 M, Shamma SA, Lalor EC (2014) Attentional Selection in a Cocktail Party
559 Environment Can Be Decoded from Single-Trial EEG. *Cereb Cortex*.
- 560 Parker AJ, Newsome WT (1998) Sense and the Single Neuron: Probing the Physiology of
561 Perception. *Annu Rev Neurosci* 21:227–277.
- 562 Phillips DP, Irvine DR (1981) Responses of single neurons in physiologically defined
563 primary auditory cortex (AI) of the cat: frequency tuning and responses to intensity. *J*
564 *Neurophysiol* 45:48–58.
- 565 Rimmele JM, Morillon B, Poeppel D, Arnal LH (2018) Proactive Sensing of Periodic and
566 Aperiodic Auditory Patterns. *Trends Cogn Sci* 22:870–882.
- 567 Theunissen FE, David SV, Singh NC, Hsu A, Vinje WE, Gallant JL (2001) Estimating
568 spatio-temporal receptive fields of auditory and visual neurons from their responses to
569 natural stimuli. *Netw Comput Neural Syst* 12:289–316.
- 570 Wu MC-K, David SV, Gallant JL (2006) Complete Functional Characterization of Sensory
571 Neurons by System Identification. *Annu Rev Neurosci* 29:477–505.
- 572

573 **Figures**

574

575 **Figure 1:** *A,B* – Example segments of AM BBN and speech stimuli, respectively. *C,D* –
576 Power spectral densities (PSDs) of AM BBN and speech stimuli, respectively. The AM BBN
577 had a broadband frequency distribution by design, while the male speaker had a frequency
578 distribution that was dominated by frequencies below 5000 Hz. *E,F* – Amplitude histograms
579 of AM BBN and speech envelopes, respectively. Both envelopes had quite broadly distributed
580 amplitude distributions. Please note that the amplitude distribution of the AM BBN envelope
581 was uniform by design, but after extracting the envelope from the AM BBN signal using the
582 Hilbert transform, it was less so. Please also note that the amplitude distribution of the
583 speech envelope was more skewed, with a higher percentage of samples in the higher
584 amplitude bins. *G,H* – PSDs of AM BBN and speech envelopes, respectively. Both signals
585 had envelopes with a bottom-heavy (right-skewed) frequency distribution indicating that their
586 modulation rates were dominated by low frequencies.

587

588 **Figure 2:** *A* – Example segments of the envelope, SPL envelope, and onset envelope stimulus
589 representations. *B* – Corresponding segment of the AB envelope.

590

591 **Figure 3:** *A* – Topographic plot displaying prediction accuracies for the envelope model for
592 the AM BBN dataset and highlighting the channels chosen for analysis. *B* – Prediction
593 accuracies for each model and subject, including null hypotheses for the envelope model as
594 determined from the permutation tests, and indications of significance as determined from the
595 *t*-tests.

596

597 **Figure 4:** Analysis of amplitude-dependent changes at a single representative channel over
598 left central scalp for the AM BBN dataset. **A** – Group average AB envelope TRF, plotted to
599 minimize the difference between adjacent traces. **B** – Image plot of group average AB
600 envelope TRF. **C** – NI peak latencies across group average AB envelope TRF bins. **D** – PI-
601 NI peak-to-peak amplitudes across group average AB envelope TRF bins.

602

603 **Figure 5:** **A** – Topographic plot displaying prediction accuracies for the envelope model for
604 the speech dataset and highlighting the channels chosen for analysis. **B** – Prediction
605 accuracies for each model and subject, including null hypotheses for the envelope model as
606 determined from the permutation tests, and indications of significance as determined from the
607 *t*-tests.

608

609 **Figure 6:** Analysis of amplitude-dependent changes at a single representative channels over
610 left central scalp for the speech dataset. **A** – Group average AB envelope TRF, plotted to
611 minimize the difference between adjacent traces. **B** – Image plot of group average AB
612 envelope TRF. **C** – NI peak latencies across group average AB envelope TRF bins. **D** – PI-
613 NI peak-to-peak amplitudes across group average AB envelope TRF bins.

614

615 **Figure 7:** **A** – Prediction accuracies for each model and subject for the AM BBN dataset. **B** –
616 Prediction accuracies for each model and subject for the speech dataset. **C** – Group average
617 AB onset envelope TRF for the AM BBN dataset, plotted to minimize the difference between

618 adjacent traces. **D** – Image plot of group average AB onset envelope TRF for the AM BBN
619 dataset.













