

---

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

## **Auditory Brainstem Responses to Continuous Natural Speech in Human Listeners**

**Ross K. Maddox<sup>1,2,3</sup> and Adrian K. C. Lee<sup>4,5</sup>**

<sup>1</sup>*Department of Biomedical Engineering, University of Rochester, Rochester, NY 14627, USA*

<sup>2</sup>*Department of Neuroscience, University of Rochester, Rochester, NY 14642, USA*

<sup>3</sup>*Del Monte Institute for Neuroscience, University of Rochester, Rochester, NY 14642, USA*

<sup>4</sup>*Institute for Learning & Brain Sciences, University of Washington, Seattle, WA 98195, USA*

<sup>5</sup>*Department of Speech & Hearing Sciences, University of Washington, Seattle, WA 98105, USA*

DOI: 10.1523/ENEURO.0441-17.2018

Received: 19 December 2017

Revised: 22 January 2018

Accepted: 24 January 2018

Published: 31 January 2018

---

**Author contributions:** RKM and AKCL designed research and wrote the paper; RKM performed research and analyzed data.

**Funding:** <http://doi.org/10.13039/100000055HHS> | NIH | National Institute on Deafness and Other Communication Disorders (NIDCD)  
R00DC014288  
R01DC013260

**Conflict of Interest:** Authors report no conflict of interest.

This work was funded by NIH grants R00DC014288 awarded to RKM and R01DC013260 awarded to AKCL.

**Correspondence should be addressed to** To whom correspondence should be addressed Ross K Maddox, University of Rochester, 201 Robert B. Goergen Hall, P.O. Box 270168, Rochester, NY 14627. Email: [ross.maddox@rochester.edu](mailto:ross.maddox@rochester.edu)

**Cite as:** eNeuro 2018; 10.1523/ENEURO.0441-17.2018

**Alerts:** Sign up at [eneuro.org/alerts](http://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 © 2018 Maddox and Lee

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.

## **Auditory brainstem responses to continuous natural speech in human listeners**

### Abbreviated title

Auditory brainstem responses to natural speech

### Authors

Ross K Maddox

Department of Biomedical Engineering, University of Rochester, Rochester, NY 14627, USA  
Department of Neuroscience, University of Rochester, Rochester, NY 14642, USA  
Del Monte Institute for Neuroscience, University of Rochester, Rochester, NY 14642, USA

Adrian KC Lee

Institute for Learning & Brain Sciences, University of Washington, Seattle, WA 98195, USA  
Department of Speech & Hearing Sciences, University of Washington, Seattle, WA 98105, USA

### Author contributions

RKM and AKCL designed research and wrote the paper; RKM performed research and analyzed data.

### To whom correspondence should be addressed

Ross K Maddox

University of Rochester  
201 Robert B. Goergen Hall  
P.O. Box 270168  
Rochester, NY 14627

Email: ross.maddox@rochester.edu

Number of figures: 7

Number of tables: 0

Number of multimedia: 0

Abstract words: 238

Significance statement words: 117

Introduction words: 717

Discussion words: 2,156

### Acknowledgments

A preliminary version of this work was presented at the MidWinter Meeting of the Association for Research in Otolaryngology in February 2017. The authors wish to thank Susan McLaughlin and Tiffany Waddington for assistance with data collection.

### Conflict of interest

Authors report no conflict of interest.

### Funding sources

This work was funded by NIH grants R00DC014288 awarded to RKM and R01DC013260 awarded to AKCL.

## 1 ABSTRACT

2 Speech is an ecologically essential signal whose processing crucially involves the subcortical nuclei of  
3 the auditory brainstem, but there are few experimental options for studying these early responses in  
4 human listeners under natural conditions. While encoding of continuous natural speech has been  
5 successfully probed in the cortex with neurophysiological tools such as electro- and  
6 magnetoencephalography, the rapidity of subcortical response components combined with unfavorable  
7 signal-to-noise ratios has prevented application of those methods to the brainstem. Instead,  
8 experiments have used thousands of repetitions of simple stimuli such as clicks, tone-bursts, or brief  
9 spoken syllables, with deviations from those paradigms leading to ambiguity in the neural origins of  
10 measured responses. In this study we developed and tested a new way to measure the auditory  
11 brainstem response to ongoing, naturally uttered speech, using electroencephalography to record from  
12 human listeners. We found a high degree of morphological similarity between the speech-derived  
13 auditory brainstem responses (ABR) and the standard click-evoked ABR—in particular, a preserved  
14 wave V, the most prominent voltage peak in the standard click-evoked ABR. Because this method  
15 yields distinct peaks that recapitulate the canonical ABR, at latencies too short to originate from the  
16 cortex, the responses measured can be unambiguously determined to be subcortical in origin. The use  
17 of naturally uttered speech to measure the ABR allows the design of engaging behavioral tasks,  
18 facilitating new investigations of the potential effects of cognitive processes like language and attention  
19 on brainstem processing.

20

## 21 SIGNIFICANCE STATEMENT

22 The brainstem is crucial to speech processing, yet a majority of speech studies have focused on the  
23 cortex. This is in large part because practical limitations have made elusive a paradigm for studying  
24 brainstem processing of continuous natural speech in human listeners. Here we adapt methods that  
25 have been employed for studying cortical activity to the auditory brainstem. We measure the response  
26 to continuous natural speech and show that it recapitulates important aspects of the click-evoked  
27 response. The method also allows simultaneous investigation of cortical activity with no added  
28 recording time. This discovery paves the way for studies of speech processing in the human brainstem,  
29 including its interactions with higher order cognitive processes originating in the cortex.

30

## 31 INTRODUCTION

32 When speech enters the ear and is encoded by the cochlea, it goes on to be processed by an  
33 ascending pathway that spans the auditory nerve, brainstem, and thalamus before reaching the cortex.  
34 Far from being relays, these subcortical nuclei perform a dazzling array of important functions, from  
35 sound localization (Grothe and Pecka, 2014) to vowel coding (Carney et al., 2015), making their  
36 function essential to understand. In humans, the primary method for measuring activity in subcortical  
37 nuclei is the auditory brainstem response (ABR): a highly stereotyped scalp potential in the first ~10 ms  
38 following a very brief stimulus such as a click, recorded through electroencephalography (EEG)  
39 (Burkard et al., 2006). The potential comprises components referred to as waves, given Roman  
40 numerals I–VII according to their latency. Individual waves have been tied to activity in specific parts of  
41 the ascending pathway: wave I (~2 ms latency) is driven by auditory nerve activity, wave III (~4 ms) by  
42 the cochlear nucleus, and wave V (~6 ms) principally by the lateral lemniscus (Møller et al., 1995).  
43 However, because the waves are so rapid, and the signal-to-noise ratio (SNR) so low, the ABR must be  
44 measured by presenting thousands of repeated punctate stimuli. Thus, while there are important  
45 neuroscience questions regarding how subcortical nuclei process natural stimuli like speech, or how  
46 they might be affected by cognitive processes through efferent feedback (Terreros and Delano, 2015),  
47 the practical limitations of the ABR paradigm make it primarily a clinical tool.

48 One common method for measuring the brainstem response to speech is the complex ABR (cABR)  
49 (Skoie and Kraus, 2010). The cABR represents the averaged response to repetitions of a short spoken  
50 syllable (e.g., a ~40 ms “da”). The onset response can be analyzed in the time domain, but because the  
51 stimulus is longer than the response, ambiguity about the origin of response components arises for all  
52 but the earliest latencies. The voiced part of the speech elicits a frequency following response (FFR)

53 that can be analyzed in the frequency domain. The FFR has been shown to be primarily driven by the  
54 inferior colliculus of the brainstem, but it results from a mixture of sources including the superior olive  
55 (Smith et al., 1975) and also may include small contributions from the cortex (Coffey et al., 2016).

56 A different method, used for studying cortical activity, treats the auditory evoked potential as the  
57 impulse response of a linear system, which can be mathematically derived from known input and output  
58 signals (Aiken and Picton, 2008; Ding and Simon, 2012a, 2012b; Lalor et al., 2009; Lalor and Foxe,  
59 2010). Continuous natural speech is presented (input) while EEG is recorded (output), and the brain's  
60 response is calculated through linear regression. Rather than raw audio, the regressor (i.e., input) used  
61 is the amplitude envelope, which by construction contains no fast fluctuations, making it too slow for  
62 studying subcortical nuclei. A recent study aimed at the brainstem uses the cross-correlation of the  
63 speech stimulus's fundamental waveform with the EEG recording (Forte et al., 2017). The response is  
64 a single peak with a latency of 9 ms but a relatively broad width. As with the FFR, interpreting this  
65 response is complicated by the likelihood that it is dominated by the inferior colliculus but represents a  
66 mixture of sources.

67 Here we measured auditory brainstem activity in response to natural speech using a new paradigm.  
68 The methods were based on cortical studies, with an important difference: the regressor was the  
69 rectified speech audio, meaning that fine structure was largely preserved. The speech-derived  
70 responses recapitulated important aspects of the click-evoked ABR, most notably in the presence of a  
71 distinct wave V. The speech-derived wave V latency and amplitude were both highly correlated with the  
72 click-evoked response across subjects, demonstrating common neural generators. Thus by preserving  
73 the latencies of individual response components, the speech-derived ABR allows the experimenter to  
74 assess neural activity at separate stages along the auditory pathway.

75 The goal of this study was to develop a method for measuring the ABR to natural speech in  
76 experiments where clicks and other standard stimuli are disadvantageous, inappropriate, or impossible  
77 to use. The results show that it is possible to use natural speech stimuli to study speech processing in  
78 the human brainstem, paving the way for subcortical studies of attention, language, and other cognitive  
79 processes.

80

## 81 **MATERIALS AND METHODS**

### 82 Experimental design and statistical analysis

83 Our goal was to measure the speech-derived ABR in human listeners and validate it against the click-  
84 evoked response. We first recorded click-evoked responses to pseudorandomly timed click trains and  
85 then validated them against the responses evoked by standard, periodic click trains. We then compared  
86 the speech-derived response to the pseudorandom click-evoked response. We validated by the  
87 speech-derived response by comparing its overall morphology and wave V latency and amplitude to  
88 those of the click-evoked response.

89 All subjects' click- and speech-derived responses were plotted individually. To compare the similarity of  
90 two responses from a single subject (e.g., the click-evoked response to the speech-derived response),  
91 Pearson's product-moment correlation was used. The median and interquartile range of each  
92 distribution of correlation coefficients across subjects was reported, in addition to plotting its histogram.  
93 Two distributions of correlation coefficients were compared using Wilcoxon's signed-rank test for non-  
94 normal distributions.

95

### 96 Subjects

97 All experiments were done under a protocol approved by the University of Washington Institutional  
98 Review Board. All subjects gave informed consent prior to participation, and were compensated for  
99 their time. We collected data from 24 subjects (17 females). The mean age was 27.8 years, with a  
100 standard deviation of 6.9 and a range of 19–45. Subjects had normal hearing, defined as audiometric  
101 thresholds of 20 dB HL or better in both ears at octave frequencies ranging from 250 to 8000 Hz. All

102 subjects identified English as their first language except for two, who identified a different language but  
103 had been speaking English daily for over twenty years.

104

#### 105 EEG recording

106 Scalp potentials were recorded with passive Ag/AgCl electrodes, with the positive and negative  
107 electrodes connected to a differential preamplifier (Brainvision LLC, Greenboro, SC). The positive  
108 electrode was at location FCz in the standard 10-20 coordinate system. The negative (reference)  
109 electrode was clipped onto the subject's left earlobe. The ground electrode was placed at Fpz. Data  
110 were high-pass filtered at 0.1 Hz during recording (additional filtering occurred offline).

111 Subjects were seated in a comfortable chair in a sound-treated room (IAC, North Aurora, IL). They were  
112 not asked to attend the stimuli. Instead, they faced a computer monitor showing silent episodes of  
113 "Shaun the Sheep" (Starzak and Sadler, 2007), an animated show that has no talking, making subtitles  
114 unnecessary. They were first presented with 40 epochs of speech stimuli for calculating the speech  
115 ABR, and then were presented with 10 minutes of click stimuli (twenty repetitions of a frozen 30 s  
116 epoch). All stimuli were presented over insert earphones (ER-2, Etymotic Research, Elk Grove, IL)  
117 which were plugged into a stimulus presentation system consisting of a real-time processor and a  
118 headphone amplifier (RP2.1 and HB7, respectively, Tucker Davis Technologies, Alachua, FL). Stimulus  
119 presentation was controlled with a python script using publicly available software (available at  
120 <https://github.com/LABSN/expyfun>).

121

#### 122 Speech stimuli

123 Speech stimuli were taken from two audiobooks. The first was *A Wrinkle in Time* (L'Engle, 2012), read  
124 by a female narrator. The second was *The Alchemyst* (Scott, 2007), read by a male narrator. The  
125 audiobooks were purchased on compact disc and ripped to uncompressed wav files to avoid data  
126 compression artifacts. They were resampled to 24,414 Hz, the native rate of the RP2 presentation  
127 system. They were then processed so that any silent pauses in the speech longer than 0.5 s were  
128 truncated to 0.5 s. Because the ABR is principally driven by higher stimulus frequencies (Abdala and  
129 Folsom, 1995), the speech was gently high-passed with a first-order Butterworth filter with a cutoff of  
130 1,000 Hz and a slope of 6 dB / octave. The speech was still natural sounding and completely  
131 intelligible. This filter also helped to compensate for low-frequency spectral differences between the  
132 male and female narrator around their fundamental frequencies. After that, the speech was normalized  
133 to an average root-mean-square amplitude that matched that of a 1 kHz tone at 75 dB SPL. Figures  
134 1A,D,G show the pressure waveform of the word "Thursday" spoken by the male narrator, the  
135 spectrogram of that word's first syllable, and the power spectral density (PSD) of a 30 s segment of the  
136 female and male speech stimuli. It is evident from Fig. 1D that the filtering did not affect the presence of  
137 pitch information (glottal pulses at the fundamental frequency are easily visible as vertical striations,  
138 even well below 1,000 Hz), and from Fig. 1G that the lowest speech formants were still present (plenty  
139 of energy remaining in 300–500 Hz region).

140 The audiobooks were then sectioned into epochs of 64 s, including a 1 s raised cosine fade-in and  
141 fade-out. The last four seconds of each epoch were repeated as the first four seconds of the next one,  
142 so that subjects could pick up where they left off in the story (if they were listening), meaning that 60 s  
143 of novel speech were presented in each epoch. The stimuli were not new to the subjects—before this  
144 passive listening task, they had completed a session using the same stimuli where they had to answer  
145 questions about the content they had just heard. Data from that task were for a different scientific  
146 question and do not appear here. These minute-long excerpts were presented in sequence, two from  
147 one story and then in alternating sets of four, finishing with two epochs from the second story. Speech  
148 stimuli were presented diotically.

149

#### 150 Click stimuli

151 Click stimuli were aperiodic trains of rarefaction clicks lasting 82  $\mu$ s (representing two samples at the  
152 24,414 Hz sampling rate, which was closest possible to the standard 100  $\mu$ s click duration with our  
153 hardware). Clicks were timed according to a Poisson point process with a rate of 44.1 clicks / s. The  
154 timing of one click had no correlation with the timing of any other click in the train, rendering the  
155 sequence spectrally white in the statistical sense. A pair of 30 s sequences was created and presented  
156 dichotically 20 times to each subject, meaning that 26,460 clicks contributed to each ear's response.  
157 The responses presented herein are the sum of the monaural responses. Clicks were presented at 75  
158 dB peak-to-peak equivalent SPL (i.e., the amplitude of clicks matched the peak-to-peak amplitude of a  
159 1 kHz sinusoid presented at 75 dB SPL).

160 While no previous study has used exactly this type of click timing, several have used various types of  
161 pseudorandom sequences (Burkard et al., 1990; Delgado and Ozdamar, 2004; Holt and Özdamar,  
162 2014; Thornton and Slaven, 1993). Uniformly, these studies find that the ABRs from randomized versus  
163 periodic click trains are highly similar at the same stimulation rates. Random timing has two main  
164 benefits over the much more common periodic timing: 1) the analysis window for the response can be  
165 extended arbitrarily to any beginning and end point without fear of temporal wrapping, and 2) no high-  
166 pass filtering is necessary to remove the strong frequency component at the (periodic) presentation  
167 rate, because it does not exist. A third benefit, specific to this study, is that the same linear systems  
168 analysis could be done to compute the speech-derived and the click-evoked ABR, yielding a more  
169 direct comparison between the two. Figures 1B,E,H show part of a Poisson click train in the same  
170 manner that Figs. 1A,D,G do for speech.

171 To be sure that the click paradigm we used yielded results matching standard ABRs evoked with  
172 periodic click trains, we also collected ABRs using periodic click trains of the same rate of 44.1 clicks /  
173 s, presented diotically. Periodic trains were also presented in twenty epochs of 30 s, yielding the same  
174 total sweep count of 26,460. The periodic click train stimulus is shown in Figs. 1C,F,I.

175

#### 176 Data analysis

177 Responses to both speech and click train stimuli were found through deconvolution, in a manner  
178 broadly similar to previous papers focused on cortical activity (Lalor et al., 2009; Lalor and Foxe, 2010).  
179 The essence of deconvolution is determining the impulse response of a linear time-invariant system  
180 given a known input (here, the processed continuous speech signal) and a known output (here, the  
181 recorded scalp potential). The methods in this study vary from previous ones in the recording  
182 parameters and preprocessing steps, but otherwise utilize essentially the same mathematical  
183 principles.

184

#### 185 *Speech stimuli preprocessing*

186 Before we could derive the speech response, we needed to calculate the regressor from the audio  
187 data. The auditory brain is mostly agnostic to the sign of an acoustic input, as evidenced by the high  
188 degree of similarity between evoked responses to compression versus rarefaction clicks (Møller et al.,  
189 1995). For this reason, some sort of rectifying nonlinearity applied to the input speech is needed as a  
190 preprocessing step. We used half-wave rectification. Specifically, we performed all analyses twice—  
191 once keeping the positive peaks, and then a second time keeping the inverted negative peaks—and  
192 then averaged the resulting responses, in a process akin to the compound peristimulus time histogram  
193 used by Pfeiffer and Kim (1972). This alone significantly reduced, but did not eliminate, stimulus  
194 artifacts, similar to the common technique of alternating polarity in the click-evoked ABR (Hall III, 2006).  
195 Further artifact reduction steps are described later in this section. Following rectification, the data were  
196 downsampled from 24,414 Hz to the EEG recording rate of 10,000 Hz.

197

#### 198 *Click train preprocessing*

199 Owing to its extreme sparsity, downsampling a click train using standard methods would result in  
200 significant signal processing artifact, viz., Gibbs ringing. We instead used the list of click times from the

201 original click train (24,414 Hz sampling rate) and created a click train at 10,000 Hz sampling rate by  
 202 placing unit-height single-sample impulses at the closest integer indices corresponding the original click  
 203 times.

204 When the input to a system has a white power spectrum, the system's impulse response can be  
 205 determined as the cross-correlation of the input and output. For a click train, which is essentially a  
 206 series of unit-height single-sample impulses, the deconvolved impulse response becomes equivalent to  
 207 the click-triggered average, which is how ABRs are usually calculated. This results in a convenient  
 208 parity between the typical averaging methods used for ABR and the deconvolution used here. In other  
 209 words: rather than using a completely new mode of analysis for ABR (deconvolution), we have instead  
 210 generalized the methods already in use to be appropriate for arbitrary stimuli, beyond click trains.

211

#### 212 *EEG preprocessing*

213 EEG data were first high-pass filtered at 1 Hz (first-order Butterworth), and then notch filtered at 60,  
 214 180, and 300 Hz with 5 Hz wide second-order infinite impulse response notch filters, designed with the  
 215 *iirnotch* function of the SciPy python package (RRID:SCR\_008058). Because of the continuous nature  
 216 of the stimuli, no epoch rejection was done. Instead, any time the EEG potential crossed  $\pm 100 \mu\text{V}$ , a 1 s  
 217 segment of the response was zeroed, centered around the offending sample, removing it from the  
 218 calculation. 100  $\mu\text{V}$  is a larger rejection threshold than most EEG studies use, which was necessary  
 219 because the EEG data had higher power due to the minimal filtering that was applied (high-pass at 1  
 220 Hz). Zeroing portions of an epoch slightly reduces its energy. So that the amplitude of the calculated  
 221 response was not affected, the EEG data for each epoch was multiplied by a corrective gain factor  $g_r$ :

$$222 \quad g_r = N / (N - N_r),$$

223 where  $N$  is the total number of samples in the epoch and  $N_r$  is the number of rejected samples. After  
 224 filtering and resampling, the data were segmented into epochs that started with the stimulus onset and  
 225 ended 100 ms after the stimulus (epochs were thus 64.1 s long for speech stimuli and 30.1 s long for  
 226 clicks). With these parameters, a median of 1.1% (0.3%–2.3% interquartile range) of data was rejected  
 227 from each subject's EEG recordings.

228

#### 229 *Stimulus artifact removal*

230 EEG recordings from some subjects showed stimulus artifacts, resulting from electromagnetic  
 231 "leakage" of the headphone driver to the EEG system. We developed a protocol for removing these  
 232 artifacts which involved estimating the artifact and then subtracting it from the EEG recording. For each  
 233 epoch, we first computed the discrete Fourier transform (DFT) of the stimulus (NB: the raw stimulus  
 234 audio, not the rectified stimuli) and the EEG recording. We then divided the EEG DFT by the stimulus  
 235 DFT and, computed the inverse DFT of the quotient. We then cropped that signal (which is effectively  
 236 the estimated impulse response that describes the electromagnetic leakage as a system) to the lags in  
 237 a 10 ms window centered around the artifact at  $-0.9$  ms. Because this 10 ms impulse response was  
 238 empirically estimated from the data, and was elicited by a stimulus with very little energy below  $\sim 100$   
 239 Hz, it contained a high level of low-frequency noise. This noise was removed by fitting a sixth-order  
 240 polynomial to the estimated impulse response and subtracting that fit from the signal itself. The  
 241 polynomial order was chosen as the lowest that removed artifacts clearly (by visual inspection)  
 242 unrelated to the impulse response. Because the noise was so much larger than the true impulse  
 243 response (and was largely made of frequencies absent from the speech signal), the result of this  
 244 subtraction was a clear impulse response with little low-frequency noise remaining. This process was  
 245 completed for every epoch individually, the average impulse response for each subject was computed  
 246 across all epochs, and that average was multiplied by a 10 ms Hann window. Finally, for each epoch,  
 247 the stimulus was convolved with the computed artifact impulse response, and the resulting signal was  
 248 subtracted from the EEG recording.

249 The stimulus artifact was only removed from the speech-derived responses. These steps were not  
 250 necessary for click-evoked responses because those artifacts manifest as single sharp spikes before 0  
 251 latency.

252 It should be noted that there exist simpler ways to eliminate or mitigate stimulus artifacts. The simplest  
 253 is electromagnetic shielding around the headphone drivers. Alternating the polarity of the speech  
 254 stimulus should also significantly reduce stimulus artifacts in future experiments. This could be done at  
 255 the level of the 64 s epochs, or it could be done at the word or phrase level, as long as the phase  
 256 inversions were hidden by silent gaps in the speech. However, these methods must be implemented at  
 257 the time of the recording, and were not here, which is why the signal processing steps described above  
 258 were used.

259

#### 260 *Response calculation*

261 We used linear least-squares regression to calculate the responses, as in previous work (Lalor et al.,  
 262 2009). The response was considered to be the weights over a range of time lags that best  
 263 approximated the EEG output as the weighted sum of the input stimulus regressor over those lags. For  
 264 the sake of computational and memory efficiency, the stimulus autocorrelation matrix and stimulus-  
 265 response cross-correlation were both calculated via their Fourier counterparts using frequency-domain  
 266 multiplication. These specific methods have been incorporated into the mne-python package (Gramfort  
 267 et al., 2013) (RRID:SCR\_005972). The stimulus regressors were sufficiently broadband such that no  
 268 regularization was necessary, so none was used (had there been near-zeros in their amplitude spectra  
 269 this would not have been the case). The response weights were calculated over the range of lags  
 270 spanning  $-150$  to  $350$  ms. After the response was calculated, it was low-pass filtered at  $2,000$  Hz (first-  
 271 order Butterworth), and then baseline corrected by subtracting the mean potential between  $-10$  and  $0$   
 272 ms from the whole response (note that this baseline step affects the plots but none of the measures  
 273 used here, and can thus be considered optional). For the speech stimuli, the response to each narrator  
 274 was calculated separately, and then averaged to calculate each subject's speech-derived response.  
 275 The stimulus regressors were sufficiently broadband such that no regularization was necessary, so  
 276 none was used (had there been near-zeros in their amplitude spectra this would not have been the  
 277 case).

278

#### 279 *Speech-derived response amplitude normalization*

280 Auditory onsets elicit much larger responses than ongoing stimulus energy due to (Thornton and  
 281 Slaven, 1993). However, this non-linear adaptation is not accounted for by the linear regression. For  
 282 that reason, the raw speech-derived responses, for which the majority of the stimulus energy can be  
 283 considered "ongoing," were much smaller than the click-evoked responses, whose stimuli are  
 284 essentially a series of onsets. To correct for this, we computed a single empirical subject-specific  
 285 normalization factor,  $g_n$ , that put the speech-derived responses in a similar amplitude range as the  
 286 click-evoked ones:

287

$$g_n = E_i(\sigma_{c,i}) / E_i(\sigma_{s,i}),$$

288 where  $\sigma_{c,i}$  is the standard deviation of subject  $i$ 's click-evoked response in the range of  $0-20$  ms,  $\sigma_{s,i}$  is  
 289 the same for the speech-derived response, and  $E_i$  represents the mean over subjects. All speech-  
 290 derived responses shown in microvolts have been multiplied by  $g_n$ . In our study  $g_n$  had a value of  $28.2$ ,  
 291 but it must be stressed that this value depends on the unitless scale chosen for storing the digital audio  
 292 (ours had a root-mean-square amplitude of  $0.01$ ), and is thus not suitable for use in other studies. For  
 293 this reason no direct amplitude comparisons were made between click- and speech-derived responses.  
 294 Instead, we computed correlations (which do not depend on scaling factors) of their morphologies  
 295 within subjects, as well as their wave V latencies and amplitudes across subjects.

296

#### 297 Standard ABR measurement

298 The ABR to the periodic click trains was calculated through traditional averaging rather than regression.  
 299 The raw data were notch filtered to remove line noise and low-pass filtered at  $2,000$  Hz as described  
 300 above. However, the high-pass filter was different: a causal second order Butterworth filter with a cutoff  
 301 of  $150$  Hz was used to be consistent with standard practice and to generate a canonical waveform

302 (Burkard et al., 2006; Hall III, 2006). The response to each click presentation was then epoched from  
303  $-3$  ms to 19.7 ms, which was the longest window allowed by the periodic click rate of 44.1 clicks / s  
304 before temporal wrapping occurred. Filtered epochs were rejected if the peak-to-peak amplitude  
305 exceeded 100  $\mu$ V.

306

## 307 RESULTS

### 308 Poisson click trains yield canonical ABRs

309 Responses to Poisson click trains were used as the benchmark to which the speech-derived responses  
310 were compared. Even though similar types of pseudorandom stimuli have been used in the past, it was  
311 important to confirm that these specific stimuli used here provided canonical ABR waveforms. The  
312 grand average periodic and Poisson click trains are shown overlaid in Fig. 2A (both shown high-pass  
313 filtered at 150 Hz). To quantify their similarity, we computed Pearson's correlation coefficient between  
314 the two waveforms for each subject between lags of 0 and 19.7 ms. The median correlation was 0.89  
315 (interquartile range 0.82–0.92), indicating a very high degree of similarity. The histogram of correlations  
316 is shown in Fig. 2B.

317 Figure 2C shows the average Poisson click-evoked response under two filtering conditions: 1) high-  
318 pass filtered at 150 Hz as in Fig. 2A, and 2) broadband (high-passed at 1 Hz as described in the EEG  
319 pre-processing methods section above). The latter will be used henceforth as the click-evoked ABR to  
320 which the speech-derived ABR is compared. It is thus important to note that even though these  
321 responses seem to have morphological differences from the "standard" ABR, that is simply because  
322 using pseudorandom click timing obviates the need for high-pass filtering, and that filtering was  
323 bypassed in the interest of comparing the whole responses. The wideband responses we obtained here  
324 using Poisson click trains were highly similar in shape, amplitude, and latency to previous wideband (5  
325 Hz high-pass) ABRs obtained using low rate (11 Hz) periodic clicks (Gu et al., 2012), and were much  
326 more efficient to obtain.

327

### 328 Early speech-derived responses exhibit brainstem response characteristics

329 Broadly speaking, there were strong similarities between the early ( $< 20$  ms) click-evoked and speech-  
330 derived responses (Fig. 3A). In this latency range, responses are likely to progress from brainstem to  
331 thalamus and primary auditory cortex as latency increases. We will first make whole-waveform  
332 comparisons, and then consider specific canonical ABR components.

333 To compare the overall waveforms, we computed Pearson's correlation coefficient of the speech- and  
334 click-evoked waveforms for each subject in the range of 0–20 ms (Fig. 3B). The median correlation  
335 coefficient was 0.82 (interquartile range 0.77–0.89). Figure 3C shows each subject's click- and speech-  
336 derived response, in descending correlation order. In our speech-derived responses, waves I–IV were  
337 "smeared" together. However, we found a clear wave V in individual subjects' responses as well as the  
338 grand average. Wave VI was also visible in the grand average, but was less consistent at the  
339 individual-subject level.

340 We identified wave V by low-pass filtering at 1,000 Hz with a zero-phase filter and finding the peak of  
341 the waveform in the 5–7 ms range. For the click-evoked responses, wave V was present for all  
342 subjects, with a latency of  $6.50 \pm 0.25$  ms (mean  $\pm$  standard deviation). For speech-derived responses,  
343 wave V was present for all subjects, with a latency of  $6.17 \pm 0.31$  ms. As shown in Fig. 4A, the click-  
344 evoked and speech-derived wave V latencies were highly correlated across subjects ( $r = 0.78$ ,  $p =$   
345  $7 \times 10^{-6}$ , Pearson's product-moment). The peak amplitudes of the speech-derived and click-evoked  
346 wave V were also correlated ( $r = 0.62$ ,  $p = 0.0011$ ; Fig. 4B). These correlations strongly suggest that  
347 the click-evoked and speech-derived ABR have common neural generators.

348

### 349 Speech responses across talkers are similar but not identical

350 One important question is whether the speech-derived response maintains its morphology independent  
 351 of the specific input stimulus, or if it depends on the specific narrator. To investigate this, we split the  
 352 responses to male- and female-narrated trials and compared them to determine the role that the  
 353 difference in the narrators' input spectra might play. The grand average waveforms for the two narrators  
 354 are of the same magnitude and overall shape, despite the differing spectra of their input stimuli (Fig.  
 355 5A). The median female-male correlation coefficient was 0.81 (interquartile range 0.68–0.90; Fig. 5B).

356 While perfect overlap would be indicated by correlation coefficients of 1.0, splitting the data in half (viz.,  
 357 into male- and female-narrated epochs) adds noise to each of the responses. To put the male-female  
 358 correlation coefficients in context, we can split the data a different way and compare. We split the data  
 359 into halves consisting of the even versus odd trials, which contained the same number of male and  
 360 female epochs (i.e., each split contained 10 male and 10 female trials, distributed evenly in time across  
 361 the recording session). We then compared those waveforms as above. The median correlation  
 362 coefficient between splits was 0.89 (interquartile range 0.78–0.92). We compared the male-female split  
 363 coefficients to these arbitrarily split coefficients, and found a significant difference ( $T(23) = 68$ ,  $p =$   
 364  $0.019$ , Wilcoxon signed-rank test). This indicates that while the responses to female and male-uttered  
 365 speech are indeed similar, there is still some dependence on the stimulus.

366

#### 367 Sufficient SNR was attained for all subjects

368 A measure's usefulness decreases with the amount of time required to obtain it. SNR generally  
 369 increases with additional data, and is thus a function of recording time. To assess SNR, we first  
 370 computed the cumulative average response up to each of the 40 recording epochs, for each subject.  
 371 We then computed the SNR, in decibels, as

$$372 \quad \text{SNR} = 10 \log_{10}[(\sigma_{\text{ABR}}^2 - \sigma_{\text{noise}}^2) / \sigma_{\text{noise}}^2],$$

373 where  $\sigma_{\text{ABR}}^2$  and  $\sigma_{\text{noise}}^2$  are the variances of the response in the lag intervals from 0 to 20 ms and –125  
 374 to –10 ms, respectively. The results are plotted in Fig. 6.

375 Experimental demands differ, but an SNR of 0 dB or better typically allows a response to be easily seen  
 376 and inspected. For these 24 subjects, 50% achieved that threshold after 9 epochs, 75% after 17  
 377 epochs, 95% by 27 epochs, and all by the end of the thirty-third epoch. These data are shown in Figure  
 378 6. Taken as a whole, they confirm that the speech-derived ABR can be measured to a useful SNR in  
 379 reasonable durations. Recording times should also be short enough that multiple conditions can be  
 380 tested in a single experimental session. The median SNR's evolution over time also generally aligned  
 381 with the theoretical expectation of +3 dB per doubling of recording time: after 10, 20, and 40 epochs it  
 382 was 1.1 dB, 4.4 dB, and 7.8 dB respectively.

383 We also sought to address the stability of the speech-derived responses over time. To do so, we split  
 384 the data into halves—epochs 1 through 20, and 21 through 40—and compared the responses (female-  
 385 and male-narrated trials distributed evenly throughout each of these halves). Across subjects, the  
 386 median correlation was 0.86 (interquartile range 0.81–0.92). These high correlations suggest  
 387 responses were stable over the session. They are very similar to the even-odd trial split correlations,  
 388 which would be less affected by a response that changes or drifts over the recording session, which  
 389 had a median of 0.89 (0.78–0.92).

390

## 391 **DISCUSSION**

### 392 Early speech responses are interpretable as ABRs

393 The major goal of this work was to study the response of the human auditory brainstem to naturally  
 394 spoken, continuous speech. We computed the speech-derived responses using regression and  
 395 validated them against click-evoked responses. Comparison of the speech-derived and click-evoked  
 396 ABR revealed a high degree of morphological similarity between waveforms, similar overall wave V  
 397 latencies, and a strong correlation between speech-derived and click-evoked wave V latency and  
 398 amplitude across subjects. Taken together, these results show that the speech-derived ABR developed  
 399 here is just that—the response of the auditory brainstem to naturally uttered speech. Note, however,

400 that the goal of this study was not to replace the click-evoked ABR—it was to allow the ABR to be  
401 measured in response to natural speech stimuli presented in the context of an engaging behavioral  
402 task.

403 Incoming acoustic information travels up the auditory pathway in an initial feedforward sweep, from  
404 brainstem to thalamus to cortex. Because the response calculated here is broadband, distinct  
405 components over the range of latencies were preserved. We can thus “localize through latency” and  
406 logically conclude that the peak in the response at ~6 ms has subcortical origins, because it is too soon  
407 after the stimulus to be cortical, where the earliest estimated latencies are 11–14 ms (Wassenhove and  
408 Schroeder, 2012). This eschews the problem of source mixing when attempting to determine brainstem  
409 activity through spatial means, such as beamforming and dipole fits. However, as discussed below, our  
410 method does not preclude those analyses—rather it complements them and facilitates their use,  
411 particularly at longer latencies where sources have cortical origins more appropriate for spatial filtering.

412

#### 413 Speech-derived ABR facilitates new studies of brainstem processing under natural conditions

414 The motivation for this work was to allow the study of auditory brainstem activity in human listeners  
415 without stimulus limitations. To demonstrate the technique’s utility, in this section we propose two  
416 important questions whose corresponding experiments are specifically facilitated by the new methods.

417 Natural speech has rich spectrotemporal structure. It also has semantic and connotative contents that  
418 transcend its acoustics, which cannot be ascertained by a listener unfamiliar with the language. But  
419 until now, studies of brainstem speech processing have been limited to simple repeated stimuli, and  
420 studies using natural speech have been mostly limited to the later potentials corresponding to the  
421 cortex. One exception to this is a recent study by Forte et al. (2017), whose central question is the first  
422 one discussed below.

423

424 *Example experiment 1: does selective attention to one of two competing speech streams affect the*  
425 *brainstem’s response to those streams?*

426 A fundamental question in auditory neuroscience is how the brain selects one sound of interest from a  
427 mixture of several (first described by Cherry (1953) as the Cocktail Party Problem). A number of studies  
428 have shown that in the cortex, the representation of a natural speech stream attended by subjects is  
429 greater than that of one to be ignored (Mesgarani and Chang, 2012; O’Sullivan et al., 2014). Many  
430 brainstem nuclei receive efferent projections from the cortex, or from higher nuclei of the brainstem,  
431 which suggests the possibility of cortical modulation of subcortical processing playing a role in selective  
432 attention (Terrerros and Delano, 2015).

433 The present technique can be applied to studies using the same paradigms as previous cortically-  
434 focused ones: present two speech streams, ask the listener to attend to one, and calculate the speech-  
435 derived response from each stream. An effect of attention in the brainstem would manifest as a  
436 stronger response to the attended stream.

437 The strength of such a paradigm is that it is much more likely to engage attentional mechanisms due to  
438 its ecological validity—whereas asking subjects to attend to one stream of clicks but not the other may  
439 not. The present technique allows specific measures of brainstem processing (e.g., wave V amplitude  
440 and latency) to be reported as a function of attention. One recent study using natural speech has  
441 suggested a subcortical effect of attention (Forte et al., 2017). These results are promising, but  
442 because the responses are akin to a continuous FFR and the waveforms do not recapitulate the  
443 canonical ABR, the specific origins of the effect are harder to verify.

444

445 *Example experiment 2: does understanding speech affect its representation in the brainstem?*

446 The brainstem, in particular the inferior colliculus, is very important to speech processing (Carney et al.,  
447 2015). It is not known if the processing in the brainstem is purely acoustical, or if it is affected by higher  
448 order (e.g., semantic) processing of the speech. With the speech-derived ABR it is possible to test that.

449 A basic design to do so would involve two sets of subjects that spoke and understood only one of two  
450 separate languages (*A* and *B*). Speech stimuli from both languages would be presented to both groups,  
451 and the speech-derived ABR measured.

452 If, e.g., language *A* speakers showed a wave *V* amplitude that was bigger for language *A* than  
453 language *B*, and language *B* speakers showed the opposite, then an effect of speech comprehension  
454 on brainstem processing would be indicated. An additional advantage of the technique is that if the  
455 speech-derived ABR did not show an effect, then the analysis window could be extended to show  
456 cortical processing as well (see next section), so that the first level at which comprehension does affect  
457 the response could still be investigated.

458 It is not possible to test effects of comprehension with click stimuli, because clicks cannot be  
459 understood. This example makes manifest the distinction between the click-evoked ABR and the  
460 speech-derived ABR developed here. While the present study has been aimed at showing the two  
461 responses' neural origins are the same, their experimental applications are quite different.

462

#### 463 Subcortical and cortical responses are available simultaneously

464 While the focus of this work is on the brainstem and midbrain responses, these methods can be used to  
465 measure both subcortical and cortical activity. Simultaneous subcortical and cortical measurements are  
466 possible with the cABR (Skoie and Kraus, 2010), but the differing parameters for number of trials and  
467 inter-stimulus interval needed mean that recording paradigms can be very long. Work aimed at optimal  
468 parameters for simultaneous subcortical-cortical recordings has been successful (Bidelman, 2015), but  
469 still necessarily results in compromises. The present methods allow simultaneous measurement with no  
470 additional recording time and no limitations on the response window due to inter-stimulus interval.

471 This flexibility is illustrated in Fig. 7, where the same average speech-derived response measured here  
472 is plotted three different ways. Figure 7A shows the speech-derived ABR, Fig. 7B extends the window  
473 and employs a low-pass filter appropriate for viewing the middle latency response, and Fig. 7C extends  
474 the time window further and lowers the low-pass frequency to accentuate late auditory evoked  
475 potentials of cortical origin. It is interesting to note in Fig. 7C that, while cortical response amplitudes  
476 are generally thought of as being larger than the ABR, this is not the case when using continuous  
477 speech as a stimulus. This likely stems from the fact that there is significant cortical adaptation to a  
478 continuous stimulus, where typical event-related potential designs are careful to allow enough time  
479 between stimulus onsets to prevent adaptation. While the later peaks in Fig. 7C are surely cortical in  
480 origin, their specific latencies do not perfectly match the canonical latencies of N1 and P2. It is not  
481 entirely clear why this would be the case.

482 While only one EEG channel was used here, there is no reason a full electrode montage could not be  
483 used, assuming one is available (along with considerable hard drive space). This would allow the  
484 simultaneous study of brainstem and cortical processing under natural conditions. Additionally,  
485 interactions between the two are also possible to study by adding interaction terms to the linear model.  
486 For example, a significant interaction between time-varying parietal alpha power and the size of the  
487 ABR could indicate a functional relationship between those areas.

488

#### 489 Filtering must be done carefully

490 It is common practice in EEG experiments to use zero-phase filters whose impulse responses are non-  
491 causal and symmetric about zero lag. This is done to preserve the latencies of the peaks and is  
492 appropriate in many cases. However, the strength of the present approach lies in using the latency of  
493 the response peaks to confirm their subcortical origin. If a non-causal filter is used to filter the EEG  
494 data, then it is possible that a peak at a latency corresponding to cortical activity (e.g., 25 ms) could  
495 affect the response waveform at brainstem latencies (e.g., 6 ms). This could have the result of  
496 erroneous findings that attribute cortical phenomena to subcortical nuclei. Thus, the following two  
497 guidelines are recommended for experiments specifically aimed at the auditory brainstem. First, EEG  
498 data should be filtered with causal filters. Second, when calculating regressors, any filtering that is done  
499 to the input stimulus should be anti-causal (i.e., with an impulse response that has non-zero values only

500 at negative lags). The latter can be practically accomplished by reversing the signal in time, filtering it  
501 with a standard causal filter, and then reversing that result. Using causal filters will inevitably affect the  
502 latencies of peaks, but this can be mitigated by filtering sparingly (i.e., as broadband as the specific  
503 analyses will allow) with low-order filters, as was done here.

504

#### 505 Responses to arbitrary stimuli can be measured

506 For a spectrally rich but non-white stimulus like speech, an important step in deconvolution is whitening  
507 the input stimulus. For a linear system, two broadband stimuli with different spectra should yield the  
508 same impulse response. However, there is no such guarantee for a non-linear system like the auditory  
509 system.

510 The present study suggests that it would be possible to use a range of stimuli to evoke responses with  
511 similar morphologies. First, we consider the main comparison: speech-derived to click-evoked ABR.  
512 Natural speech is different by almost any metric from Poisson click trains, and yet the responses that  
513 we find through regression are very similar (Fig. 3A,B). Second, we consider the responses to female  
514 versus male speech. Males typically speak at a fundamental frequency about half that of females, due  
515 to relatively larger vocal folds. Such a difference, when estimating the response of a highly non-linear  
516 system using linear methods, could have resulted in major differences in the response waveforms, but  
517 this was not the case (Fig. 5A,B). Taken together, it is reasonable to expect that the present technique  
518 could be applied to other real-world non-speech stimuli such as music or environmental sounds, as well  
519 any spectrally rich synthetic stimulus of interest in the lab.

520 Despite the similarity between responses to different stimuli, the differences (e.g. between the female  
521 and male speech-derived responses) do represent a caveat. In future studies, experimenters must be  
522 careful in making comparisons between responses across conditions that did not use identical stimuli.  
523 We suggest that these methods will be most useful in cases where the acoustic stimuli can be  
524 counterbalanced across conditions. While this is good practice in most studies, it is especially important  
525 here for drawing strong conclusions.

526

#### 527 Other regressors may offer improvements

528 An important difference between this study and those that came before it is choice of the regressor.  
529 Because the auditory system is fundamentally nonlinear (viz., it responds with the same sign to both  
530 compression (positive) and rarefaction (negative) clicks), some sort of manipulation of the audio into an  
531 all-positive signal is needed. Previous studies have used the amplitude envelope (Aiken and Picton,  
532 2008; Lalor and Foxe, 2010), spectrotemporal representations (Ding and Simon, 2009), and even  
533 dynamic higher-order features of speech (Di Liberto and Lalor, 2017).

534 Critically, the rectified speech audio used here is a broadband signal, which is what allows distinct ABR  
535 components at short latencies to be resolved in the derived response. There are many other  
536 transformations one could do, which will have important effects on the response waveform obtained.  
537 We piloted several (for example, “raising” the audio to be all-positive by adding it to its Hilbert amplitude  
538 envelope), but decided on the half-wave rectified audio due to its simplicity and the robustness of the  
539 responses it yielded. It is possible—likely, even—that there are better transformations. One  
540 shortcoming of our approach is that no distinct wave I was found, and all of waves I–V were smeared  
541 together. An improvement in the regressor is the most likely route to addressing this, if it is indeed  
542 addressable, and will be a focus of future work.

543

#### 544 Conclusions and future directions

545 Here we present and validate a method for determining the response of the auditory brainstem to  
546 continuous, naturally uttered, non-repeated speech. Speech processing involves a complex network  
547 that ranges from the earliest parts of the auditory pathway to auditory and association cortices. The  
548 technique described here facilitates new neuroscience experiments by making it possible to measure  
549 activity across the auditory neuraxis while human subjects perform natural and engaging tasks. These

550 paradigms will allow study of the subcortical effects of language learning and understanding, attention,  
551 multisensory integration, and many other cognitive processes.

552

## 553 REFERENCES

- 554 Abdala C, Folsom RC (1995) The development of frequency resolution in humans as revealed by the  
555 auditory brain-stem response recorded with notched-noise masking. *J Acoust Soc Am* 98:921–  
556 930.
- 557 Aiken SJ, Picton TW (2008) Human cortical responses to the speech envelope. *Ear Hear* 29:139–157.
- 558 Bidelman GM (2015) Towards an optimal paradigm for simultaneously recording cortical and brainstem  
559 auditory evoked potentials. *J Neurosci Methods* 241:94–100.
- 560 Burkard R, Shi Y, Hecox KE (1990) A comparison of maximum length and Legendre sequences for the  
561 derivation of brain-stem auditory-evoked responses at rapid rates of stimulation. *J Acoust Soc*  
562 *Am* 87:1656–1664.
- 563 Burkard RF, Don M, Eggermont JJ (2006) *Auditory Evoked Potentials: Basic Principles and Clinical*  
564 *Application*, 1st ed. Philadelphia: Lippincott Williams & Williams.
- 565 Carney LH, Li T, McDonough JM (2015) Speech Coding in the Brain: Representation of Vowel  
566 Formants by Midbrain Neurons Tuned to Sound Fluctuations. *eneuro* ENEURO.0004-15.2015.
- 567 Cherry EC (1953) Some experiments on the recognition of speech, with one and with two ears. *J*  
568 *Acoust Soc Am* 25:975–979.
- 569 Coffey EBJ, Herholz SC, Chepesiuk AMP, Baillet S, Zatorre RJ (2016) Cortical contributions to the  
570 auditory frequency-following response revealed by MEG. *Nat Commun* 7:11070.
- 571 Delgado RE, Ozdamar O (2004) Deconvolution of evoked responses obtained at high stimulus rates. *J*  
572 *Acoust Soc Am* 115:1242–1251.
- 573 Di Liberto GM, Lalor EC (2017) Indexing cortical entrainment to natural speech at the phonemic level:  
574 Methodological considerations for applied research. *Hear Res* 348:70–77.
- 575 Ding N, Simon JZ (2012a) Neural coding of continuous speech in auditory cortex during monaural and  
576 dichotic listening. *J Neurophysiol* 107:78–89.
- 577 Ding N, Simon JZ (2012b) Emergence of neural encoding of auditory objects while listening to  
578 competing speakers. *Proc Natl Acad Sci* 109:11854–11859.
- 579 Ding N, Simon JZ (2009) Neural Representations of Complex Temporal Modulations in the Human  
580 Auditory Cortex. *J Neurophysiol* 102:2731–2743.
- 581 Forte AE, Etard O, Reichenbach T (2017) The human auditory brainstem response to running speech  
582 reveals a subcortical mechanism for selective attention. *eLife* 6:e27203.
- 583 Gramfort A, Luessi M, Larson E, Engemann DA, Strohmeier D, Brodbeck C, Goj R, Jas M, Brooks T,  
584 Parkkonen L, Hämäläinen M (2013) MEG and EEG data analysis with MNE-Python. *Front*  
585 *Neurosci* 7.
- 586 Grothe B, Pecka M (2014) The natural history of sound localization in mammals – a story of neuronal  
587 inhibition. *Front Neural Circuits* 8:116.
- 588 Gu JW, Herrmann BS, Levine RA, Melcher JR (2012) Brainstem auditory evoked potentials suggest a  
589 role for the ventral cochlear nucleus in tinnitus. *J Assoc Res Otolaryngol JARO* 13:819–833.
- 590 Hall III JW (2006) *New Handbook for Auditory Evoked Responses*, 1st ed. Boston: Pearson.
- 591 Holt FD, Özdamar Ö (2014) Simultaneous acquisition of high-rate early, middle, and late auditory  
592 evoked potentials In: 2014 36th Annual International Conference of the IEEE Engineering in  
593 Medicine and Biology Society , Presented at the 2014 36th Annual International Conference of  
594 the IEEE Engineering in Medicine and Biology Society pp1481–1484.
- 595 Lalor EC, Foxe JJ (2010) Neural responses to uninterrupted natural speech can be extracted with  
596 precise temporal resolution. *Eur J Neurosci* 31:189–193.
- 597 Lalor EC, Power AJ, Reilly RB, Foxe JJ (2009) Resolving Precise Temporal Processing Properties of  
598 the Auditory System Using Continuous Stimuli. *J Neurophysiol* 102:349–359.
- 599 L'Engle M (2012) *A Wrinkle in Time*. Listening Library.
- 600 Mesgarani N, Chang EF (2012) Selective cortical representation of attended speaker in multi-talker  
601 speech perception. *Nature* 485:233–236.

- 602 Møller AR, Jho HD, Yokota M, Jannetta PJ (1995) Contribution from crossed and uncrossed brainstem  
 603 structures to the brainstem auditory evoked potentials: A study in humans. *The Laryngoscope*  
 604 105:596–605.
- 605 O’Sullivan JA, Power AJ, Mesgarani N, Rajaram S, Foxe JJ, Shinn-Cunningham BG, Slaney M,  
 606 Shamma SA, Lalor EC (2014) Attentional Selection in a Cocktail Party Environment Can Be  
 607 Decoded from Single-Trial EEG. *Cereb Cortex* bht355.
- 608 Pfeiffer RR, Kim DO (1972) Response Patterns of Single Cochlear Nerve Fibers to Click Stimuli:  
 609 Descriptions for Cat. *J Acoust Soc Am* 52:1669–1677.
- 610 Scott M (2007) *The Alchemyst: The Secrets of the Immortal Nicholas Flamel, Book 1*. Listening Library.
- 611 Skoe E, Kraus N (2010) Auditory brainstem response to complex sounds: a tutorial. *Ear Hear* 31:302–  
 612 324.
- 613 Smith JC, Marsh JT, Brown WS (1975) Far-field recorded frequency-following responses: Evidence for  
 614 the locus of brainstem sources. *Electroencephalogr Clin Neurophysiol* 39:465–472.
- 615 Starzak R, Sadler C (2007) *Shaun the Sheep (Season 1)*. Aardman Animations.
- 616 Terreros G, Delano PH (2015) Corticofugal modulation of peripheral auditory responses. *Front Syst*  
 617 *Neurosci* 9.
- 618 Thornton ARD, Slaven A (1993) Auditory brainstem responses recorded at fast stimulation rates using  
 619 maximum length sequences. *Br J Audiol* 27:205–210.
- 620 Wassenhove V van, Schroeder CE (2012) Multisensory Role of Human Auditory Cortex In: *The Human*  
 621 *Auditory Cortex*, Springer Handbook of Auditory Research , pp295–331. Springer, New York,  
 622 NY.
- 623

624

#### 625 **FIGURE LEGENDS**

626 Figure 1. Acoustic stimuli. (A,B,C) Pressure waveforms for one second of speech, Poisson click train,  
 627 and standard periodic click train, respectively. Vertical scale is arbitrary but consistent across plots.  
 628 (D,E,F) Spectrograms of a smaller excerpt of the above stimuli, with darker colors corresponding to  
 629 higher power. (G,H,I) Power spectral density plots of the above stimuli, calculated from 30 s of data  
 630 using Welch’s method with a segment length of 5.67 ms, segment overlap of 50%, and Hann window.  
 631 Note that even though the speech recordings were gently high-pass filtered at 1000 Hz, there remains  
 632 plenty of power in the 125–1000 Hz range (G) and pitch information is clearly preserved (D; vertical  
 633 striations between 200 and 300 ms correspond to individual glottal pulses).

634

635 Figure 2. Comparison of ABR to standard periodic click trains and Poisson click trains. (A) The average  
 636 ABR waveform evoked by the standard, periodic click train at 44.1 clicks / s (black) and the  
 637 pseudorandom Poisson click train (gray; 44.1 clicks / s overall rate). Areas show  $\pm 1$  SEM. Both  
 638 responses are high-pass filtered at 150 Hz. The spike at  $-1$  ms is a stimulus artifact, and occurs before  
 639 0 ms to compensate for the 1 ms tube delay of the earphones. (B) The histogram of per-subject  
 640 correlation coefficients between the standard and Poisson click-evoked ABRs. Solid/dotted black lines  
 641 show median/quartiles. (C) Comparison of the Poisson click-evoked ABR with 150 Hz high-pass  
 642 filtering (gray) and without (i.e., broadband; blue). The latter is used as the benchmark response for the  
 643 remainder of the study.

644

645 Figure 3. Comparison of click-evoked responses (blue) with speech-derived responses (red). (A) The  
 646 average waveform across subjects (areas show  $\pm 1$  SEM). (B) The histogram of correlation coefficients  
 647 between the click-evoked and speech-derived stimuli for each subject. Solid/dotted black lines show  
 648 median/quartiles. (C) Individual subject responses, sorted by descending correlation coefficient. The  
 649 correlation is shown in the upper right corner.

650

651 Figure 4. Correlation of speech-derived and click-evoked wave V latencies (A) and amplitudes (B)  
652 across subjects. Because the click-evoked wave V is known to be subcortical, the strong correlations  
653 across subjects points to brainstem neural generators for the speech-derived response as well. Points  
654 have been jittered slightly to prevent visual overlap. Regression lines are shown with the 95%  
655 confidence interval shaded.

656

657 Figure 5. Comparison of female-narrated responses (green) with male-narrated responses (purple). (A)  
658 The average waveforms across subjects (areas show  $\pm 1$  SEM). (B) The histogram of correlation  
659 coefficients between the female-evoked and male-evoked stimuli for each subject. Solid/dotted black  
660 lines show median/quartiles. The speech-derived ABRs from the male and female narrators show  
661 strong similarities, but are not identical, indicating some talker dependence.

662

663 Figure 6. SNR as a function of data accumulation. SNR was calculated for each subject using the mean  
664 of all the data up to each recording epoch by computing the variance in the ABR time interval 0 to 20  
665 ms, and in the pre-stimulus noise interval  $-125$  to  $-10$  ms.

666

667 Figure 7. Changes to the range of lags and filtering parameters allows early, middle, and late  
668 responses to be analyzed from the same recording. (A) The average speech-derived auditory  
669 brainstem response with canonical waves V and VI labeled. (B) The middle latency response with its  
670 canonical waves labeled (low-pass frequency: 200 Hz). (C) The late auditory evoked potential (low-  
671 pass frequency: 20 Hz). Due to adaptation, the amplitudes for the later waves are much smaller than  
672 typically seen in the event-related potential literature. These peaks are also not given canonical labels  
673 as in A and B because their latencies do not directly correspond to the standard N1 and P2 peaks.  
674 Shaded areas show  $\pm 1$  SEM.













