

Auditory brainstem responses to continuous natural speech in human listeners

Abbreviated title

Auditory brainstem responses to natural speech

Authors

Ross K Maddox

Department of Neuroscience, University of Rochester, Rochester, NY 14642, USA
Department of Biomedical Engineering, University of Rochester, Rochester, NY 14627, 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

Corresponding Author

Ross K Maddox

Departments of Neuroscience and Biomedical Engineering, Del Monte Institute for Neuroscience
University of Rochester
201 Robert B. Goergen Hall
P.O. Box 270168
Rochester, NY 14627

Number of pages

17

Number of figures

6

Number of words

Abstract: 224

Introduction: 650

Discussion: 1,240

Conflict of interest

The authors declare no competing financial interests.

Acknowledgements

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

1 **ABSTRACT**

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

18

19 **SIGNIFICANCE STATEMENT**

20 Speech processing is usually studied in the cortex, but it starts in the auditory brainstem. However, a
21 paradigm for studying brainstem processing of continuous natural speech in human listeners has been
22 elusive due to practical limitations. Here we adapt methods that have been employed for studying
23 cortical activity to the auditory brainstem. We measure the response to continuous natural speech and
24 show that it is highly similar to the click-evoked response. The method also allows simultaneous
25 investigation of cortical activity with no added recording time. This discovery paves the way for studies
26 of speech processing in the human brainstem, including its interactions with higher order cognitive
27 processes originating in the cortex.

28

29 **INTRODUCTION**

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

46 One common method for measuring the brainstem response to speech is the complex ABR (cABR)
47 (Skoe and Kraus, 2010). The cABR represents the averaged response to repetitions of a short spoken
48 syllable (e.g., a ~40 ms “da”). It can be analyzed in the time domain, but because the stimulus is longer
49 than the response, ambiguity about the origin of response components arises. The voiced part of the
50 speech elicits a frequency following response (FFR) that can be analyzed in the frequency domain. The
51 FFR at the stimulus’s harmonics is reasoned to have subcortical origins because of the lower frequency
52 phase-locking limit in the auditory cortex (Joris et al., 2004), but a recent magnetoencephalography

53 study showed cortical contributions to the FFR (Coffey et al., 2016), rendering strong conclusions about
54 exclusively subcortical phenomena difficult to make.

55 A different method, used for studying cortical activity, treats the auditory evoked potential as the
56 impulse response of a linear system, which can be mathematically derived from known input and output
57 signals (Aiken and Picton, 2008; Lalor et al., 2009; Lalor and Foxe, 2010; Ding and Simon, 2012a,
58 2012b). Continuous natural speech is presented (input) while EEG is recorded (output), and the brain's
59 response is calculated through linear regression. Rather than raw audio, the regressor (i.e., input) used
60 is the amplitude envelope, which by construction contains no fast fluctuations, making it too slow for
61 studying subcortical nuclei. A recent study aimed at the brainstem used the amplitude envelope of a
62 speech stimulus's fundamental frequency as input, and the envelope of the EEG at that frequency as
63 output (Reichenbach et al., 2016). The response is a single wave with a peak latency of 10 ms,
64 suggesting brainstem involvement, but a width of 100 ms, making it impossible to exclude cortical
65 contributions.

66 Here we measured auditory brainstem activity in response to natural speech using a new paradigm.
67 The methods were based on cortical studies, with an important difference: the regressor was the
68 rectified speech audio, meaning that fine structure was largely preserved. The speech-evoked
69 responses were very similar to click-evoked ABRs, most notably in the presence of a distinct wave V.
70 Because the latency of the wave V peak is shorter than a cortical source could produce, it can be
71 unambiguously attributed to subcortical generators. We show that it is possible to study speech
72 processing in the human brainstem, paving the way for subcortical studies of attention, language, and
73 other cognitive processes.

74

75 **MATERIALS AND METHODS**

76 Experimental design and statistical analysis

77 Our goal was to measure the speech-evoked ABR in human listeners and validate it against the click-
78 evoked response. We first recorded click-evoked responses to pseudorandomly timed click trains and
79 then validated them against the responses evoked by standard, periodic click trains. We then compared
80 the speech-evoked response to the pseudorandom click-evoked response. We validated by comparing
81 the overall morphology, as well as the presence and latency of wave V in the speech-evoked response.

82 All subjects' click- and speech-evoked responses were plotted individually. To compare the similarity of
83 two responses from a single subject (e.g., the click-evoked response to the speech-evoked response),
84 Pearson's product-moment correlation was used. The median and interquartile range of each
85 distribution of correlation coefficients across subjects was reported, in addition to plotting its histogram.
86 Two distributions of correlation coefficients were compared using Wilcoxon's signed-rank test for non-
87 normal distributions. When comparing wave V latencies across stimulus conditions, a paired Student's
88 t-test was used to determine if the means differed.

89

90 Subjects

91 All experiments were done under a protocol approved by the University of Washington Institutional
92 Review Board. All subjects gave informed consent prior to participation, and were compensated for
93 their time. We collected data from 24 subjects (17 females). The mean age was 27.8 years, with a
94 standard deviation of 6.9 and a range of 19–45. Subjects had normal hearing, defined as audiometric
95 thresholds of 20 dB HL or better in both ears at octave frequencies ranging from 250 to 8000 Hz. All
96 subjects identified English as their first language except for two, who identified a different language but
97 had been speaking English daily for over twenty years.

98

99 EEG recording

100 Scalp potentials were recorded with passive Ag/AgCl electrodes, with the positive and negative
101 electrodes connected to a differential preamplifier (Brainvision LLC, Greenboro, SC). The positive
102 electrode was at location FCz in standard 10-20 coordinate system. The negative (reference) electrode

103 was clipped onto the subject's left earlobe. The ground electrode was placed at Fpz. Data were high-
104 passed at 0.1 Hz during recording (additional filtering occurred offline).

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

115

116 Speech stimuli

117 Speech stimuli were taken from two audiobooks. The first was *A Wrinkle in Time* (L'Engle, 2012), read
118 by a female narrator. The second was *The Alchemyst* (Scott, 2007), read by a male narrator. The
119 audiobooks were purchased on compact disc and ripped to uncompressed wav files to avoid data
120 compression artifacts. They were resampled to 24,414 Hz, the native rate of the RP2 presentation
121 system. They were then processed so that any silent pauses in the speech longer than 0.5 s were
122 truncated to 0.5 s. Because the ABR is principally driven by higher stimulus frequencies (Abdala and
123 Folsom, 1995), the speech was gently high-passed with a first-order Butterworth filter with a cutoff of
124 1,000 Hz and a slope of 6 dB / octave. The speech was still natural sounding and intelligible. This filter
125 also helped to compensate for low-frequency spectral differences between the male and female
126 narrator around their fundamental frequencies. After that, the speech was normalized to an average
127 root-mean-square amplitude that matched that of a 1 kHz tone at 75 dB SPL. Figures 1A,D,G show the
128 pressure waveform of the word "Thursday" spoken by the male narrator, the spectrogram of that word's
129 first syllable, and the power spectral density (PSD) of a 30 s segment of the female and male speech
130 stimuli.

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

140

141 Click stimuli

142 Click stimuli were aperiodic trains of rarefaction clicks lasting 82 μ s (representing two samples at the
143 24,414 Hz sampling rate, which was closest possible to the standard 100 μ s click duration with our
144 hardware). Clicks were timed according to a Poisson point process with a rate of 44.1 clicks / s. The
145 timing of one click had no correlation with the timing of any other click in the train, rendering the
146 sequence spectrally white in the statistical sense. A pair of 30 s sequences was created and presented
147 dichotically 20 times to each subject, meaning that 26,460 clicks contributed to each ear's response.
148 The responses presented herein are the sum of the monaural responses. Clicks were presented at 75
149 dB peak-to-peak equivalent SPL (i.e., the amplitude of clicks matched the peak-to-peak amplitude of a
150 1 kHz sinusoid presented at 75 dB SPL).

151 While no previous study has used exactly this type of click timing, several have used various types of
152 pseudorandom sequences (Burkard et al., 1990; Thornton and Slaven, 1993; Delgado and Ozdamar,
153 2004; Holt and Özdamar, 2014). Uniformly, these studies find that the ABRs from randomized versus
154 periodic click trains are highly similar at the same stimulation rates. Random timing has two main

155 benefits over the much more common periodic timing: 1) the analysis window for the response can be
156 extended arbitrarily to any beginning and end point without fear of temporal wrapping, and 2) no high-
157 pass filtering is necessary to remove the strong frequency component at the (periodic) presentation
158 rate, because it does not exist. A third benefit, specific to this study, is that the same analysis could be
159 done to compute the speech-evoked and the click-evoked ABR, yielding a more direct comparison
160 between the two. Figures 1B,E,H show part of a Poisson click train in the same manner that Figs.
161 1A,D,G do for speech.

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

166

167 Data analysis

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

174

175 *Speech stimuli preprocessing*

176 Before we could derive the speech response, we needed to calculate the regressor from the audio
177 data. The auditory brain is mostly agnostic to the sign of an acoustic input, as evidenced by the high
178 degree of similarity between evoked responses to compression versus rarefaction clicks (Møller et al.,
179 1995). For this reason, some sort of rectifying nonlinearity applied to the input speech is needed as a
180 preprocessing step. We used half-wave rectification. Specifically, we performed all analyses twice—
181 once keeping the positive peaks, and then a second time keeping the inverted negative peaks—and
182 then averaged the resulting responses, in a process akin to the compound peristimulus time histogram
183 used by Pfeiffer and Kim (1972). This significantly reduced, but did not eliminate, stimulus artifacts,
184 similar to the common technique of alternating polarity in the click-evoked ABR (Hall III, 2006).
185 Following rectification, the data were downsampled from 24,414 Hz to the EEG recording rate of 10,000
186 Hz.

187

188 *Click train preprocessing*

189 Owing to its extreme sparsity, downsampling a click train using standard methods would result in
190 significant signal processing artifact, viz., Gibbs ringing. We instead used the list of click times from the
191 original click train (24,414 Hz sampling rate) and created a click train at 10,000 Hz sampling rate by
192 placing unit-height single-sample impulses at the closest integer indices corresponding the original click
193 times.

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

201

202 *EEG preprocessing*

203 EEG data were first high-pass filtered at 1 Hz (first-order Butterworth), and then notch filtered at 60,
204 180, and 300 Hz with 5 Hz wide second-order infinite impulse response notch filters, designed with the

205 *iirnotch* function of the SciPy python package (RRID:SCR_008058). Because of the continuous nature
206 of the stimuli, no epoch rejection was done. Instead, any time the EEG potential crossed $\pm 100 \mu\text{V}$, a 1 s
207 segment of the response was zeroed, centered around the offending sample, removing it from the
208 calculation. This operation effectively reduced the energy of an epoch. So that the amplitude of the
209 calculated response was not affected, the EEG data for each epoch was multiplied by a corrective gain
210 factor g_r :

$$211 \quad g_r = N / (N - N_r),$$

212 where N is the total number of samples in the epoch and N_r is the number of rejected samples. After
213 filtering and resampling, the data were segmented into epochs that started with the stimulus onset and
214 ended 100 ms after the stimulus (epochs were thus 64.1 s long for speech stimuli and 30.1 s long for
215 clicks).

216

217 *Response calculation*

218 We used linear least-squares regression to calculate the responses, as in previous work (Lalor et al.,
219 2009). The response was considered to be the weights over a range of time lags that best
220 approximated the EEG output as the weighted sum of the input stimulus regressor over those lags. For
221 the sake of computational and memory efficiency, the stimulus autocorrelation matrix and stimulus-
222 response cross-correlation were both calculated via their Fourier counterparts using frequency-domain
223 multiplication. These specific methods have been incorporated into the mne-python package (Gramfort
224 et al., 2013) (RRID:SCR_005972) ([https://github.com/mne-tools/mne-
python/blob/8fc2a545f494de0f828b931f2285dbff426e72ad/mne/decoding/time_delaying_ridge.py](https://github.com/mne-tools/mne-python/blob/8fc2a545f494de0f828b931f2285dbff426e72ad/mne/decoding/time_delaying_ridge.py)). No
225 regularization was employed. The response weights were calculated over the range of lags spanning
226 -150 to 350 ms. After the response was calculated, it was low-pass filtered at $2,000$ Hz (first-order
227 Butterworth), and then baseline corrected by subtracting the mean potential between -10 and 0 ms
228 from the whole response. For the speech stimuli, the response to each narrator was calculated
229 separately, and then averaged to calculate each subject's speech-evoked response.
230

231

232 *Speech-evoked response amplitude normalization*

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

$$240 \quad g_n = E_i(\sigma_{c,i}) / E_i(\sigma_{s,i}),$$

241 where $\sigma_{c,i}$ is the standard deviation of subject i 's click-evoked response in the range of 0 – 30 ms, $\sigma_{s,i}$ is
242 the same for the speech-evoked response, and E_i represents the mean over subjects. All speech-
243 evoked responses shown in microvolts have been multiplied by g_n . In our study g_n had a value of 27.5 ,
244 but it must be stressed that this value depends on the unitless scale chosen for storing the digital audio,
245 and is thus not suitable for use in other studies. For this reason no direct amplitude comparisons were
246 made between click- and speech-evoked responses. Instead, their morphologies and wave V latencies
247 were compared.

248

249 Standard ABR measurement

250 The ABR to the periodic click trains was calculated through traditional averaging rather than regression.
251 The raw data were notch filtered to remove line noise and low-pass filtered at $2,000$ Hz as described
252 above. However, the high-pass filter was different: a causal second order Butterworth filter with a cutoff
253 of 150 Hz was used to be consistent with standard practice and to generate a canonical waveform
254 (Burkard et al., 2006; Hall III, 2006). The response to each click presentation was then epoched from

255 –3 ms to 19.7 ms, which was the longest window allowed by the periodic click rate of 44.1 clicks / s
256 before temporal wrapping occurred. Filtered epochs were rejected if the peak-to-peak amplitude
257 exceeded 100 μ V.

258

259 RESULTS

260 Poisson click trains yield canonical ABRs

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

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

279

280 Early speech-evoked responses exhibit brainstem response characteristics

281 Broadly speaking, there were strong similarities between the early (< 30 ms) click-evoked and speech-
282 evoked responses (Fig. 3A). In this latency range, responses are likely to progress from brainstem to
283 thalamus and primary auditory cortex as latency increases. We will first make whole-waveform
284 comparisons, and then consider specific canonical ABR components.

285 To compare the overall waveforms, we computed Pearson's correlation coefficient of the speech- and
286 click-evoked waveforms for each subject in the range of 0–30 ms (Fig. 3B). The median correlation
287 coefficient was 0.70 (interquartile range 0.63–0.75). Figure 3C shows each subject's click- and speech-
288 evoked response, in descending correlation order. In our speech-evoked responses, waves I–IV were
289 “smeared” together. However, we found a clear wave V in individual subjects' responses as well as the
290 grand average. Wave VI was also visible in the grand average, but was less consistent at the
291 individual-subject level.

292 We identified wave V by low-pass filtering at 1,000 Hz with a zero-phase filter and finding the peak of
293 the waveform in the 5–7 ms range. For the click-evoked responses, wave V was present for all
294 subjects, with a latency of 6.50 ± 0.25 ms (mean \pm standard deviation). For speech-evoked responses,
295 wave V was present for all subjects, with a latency of 6.17 ± 0.30 ms. The speech-evoked wave V
296 preceded the click-evoked by 0.26 ms ($t(23) = 6.6$, $p = 1 \times 10^{-6}$, paired t-test). As shown in Fig. 4, the
297 click-evoked and speech-evoked wave V latencies were correlated across subjects ($r = 0.75$, $p =$
298 3×10^{-5} , Pearson's product-moment). This shows a strong correspondence between the click-evoked
299 and speech-evoked ABR.

300 In some subjects' speech-evoked waveforms there are early peaks that seem to resemble waves I and
301 III. However, these are likely driven by recording artifacts (electromagnetic leakage of the earphone
302 driving signal into the EEG electrode recording). While it may have been possible to reduce these
303 artifacts further through additional signal processing, we did not do that for sake of simplicity and
304 transparency. However, it is important to note that a simple modification to the paradigm—alternating
305 the polarity of the speech stimulus—should all but remove stimulus artifacts in the future. This could be

306 done at the level of the 64 s epochs, or it could be done at the word or phrase level, as long as the
307 phase inversions were hidden by silent gaps in the speech.

308

309 Speech responses depend minimally on sex of talker stimuli

310 One important question is whether the speech-evoked response maintains its morphology independent
311 of the specific input stimulus, or if it depends on the specific narrator. To investigate this, we split the
312 responses to male- and female-narrated trials and compared them to determine the role that the
313 difference in the narrators' input spectra might play. The grand average waveforms for the two narrators
314 are of the same magnitude and overall shape, despite the differing spectra of their input stimuli (Fig.
315 5A). The median female-male correlation coefficient was 0.73 (interquartile range 0.60–0.83; Fig. 5B).
316 Figure 5C shows each subject's response to the female- and male-narrated speech, in the same order
317 as Fig. 3C to allow comparison.

318 While perfect overlap would be indicated by correlation coefficients of 1.0, splitting the data in half (viz.,
319 into male- and female-narrated epochs) adds noise to each of the responses. To put the male-female
320 correlation coefficients in context, we can split the data a different way and compare. We split the data
321 into halves that contained the same number of male and female epochs (i.e., each split contained 10
322 male and 10 female trials). We then compared those waveforms in the same way as above. The
323 median correlation coefficient between splits was 0.83 (interquartile range 0.70–0.91). We compared
324 the male-female split coefficients to these arbitrarily split coefficients, and found a significant difference
325 ($T(23) = 58$, $p = 0.009$, Wilcoxon signed-rank test). This indicates that while the responses to female
326 and male-uttered speech are very similar, there is still some dependence on the stimulus.

327

328 **DISCUSSION**

329 Early speech responses are interpretable as ABRs

330 The major goal of this work was to study the response of the human auditory brainstem to naturally
331 spoken, continuous speech. We derived the speech-evoked responses using regression and validated
332 them against click-evoked responses. Comparison of the speech-evoked and click-evoked ABR
333 revealed a high degree of morphological similarity between waveforms, similar overall wave V
334 latencies, and a strong correlation between click- and speech-evoked wave V latency across subjects.
335 Taken together, these results show that the speech-evoked ABR is just that—the response of the
336 auditory brainstem.

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

346

347 Subcortical and cortical responses are available simultaneously

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

355 This flexibility is illustrated in Fig. 6. Figure 6A shows the speech-evoked ABR, Fig. 6B extends the
356 window and employs a low-pass filter appropriate for viewing the middle latency response (Hall III,
357 2006), and Fig. 6C extends the time window further and lowers the low-pass frequency to accentuate
358 late auditory evoked potentials of cortical origin. If a full electrode montage (and sufficient hard drive
359 space) is available, the interaction of brainstem processing with any number of cortical processes is
360 now possible to investigate under natural conditions.

361

362 Filtering must be done carefully

363 It is common practice in EEG experiments to use zero-phase filters whose impulse responses are non-
364 causal and symmetric about 0 lag. This is done to preserve the latencies of the peaks and is
365 appropriate in most cases. However, the strength of the present approach lies in using the latency of
366 the response peaks to confirm their subcortical origin. If a non-causal filter is used to filter the EEG
367 data, then it is possible that a peak at a latency corresponding to cortical activity (e.g., 25 ms) could
368 affect the response waveform at brainstem latencies (e.g., 6 ms). This could have the result of
369 erroneous findings that attribute cortical phenomena to subcortical nuclei. Thus, the following two
370 guidelines should be followed for experiments specifically aimed at the auditory brainstem. First, EEG
371 data should be filtered with causal filters. Second, when calculating regressors, any filtering that is done
372 to the input stimulus should be anti-causal (i.e., with an impulse response has values only at negative
373 lags). The latter can be practically accomplished by reversing the signal in time, filtering it with a
374 standard causal filter, and then reversing that result. Using causal filters will inevitably affect the
375 latencies of peaks, but this can be mitigated by filtering sparingly (i.e., as broadband as the specific
376 analyses will allow) with low-order filters.

377

378 Responses to arbitrary stimuli can be measured

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

383 The present study suggests that a range of stimuli can be used. First, we consider the main
384 comparison: speech-evoked to click-evoked ABR. Natural speech is different by almost any metric from
385 Poisson click trains, and yet the responses that we find through regression are very similar (Fig. 3A,B).
386 Second, we consider the responses to female versus male speech. Males typically speak at a
387 fundamental frequency about half that of females. Such a difference, when estimating the response of a
388 highly non-linear system using linear methods, could have resulted in major differences in the response
389 waveforms, but this was not the case (Fig. 5A,B). Taken together, it is reasonable to expect that the
390 technique could be applied to other real-world non-speech stimuli such as music or environmental
391 sounds, as well any spectrally rich synthetic stimulus of interest in the lab.

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

398

399 Other regressors may offer improvements

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

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

415

416 Conclusions and future directions

417 Here we present and validate a method for determining the response of the auditory brainstem to
418 continuous, naturally uttered, non-repeated speech. Speech processing involves a complex network
419 that ranges from the earliest parts of the auditory pathway to auditory and association cortices. The
420 techniques described here facilitate new neuroscience experiments by making it possible to measure
421 activity across the auditory neuraxis while human subjects perform natural and engaging tasks. These
422 paradigms will allow study of the subcortical effects of language learning and understanding, attention,
423 multisensory integration, and many other cognitive processes.

424

425 **REFERENCES**

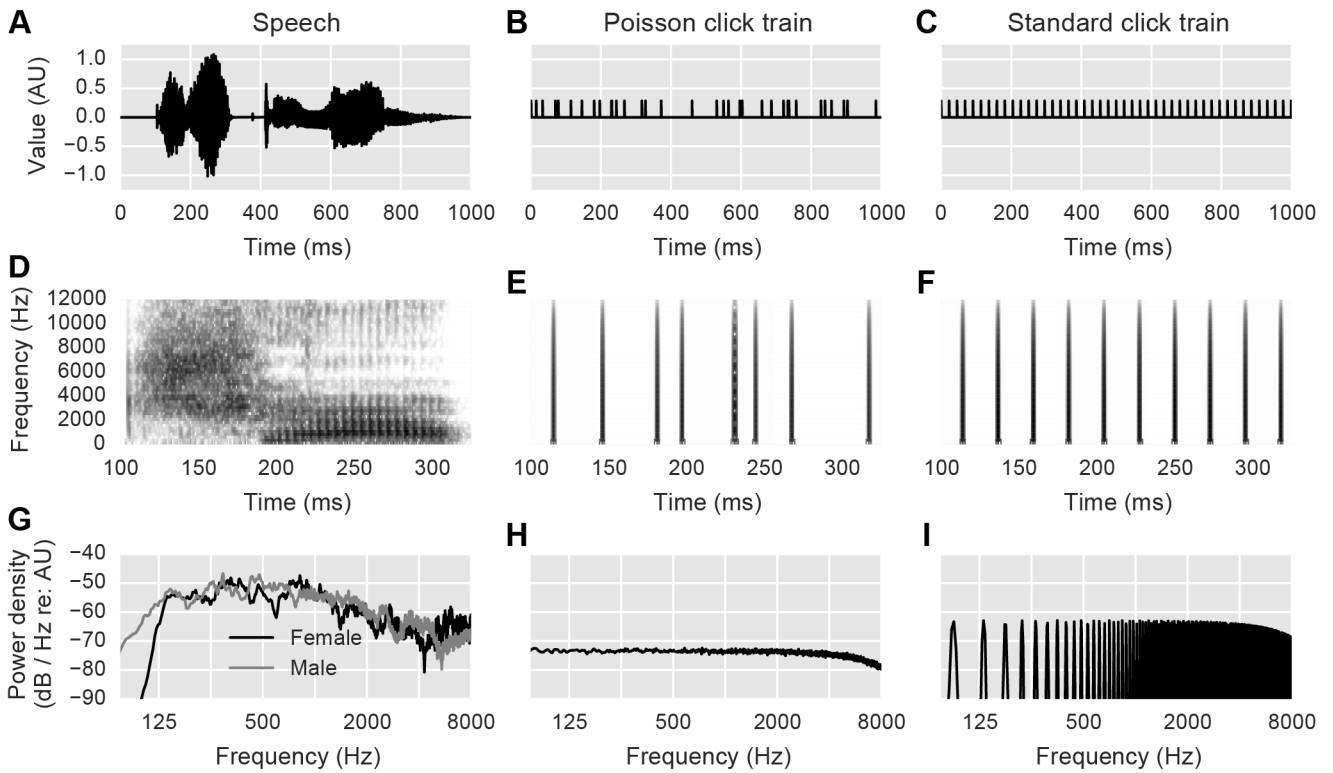
- 426 Abdala C, Folsom RC (1995) The development of frequency resolution in humans as revealed by the
427 auditory brain-stem response recorded with notched-noise masking. *J Acoust Soc Am* 98:921–
428 930.
- 429 Aiken SJ, Picton TW (2008) Human cortical responses to the speech envelope. *Ear Hear* 29:139–157.
- 430 Bidelman GM (2015) Towards an optimal paradigm for simultaneously recording cortical and brainstem
431 auditory evoked potentials. *J Neurosci Methods* 241:94–100.
- 432 Burkard R, Shi Y, Hecox KE (1990) A comparison of maximum length and Legendre sequences for the
433 derivation of brain-stem auditory-evoked responses at rapid rates of stimulation. *J Acoust Soc*
434 *Am* 87:1656–1664.
- 435 Burkard RF, Don M, Eggermont JJ (2006) *Auditory Evoked Potentials: Basic Principles and Clinical*
436 *Application*, 1st ed. Philadelphia: Lippincott Williams & Williams.
- 437 Carney LH, Li T, McDonough JM (2015) Speech Coding in the Brain: Representation of Vowel
438 Formants by Midbrain Neurons Tuned to Sound Fluctuations. *eNeuro:ENEURO*.0004-15.2015.
- 439 Coffey EBJ, Herholz SC, Chepesiuk AMP, Baillet S, Zatorre RJ (2016) Cortical contributions to the
440 auditory frequency-following response revealed by MEG. *Nat Commun* 7:11070.
- 441 Delgado RE, Ozdamar O (2004) Deconvolution of evoked responses obtained at high stimulus rates. *J*
442 *Acoust Soc Am* 115:1242–1251.
- 443 Di Liberto GM, Lalor EC (2017) Indexing cortical entrainment to natural speech at the phonemic level:
444 Methodological considerations for applied research. *Hear Res* 348:70–77.
- 445 Ding N, Simon JZ (2009) Neural Representations of Complex Temporal Modulations in the Human
446 Auditory Cortex. *J Neurophysiol* 102:2731–2743.
- 447 Ding N, Simon JZ (2012a) Neural coding of continuous speech in auditory cortex during monaural and
448 dichotic listening. *J Neurophysiol* 107:78–89.

- 449 Ding N, Simon JZ (2012b) Emergence of neural encoding of auditory objects while listening to
450 competing speakers. *Proc Natl Acad Sci* 109:11854–11859.
- 451 Gramfort A, Luessi M, Larson E, Engemann DA, Strohmeier D, Brodbeck C, Goj R, Jas M, Brooks T,
452 Parkkonen L, Hämäläinen M (2013) MEG and EEG data analysis with MNE-Python. *Front*
453 *Neurosci* 7 Available at: <http://journal.frontiersin.org/article/10.3389/fnins.2013.00267/full>
454 [Accessed August 29, 2017].
- 455 Grothe B, Pecka M (2014) The natural history of sound localization in mammals – a story of neuronal
456 inhibition. *Front Neural Circuits* 8:116.
- 457 Gu JW, Herrmann BS, Levine RA, Melcher JR (2012) Brainstem auditory evoked potentials suggest a
458 role for the ventral cochlear nucleus in tinnitus. *J Assoc Res Otolaryngol JARO* 13:819–833.
- 459 Hall III JW (2006) *New Handbook for Auditory Evoked Responses*, 1st ed. Boston: Pearson.
- 460 Holt FD, Özdamar Ö (2014) Simultaneous acquisition of high-rate early, middle, and late auditory
461 evoked potentials. In: 2014 36th Annual International Conference of the IEEE Engineering in
462 Medicine and Biology Society, pp 1481–1484.
- 463 Joris PX, Schreiner CE, Rees A (2004) Neural Processing of Amplitude-Modulated Sounds. *Physiol*
464 *Rev* 84:541–577.
- 465 Lalor EC, Foxe JJ (2010) Neural responses to uninterrupted natural speech can be extracted with
466 precise temporal resolution. *Eur J Neurosci* 31:189–193.
- 467 Lalor EC, Power AJ, Reilly RB, Foxe JJ (2009) Resolving Precise Temporal Processing Properties of
468 the Auditory System Using Continuous Stimuli. *J Neurophysiol* 102:349–359.
- 469 L'Engle M (2012) *A Wrinkle in Time*. Listening Library. Available at: [https://www.audible.com/pd/Kids/A-](https://www.audible.com/pd/Kids/A-Wrinkle-in-Time-Audiobook/B006LPK3WS)
470 *Wrinkle-in-Time-Audiobook/B006LPK3WS*.
- 471 Møller AR, Jho HD, Yokota M, Jannetta PJ (1995) Contribution from crossed and uncrossed brainstem
472 structures to the brainstem auditory evoked potentials: A study in humans. *The Laryngoscope*
473 105:596–605.
- 474 Pfeiffer RR, Kim DO (1972) Response Patterns of Single Cochlear Nerve Fibers to Click Stimuli:
475 Descriptions for Cat. *J Acoust Soc Am* 52:1669–1677.
- 476 Reichenbach CS, Braiman C, Schiff ND, Hudspeth AJ, Reichenbach T (2016) The Auditory-Brainstem
477 Response to Continuous, Non-repetitive Speech Is Modulated by the Speech Envelope and
478 Reflects Speech Processing. *Front Comput Neurosci* 10 Available at:
479 <http://journal.frontiersin.org/article/10.3389/fncom.2016.00047/full> [Accessed June 1, 2017].
- 480 Scott M (2007) *The Alchemyst: The Secrets of the Immortal Nicholas Flamel, Book 1*. Listening Library.
481 Available at: <https://www.audible.com/pd/Teens/The-Alchemyst-Audiobook/B002V1JA16>.
- 482 Skoe E, Kraus N (2010) Auditory brainstem response to complex sounds: a tutorial. *Ear Hear* 31:302–
483 324.
- 484 Terreros G, Delano PH (2015) Corticofugal modulation of peripheral auditory responses. *Front Syst*
485 *Neurosci* 9 Available at: <http://journal.frontiersin.org/Article/10.3389/fnsys.2015.00134/abstract>
486 [Accessed December 8, 2015].
- 487 Thornton ARD, Slaven A (1993) Auditory brainstem responses recorded at fast stimulation rates using
488 maximum length sequences. *Br J Audiol* 27:205–210.

489 Wassenhove V van, Schroeder CE (2012) Multisensory Role of Human Auditory Cortex. In: The Human
490 Auditory Cortex, pp 295–331 Springer Handbook of Auditory Research. Springer, New York,
491 NY. Available at: https://link.springer.com/chapter/10.1007/978-1-4614-2314-0_11 [Accessed
492 September 15, 2017].

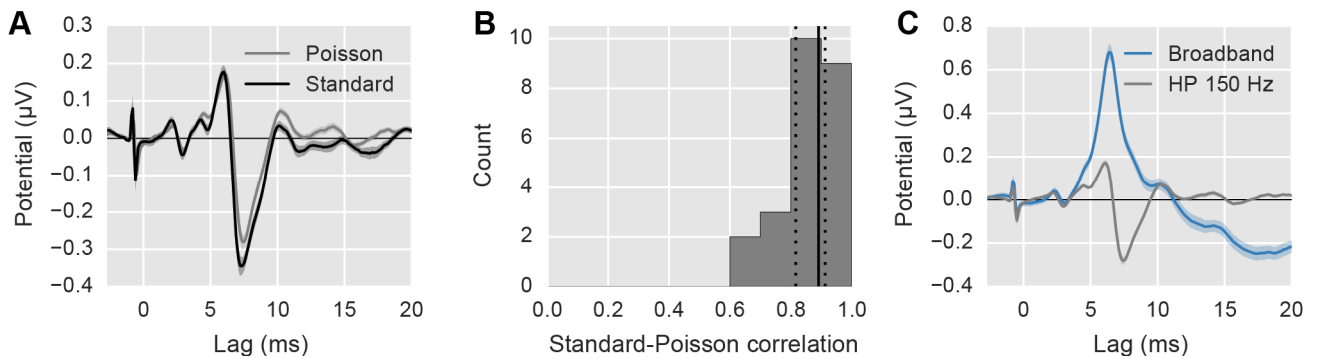
493

494 **FIGURES**



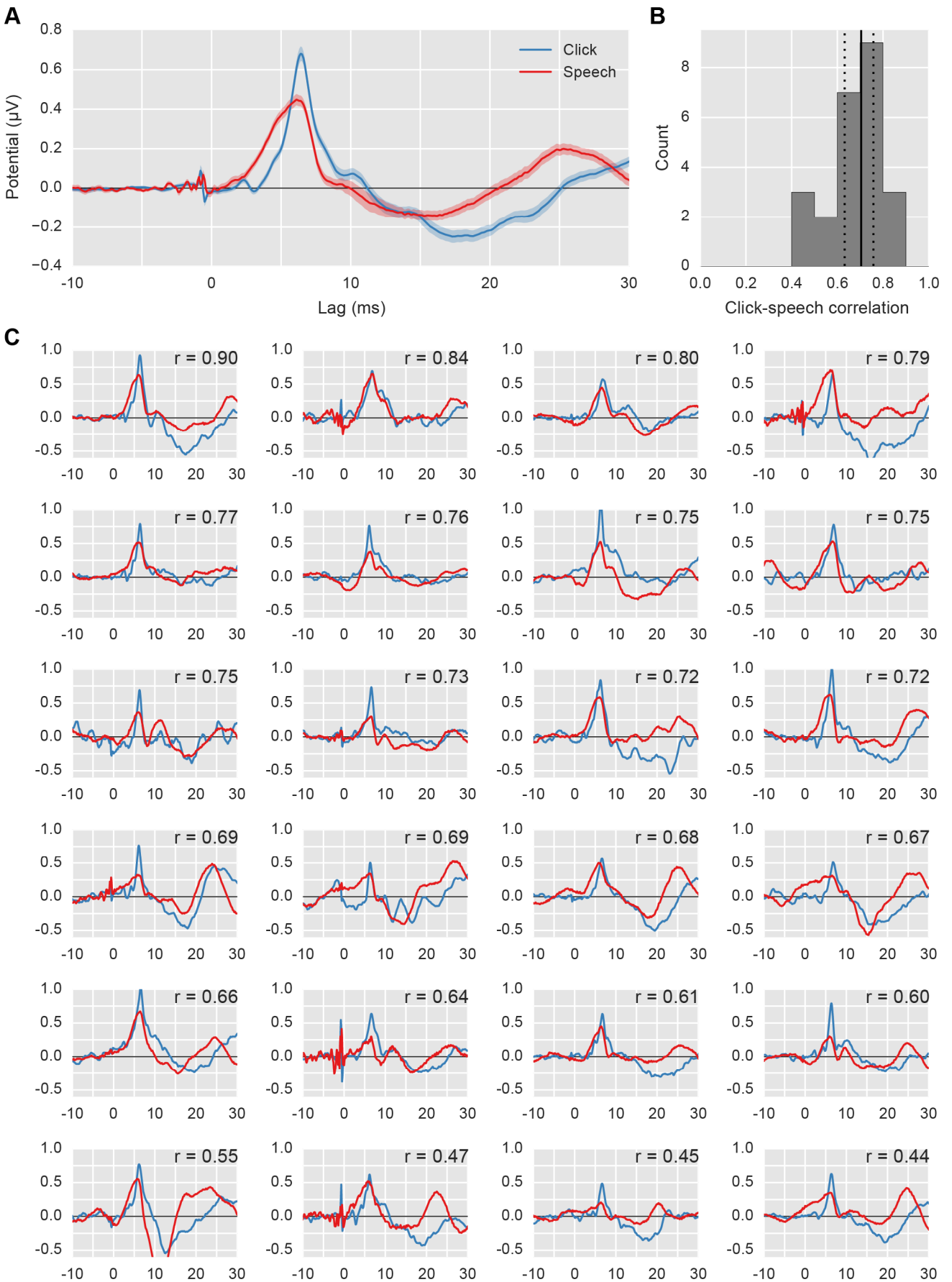
495
 496 Figure 1. Acoustic stimuli. (A,B,C) Pressure waveforms for one second of speech, Poisson click train,
 497 and standard periodic click train, respectively. Vertical scale is arbitrary but consistent across plots.
 498 (D,E,F) Spectrograms of a smaller excerpt of the above stimuli, with darker colors corresponding to
 499 higher power. (G,H,I) Power spectral density plots of the above stimuli, calculated from 30 s of data
 500 using Welch's method with a segment length of 5.67 ms, segment overlap of 50%, and Hann window.

501
 502
 503



504
 505 Figure 2. Comparison of ABR to standard periodic click trains and Poisson click trains. (A) The average
 506 ABR waveform evoked by the standard, periodic click train at 44.1 clicks / s (black) and the
 507 pseudorandom Poisson click train (gray; 44.1 clicks / s overall rate). Areas show ± 1 SEM. Both
 508 responses are high-pass filtered at 150 Hz. The spike at -1 ms is a stimulus artifact, and occurs before
 509 0 ms to compensate for the 1 ms tube delay of the earphones. (B) The histogram of correlation
 510 coefficients between the standard and Poisson click-evoked ABRs. Solid/dotted black lines show
 511 median/quartiles. (C) Comparison of the Poisson click-evoked ABR with 150 Hz high-pass filtering
 512 (gray) and without (i.e., broadband; blue). The latter is used as the benchmark response for the
 513 remainder of the study.

514

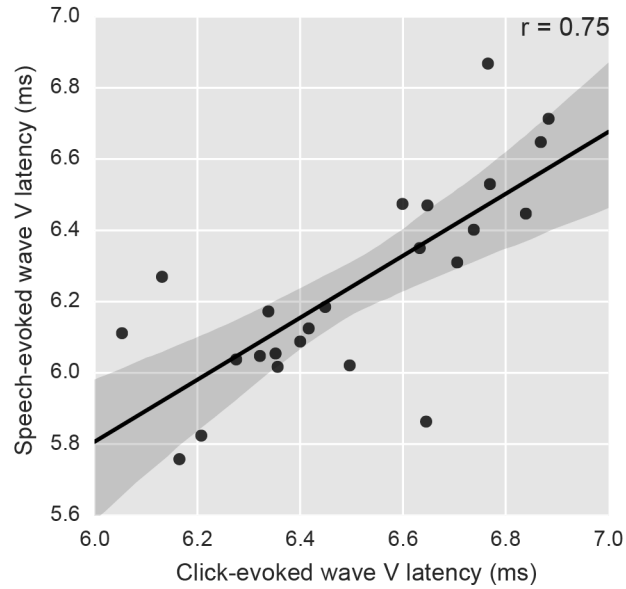


515
516 Figure 3. Comparison of click responses (blue) with speech responses (red). (A) The average
517 waveform across subjects (areas show ± 1 SEM). (B) The histogram of correlation coefficients between
518 the click-evoked and speech-evoked stimuli for each subject. Solid/dotted black lines show

519 median/quartiles. (C) Individual subject responses. The correlation coefficient is shown in the upper
520 right corner.

521

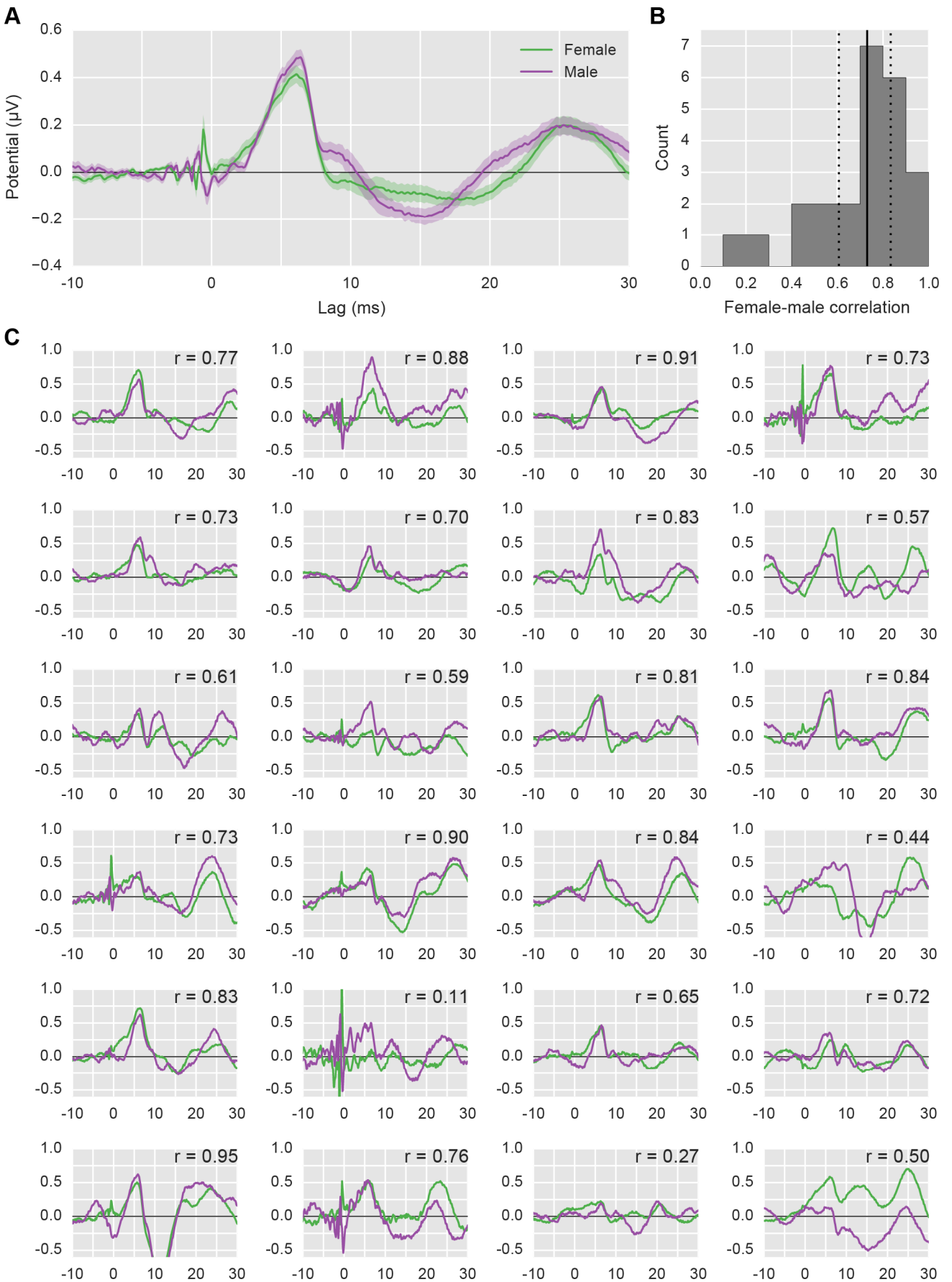
522



523

524 Figure 4. Speech-evoked versus click-evoked wave V latencies across subjects. The strong correlation
525 across subjects points to common neural generators. Points have been jittered slightly to prevent
526 overlap. Regression line is shown with the 95% confidence interval shaded.

527



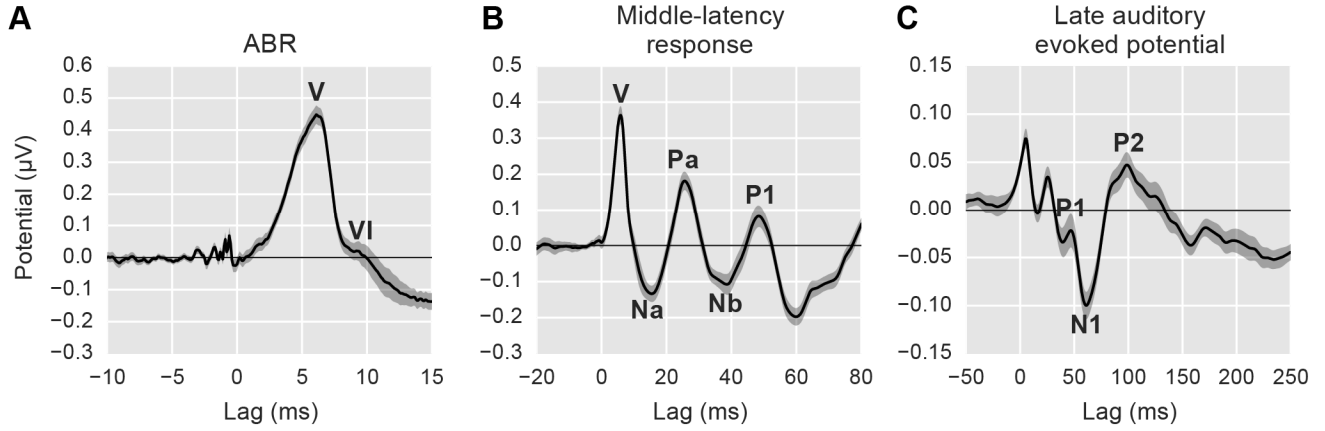
528
529 Figure 5. Comparison of female-narrated responses (green) with male-narrated responses (purple). (A)
530 The average waveform across subjects (areas show ± 1 SEM). (B) The histogram of correlation
531 coefficients between the female-evoked and male-evoked stimuli for each subject. Solid/dotted black

532 lines show median/quartiles. (C) Individual subject responses arranged in the same order as Fig. 3C for
533 easy comparison. The correlation coefficient is shown in the upper right corner. The poor correlation of
534 the worst subject ($r = 0.11$) is the result of a strong stimulus artifact.

535

536

537



538

539 Figure 6. Changes to the range of lags and filtering parameters allows early, middle, and late
540 responses to be analyzed from the same recording. (A) The speech-evoked auditory brainstem
541 response with canonical waves V and VI labeled. (B) The middle latency response with its canonical
542 waves labeled (low-pass frequency: 200 Hz). (C) The late auditory evoked potential with its canonical
543 waves labeled (low-pass frequency: 20 Hz). Shaded areas show ± 1 SEM.