

1 **EXPOSING DISTINCT SUBCORTICAL COMPONENTS OF THE AUDITORY BRAINSTEM RESPONSE**
2 **EVOKED BY CONTINUOUS NATURALISTIC SPEECH**

3 Melissa J Polonenko, Ph.D. (ORCID: 0000-0003-1914-6117)

4 Ross K Maddox, Ph.D.* (ORCID: 0000-0003-2668-0238)

5

6 Department of Biomedical Engineering

7 Department of Neuroscience

8 Del Monte Institute for Neuroscience

9 Center for Visual Science

10 University of Rochester

11

12 *correspondence:

13 University of Rochester

14 Goergen Hall

15 Box 270168

16 Rochester, NY 14627

17 Email: ross.maddox@rochester.edu

18

19 **ABSTRACT**

20 The auditory brainstem is important for processing speech, yet we have much to learn regarding the
21 contributions of different subcortical structures. These deep neural generators respond quickly, making
22 them difficult to study during dynamic, ongoing speech. Recently developed techniques have paved the
23 way to use natural speech stimuli, but the auditory brainstem responses (ABRs) they provide are
24 temporally broad and thus have ambiguous neural sources. Here we describe a new method that uses re-
25 synthesized “peaky” speech stimuli and deconvolution analysis of EEG data to measure canonical ABRs to
26 continuous naturalistic speech of male and female narrators. We show that in adults with normal hearing,
27 peaky speech quickly yields robust ABRs that can be used to investigate speech processing at distinct
28 subcortical structures from auditory nerve to rostral brainstem. We further demonstrate the versatility of
29 peaky speech by simultaneously measuring bilateral and ear-specific responses across different frequency
30 bands. Thus, the peaky speech method holds promise as a powerful tool for investigating speech
31 processing and for clinical applications.

32

33 **KEYWORDS**

34 speech, auditory brainstem response, evoked potentials, electroencephalography, assessment

35

36 INTRODUCTION

37 Understanding speech is an important, complex process that spans the auditory system from cochlea to
38 cortex. A temporally precise network transforms the strikingly dynamic fluctuations in amplitude and
39 spectral content of natural, ongoing speech into meaningful information, and modifies that information
40 based on attention or other priors (Mesgarani et al., 2009). Subcortical structures play a critical role in this
41 process – they do not merely relay information from the periphery to the cortex, but also perform important
42 functions for speech understanding, such as localizing sound (e.g., Grothe and Pecka, 2014) and encoding
43 vowels across different levels and in background noise (e.g., Carney et al., 2015). Furthermore, subcortical
44 structures receive descending information from the cortex through corticofugal pathways (Bajo et al., 2010;
45 Bajo and King, 2012; Winer, 2005), suggesting they may also play an important role in modulating speech
46 and auditory streaming. Given the complexity of speech processing, it is important to parse and understand
47 contributions from different neural generators. However, these subcortical structures are deep and respond
48 to stimuli with very short latencies, making them difficult to study during ecologically-salient stimuli such as
49 continuous and naturalistic speech. We created a novel paradigm aimed at elucidating the contributions
50 from distinct subcortical structures to ongoing, naturalistic speech.

51 Activity in deep brainstem structures can be “imaged” by the latency of waves in a surface electrical
52 potential (electroencephalography, EEG) called the auditory brainstem response (ABR). The ABR’s
53 component waves have been attributed to activity in different subcortical structures with characteristic
54 latencies: the auditory nerve contributes to waves I and II (~1.5–3 ms), the cochlear nucleus to wave III (~4
55 ms), the superior olivary complex and lateral lemniscus to wave IV (~5 ms), and the lateral lemniscus and
56 inferior colliculus to wave V (~6 ms) (Møller and Jannetta, 1983; review by Moore, 1987; Starr and
57 Hamilton, 1976). Waves I, III, and V are most often easily distinguished in the human response. Subcortical
58 structures may also contribute to the earlier P_0 (12–14 ms) and N_a (15–25 ms) waves (Hashimoto, 1982;
59 Kileny et al., 1987; Picton et al., 1974) of the middle latency response (MLR), which are then followed by
60 thalamo-cortically generated waves P_a , N_b , and P_b/P_1 (Geisler et al., 1958; Goldstein and Rodman, 1967).
61 ABR and MLR waves have a low signal-to-noise ratio (SNR) and require numerous stimulus repetitions to
62 record a good response. Furthermore, they are quick and often occur before the stimulus has ended.
63 Therefore, out of necessity, most human brainstem studies have focused on brief stimuli such as clicks,
64 tone pips, or speech syllables, rather than more natural speech.

65 Recent analytical techniques have overcome limitations on stimuli, allowing continuous naturally uttered
66 speech to be used. One such technique extracts the fundamental waveform from the speech stimulus and
67 finds the envelope of the cross-correlation between that waveform and the recorded EEG data (Forte et al.,
68 2017). The response has an average peak time of about 9 ms, with contributions primarily from the inferior
69 colliculus (Saiz-Alia and Reichenbach, 2020). A second technique considers the rectified broadband
70 speech waveform as the input to a linear system and the EEG data as the output, and uses deconvolution
71 to compute the ABR waveform as the impulse response of the system (Maddox and Lee, 2018). The
72 speech-derived ABR shows a wave V peak whose latency is highly correlated with the click response wave
73 V across subjects, demonstrating that the component is generated in the rostral brainstem. A third
74 technique averages responses to each chirp (click-like transients that quickly increase in frequency) in re-
75 synthesized “cheech” stimuli (CHirp spEECH; Miller et al., 2017) that interleaves alternating octave
76 frequency bands of speech and of chirps aligned with some glottal pulses (Backer et al., 2019). Brainstem
77 responses to these stimuli also show a wave V, but do not show earlier waves unless presented
78 monaurally over headphones (Backer et al., 2019; Miller et al., 2017). While these methods reflect
79 subcortical activity, the first two provide temporally broad responses with a lack of specificity regarding
80 underlying neural sources. None of the three methods shows the earlier canonical components such as
81 waves I and III that would allow rostral brainstem activity to be distinguished from, for example, the auditory
82 nerve. Such activity is important to assess, especially given the current interest in the potential
83 contributions of auditory nerve loss in disordered processing of speech in noise (Bramhall et al., 2019;
84 Liberman et al., 2016; Prendergast et al., 2017).

85 We asked if we could further assess underlying speech processing in multiple distinct early stages of the
86 auditory system by 1) evoking additional waves than wave V of the canonical ABR, and 2) measuring
87 responses to different frequency ranges of speech (corresponding to different places of origin on the
88 cochlea). The ABR is strongest to very short stimuli such as clicks, so we created “peaky” speech. The
89 design goal of peaky speech is to re-synthesize natural speech so that its defining spectrotemporal content

90 is unaltered but its pressure waveform consists of maximally sharp peaks so that it drives the ABR as
91 effectively as possible. The results show that peaky speech evokes canonical brainstem responses and
92 frequency-specific responses, paving the way for novel studies of subcortical contributions to speech
93 processing.

94

95 **RESULTS**

96 Broadband peaky speech yields more robust responses than unaltered speech

97 *Broadband peaky speech elicits canonical brainstem responses*

98 In previous work, brainstem responses to natural, on-going speech exhibited a temporally broad wave V
99 but no earlier waves (Maddox and Lee, 2018). We re-synthesized speech to be “peaky” with the primary
100 aim to evoke additional, earlier waves of the ABR that identify different neural generators. Indeed, Figure 1
101 shows that waves I, III, and V of the canonical ABR are clearly visible in the group average and in the
102 individual responses to broadband peaky speech. This means that broadband peaky speech, unlike the
103 unaltered speech, can be used to assess naturalistic speech processing at discrete parts of the subcortical
104 auditory system, from the auditory nerve to rostral brainstem. These responses represent weighted
105 averaged data from ~43 minutes of continuous speech (40 epochs of 64 s each), and were filtered at a
106 typical high-pass cutoff of 150 Hz to highlight the earlier ABR waves.

107 Morphology of the broadband peaky speech ABR was inspected and waves marked by a trained
108 audiologist (MJP) on 2 occasions that were 3 months apart. The intraclass correlation coefficients for
109 absolute agreement (ICC3) were ≥ 0.91 (lowest ICC3 95% confidence interval was 0.78–0.96 for wave I, p
110 < 0.01), indicating excellent reliability for chosen peak latencies. Waves I and V were identifiable in
111 responses from all subjects ($N = 22$), and wave III was clearly identifiable in 16 of the 22 subjects. These
112 waves are marked on the individual responses in Figure 1. Mean \pm SEM peak latencies for ABR waves I,
113 III, and V were 2.95 ± 0.10 ms, 5.11 ± 0.09 ms, and 6.96 ± 0.07 ms respectively. These mean peak
114 latencies are shown superimposed on the group average response in Figure 1 (bottom right). Inter-wave
115 latencies were 2.13 ± 0.05 ms ($N = 16$) for I–III, 1.78 ± 0.06 ms ($N = 16$) for III–V, and 4.01 ± 0.07 ms ($N = 22$)
116 for I–V. These peak inter-wave latencies fall within a range expected for brainstem responses but the
117 absolute peak latencies were later than those reported for a click ABR at a level of 60 dB sensation level
118 (SL) and rate between 50 to 100 Hz (Burkard and Hecox, 1983; Chiappa et al., 1979; Don et al., 1977).

119

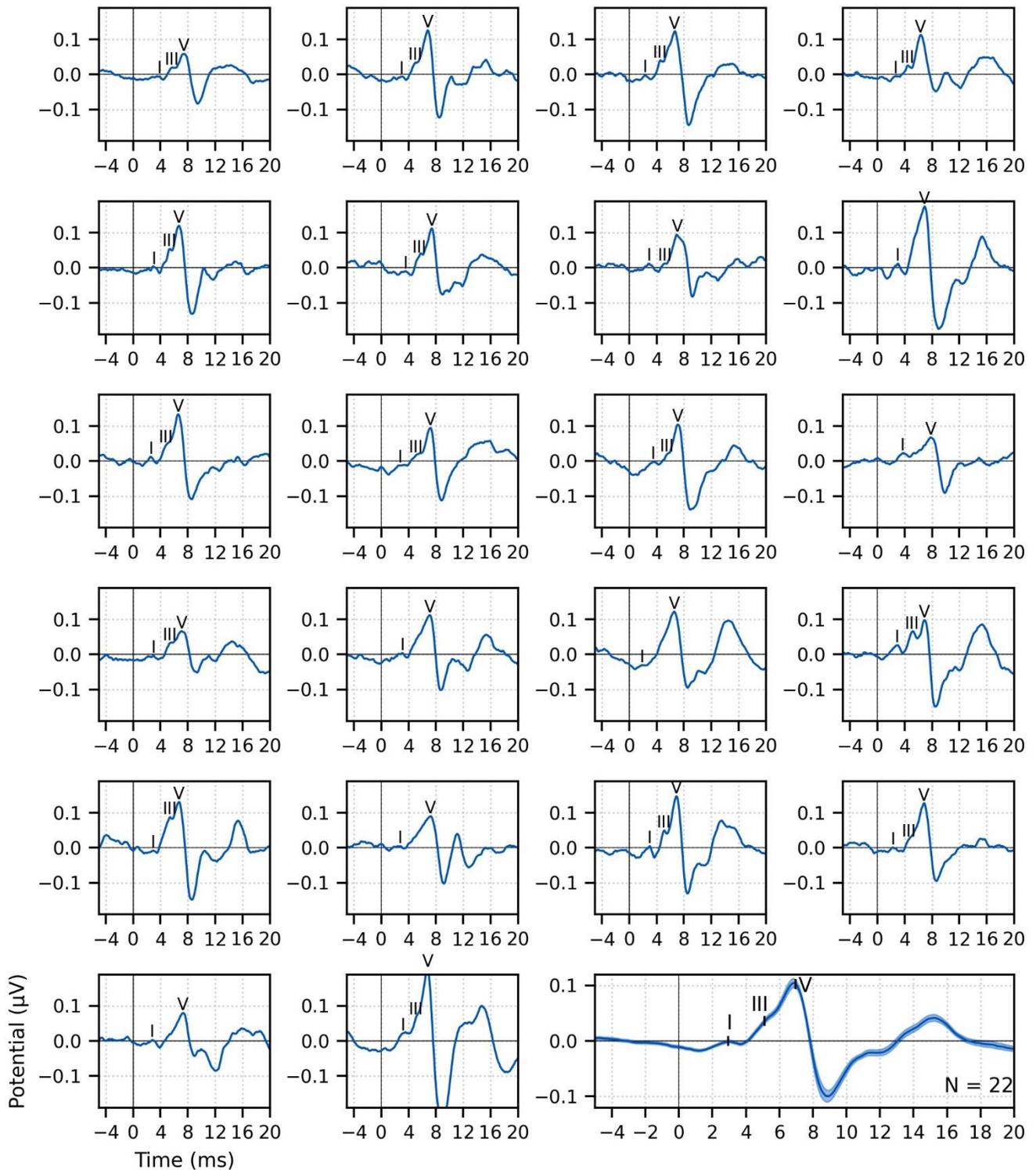


Figure 1. Single subject and group average (bottom right) weighted-average auditory brainstem responses (ABR) to ~43 minutes of broadband peaky speech. Area for the group average shows ± 1 SEM. Responses were high-pass filtered at 150 Hz using a first order Butterworth filter. Waves I, III, and V of the canonical ABR are evident in most of the single subject responses ($N = 22, 16, 22$ respectively), and are marked by the average peak latencies on the average response.

120

121 *More components of the ABR and MLR are present with broadband peaky than unaltered speech*

122 Having established that broadband peaky speech evokes robust canonical ABRs, we next compared both
123 ABR and MLR responses to those evoked by unaltered speech. To simultaneously evaluate ABR and MLR

124 components, a high-pass filter with a 30 Hz cutoff was used on the responses to ~43 minutes of each type
125 of speech. Figure 2A shows that overall there were morphological similarities between responses to both
126 types of speech; however, there were more early and late component waves to broadband peaky speech.
127 More specifically, whereas both types of speech evoked waves V, N_a and P_a, broadband peaky speech
128 also evoked waves I, often III (14–16 of 22 subjects depending if a 30 or 150 Hz high-pass filter cutoff was
129 used), and P₀. With a lower cutoff for the high-pass filter, wave III rode on the slope of wave V and was less
130 identifiable in the grand average shown in Figure 2A than that shown with a higher cutoff in Figure 1. Wave
131 V was more robust and sharper to broadband peaky speech but peaked slightly later than the broader
132 wave V to unaltered speech. For reasons unknown to us, the half-rectified speech method missed the MLR
133 wave P₀, and consequently had a broader and earlier N_a than the broadband peaky speech method,
134 though this missing P₀ was consistent with the results of Maddox and Lee (2018). These waveforms
135 indicate that broadband peaky speech is better than unaltered speech at evoking canonical responses that
136 distinguish activity from distinct subcortical and cortical neural generators.

137 Peak latencies for the waves common to both types of speech are shown in Figure 2B. Again, there was
138 good agreement in peak wave choices for each type of speech, with ICC3 ≥ 0.94 (the lowest two ICC3
139 95% confidence intervals were 0.87–0.98 and 0.92–0.99 for waves V and N_a to unaltered speech
140 respectively). As suggested by the waveforms in Figure 2A, mean ± SEM peak latency differences for
141 waves V, N_a, and P_a were longer for broadband peaky than unaltered speech by 0.29 ± 0.05 ms
142 (independents *t*-test, *t*₍₂₁₎ = 5.4, *p* < 0.01, *d* = 1.19), 4.40 ± 0.43 ms (*t*₍₂₁₎ = 9.9, *p* < 0.01, *d* = 2.16), and 0.40
143 ± 0.39 ms (*t*₍₂₀₎ = 1.0, *p* = 0.33, *d* = 0.22) respectively.

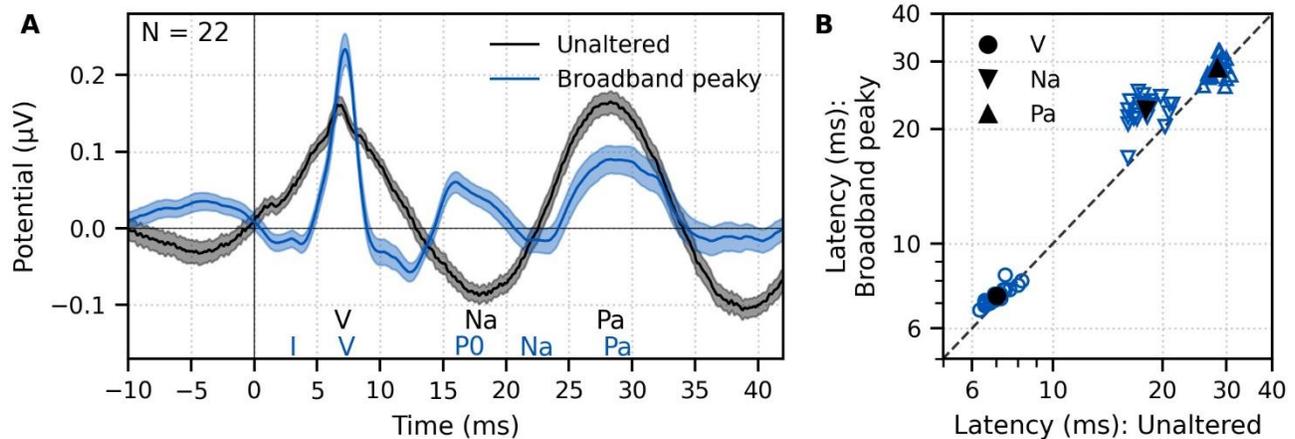


Figure 2. Comparison of auditory brainstem (ABR) and middle latency responses (MLR) to ~43 minutes each of unaltered speech and broadband peaky speech. (A) The average waveform to broadband peaky speech (blue) shows additional, and sharper, waves of the canonical ABR and MLR than the broader average waveform to unaltered speech (black). Responses were high-pass filtered at 30 Hz with a first order Butterworth filter. Areas show ± 1 SEM. (B) Comparison of peak latencies for ABR wave V (circles) and MLR waves N_a (downward triangles) and P_a (upward triangles) that were common between responses to broadband peaky and unaltered speech. Blue symbols depict individual subjects and black symbols depict the mean.

144 We also verified that the EEG data collected in response to broadband peaky speech could be regressed
145 with the half-wave rectified speech to generate a response. Details can be found in the supplemental
146 material and Supplemental Figure 1A.

147

148 Broadband peaky speech responses differ across talkers

149 We next sought to determine whether response morphologies depended on the talker identity. Responses
150 derived from unaltered speech show similar, but not identical, morphology between male and female
151 narrators, indicating some dependence on the stimulus (Maddox and Lee, 2018). To determine what extent
152 the morphology and robustness of peaky speech responses depend on a specific narrator's voice, we
153 compared waveforms and peak wave latencies for 32 minutes (30 epochs of 64 s each) each of male- and
154 female-narrated broadband peaky speech in 11 subjects. The average fundamental frequency was 115 Hz
155 for the male narrator and 198 Hz for the female narrator.

156 The group average waveforms to female- and male-narrated broadband peaky speech showed similar
157 canonical morphologies but were smaller and later for female-narrated ABR responses (Figure 3A), much
158 as they would be for click stimuli presented at higher rates (e.g., Burkard et al., 1990; Burkard and Hecox,
159 1983; Chiappa et al., 1979; Don et al., 1977; Jiang et al., 2009). All component waves of the ABR and MLR
160 were visible in the group average, although fewer subjects exhibited a clear wave III in the female-narrated
161 response (9 versus all 11 subjects). The median (interquartile range) male-female correlation coefficients
162 were 0.67 (0.60–0.77) for ABR lags of 0–15 ms with a 150 Hz high-pass filter, and 0.44 (0.35–0.61) for
163 ABR/MLR lags of 0–40 ms with a 30 Hz high-pass filter (Figure 3B).

164 To determine if this stimulus dependence was significantly different than variability introduced by averaging
165 only half the epochs (i.e., splitting by male- and female-narrated epochs), we reanalyzed the data split into
166 even and odd epochs. Each of the even/odd splits contained the same number of male- and female-
167 narrated epochs, and were evenly distributed over the entire recording session. The median (interquartile
168 range) odd-even correlation coefficients were 0.86 (0.80–0.95) for ABR lags and 0.47 (0.36–0.80) for
169 ABR/MLR lags. These odd-even coefficients were significantly higher than the male-female coefficients for
170 the ABR ($W_{(10)} = 0.0$, $p < 0.001$; Wilcoxon signed-rank test) but not when the response included the MLR
171 ($W_{(10)} = 18.0$, $p = 0.206$), indicating that the choice of narrator for using peaky speech impacts the
172 morphology of the early response.

173 As expected from the waveforms, peak latencies of component waves differed between male- and female-
174 narrated broadband peaky speech (Figure 3C). As before, $ICC3 \geq 0.83$ indicated good agreement in peak
175 wave choices (the lowest two 95% confidence intervals were 0.64–0.93 for N_a and 0.92–0.99 for I). Mean \pm
176 SEM peak latency differences (female – male) for wave I, III, and V of the ABR were 0.19 ± 0.08 ms ($t_{(10)} =$
177 2.15 , $p = 0.057$, $d = 0.68$), 0.51 ± 0.09 ms ($t_{(8)} = 5.56$, $p < 0.001$, $d = 1.97$), and 0.51 ± 0.11 ms ($t_{(10)} = 4.29$,
178 $p = 0.002$, $d = 1.36$) respectively. Latency differences were earlier for female-evoked P_0 (-1.27 ± 0.39 ms,
179 $t_{(10)} = -3.11$, $p = 0.011$, $d = -0.98$), but were not significant for later MLR peaks (N_a : -0.86 ± 0.57 ms, $t_{(8)} =$
180 -1.40 , $p = 0.197$, $d = -0.50$; P_a : -0.04 ± 0.44 ms, $t_{(9)} = -0.09$, $p = 0.933$, $d = -0.03$).

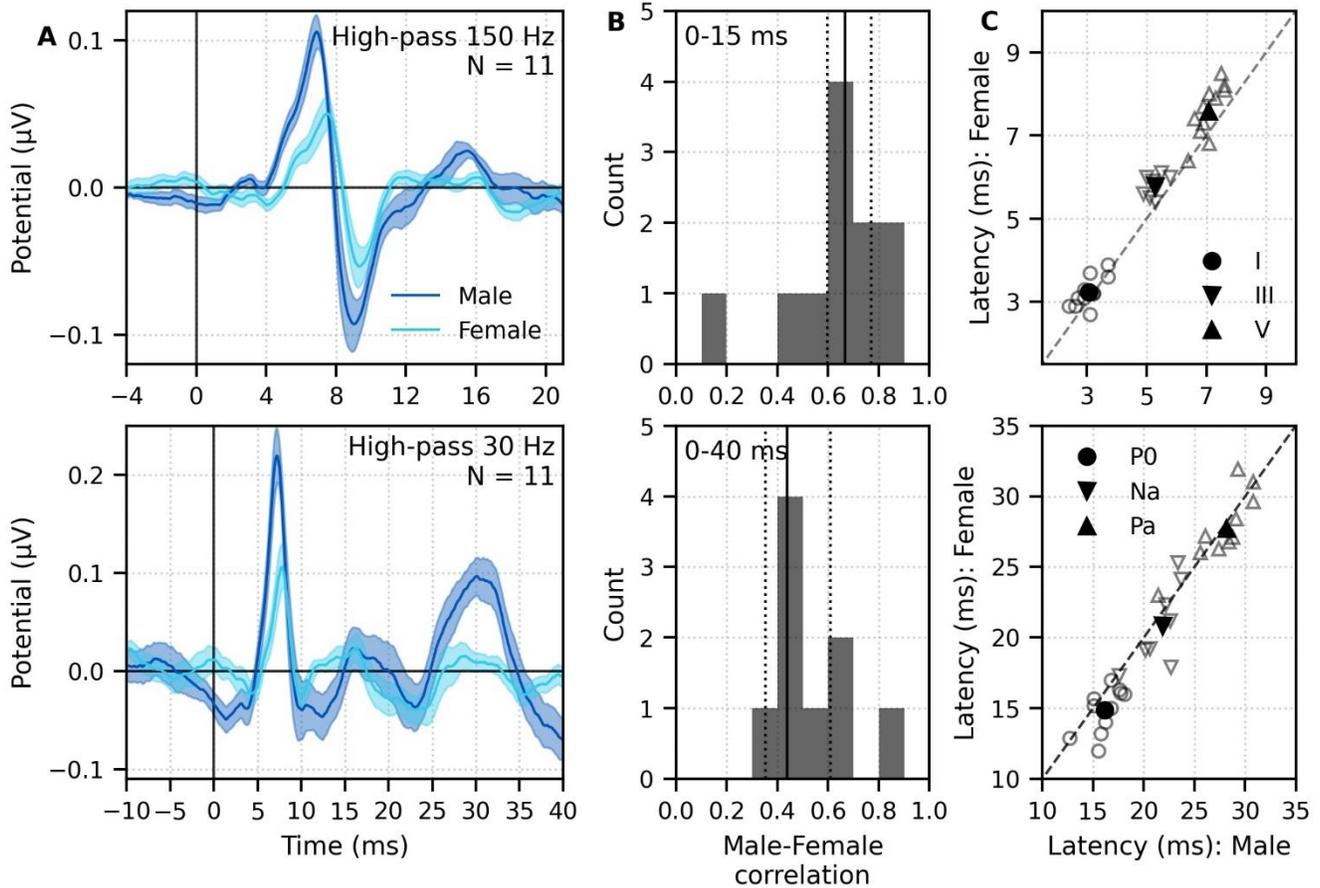


Figure 3. Comparison of responses to 32 minutes each of male (dark blue) and female (light blue) narrated re-synthesized broadband peaky speech. (A) Average waveforms across subjects (areas show ± 1 SEM) are shown for auditory brainstem response (ABR) time lags with high-pass filtering at 150 Hz (top), and both ABR and middle latency response (MLR) time lags with a lower high-pass filtering cutoff of 30 Hz (bottom). (B) Histograms of the correlation coefficients between responses evoked by male- and female-narrated broadband peaky speech during ABR (top) and ABR/MLR (bottom) time lags. Solid lines denote the median and dotted lines the inter-quartile range. (C) Comparison of ABR (top) and MLR (bottom) wave peak latencies for individual subjects (gray) and the group mean (black). ABR and MLR responses were similar to both types of input but are smaller for female-narrated speech, which has a higher glottal pulse rate. Peak latencies for female-evoked speech were delayed during ABR time lags but faster for early MLR time lags.

181

182 Multiband peaky speech yields frequency-specific brainstem responses to speech

183 *Frequency-specific responses show frequency-specific lags*

184 Broadband peaky speech gives new insights into subcortical processing of naturalistic speech. Not only are
 185 brainstem responses used to evaluate processing at different stages of auditory processing, but ABRs can
 186 also be used to assess hearing function across different frequencies. Traditionally, frequency-specific
 187 ABRs are measured using clicks with high-pass masking noise or frequency-specific tone pips. We tested
 188 the flexibility of using our new peaky speech technique to investigate how speech processing differs across
 189 frequency regions, such as 0–1, 1–2, 2–4, and 4–8 kHz frequency bands. To do this, we created new
 190 pulses trains with slightly different fundamental waveforms for each filtered frequency regions of speech,
 191 and then combined those filtered frequency bands together as multiband speech (for details, see the
 192 Multiband peaky speech subsection of Methods). Using this method, we took advantage of the fact that
 193 over time, stimuli with slightly different fundamental frequencies will be independent, yielding independent
 194 auditory brainstem responses. Therefore, the same EEG was regressed with each band's pulse train to
 195 derive the ABR and MLR to each frequency band.

196 Mean \pm SEM responses from 22 subjects to the 4 frequency bands (0–1, 1–2, 2–4, and 4–8 kHz) of \sim 43
 197 minutes of male-narrated multiband peaky speech are shown as colored waveforms with solid lines in
 198 Figure 4A. A high-pass filter with a cutoff of 30 Hz was used. Each frequency band response comprises a
 199 frequency-band-specific component as well as a band-independent common component, both of which are
 200 due to spectral characteristics of the stimuli and neural activity. The pulse trains are independent over time

201 in the vocal frequency range – thereby allowing us to pull out responses to each different pulse train and
202 frequency band from the same EEG – but they became coherent at frequencies lower than 72 Hz for the
203 male-narrated speech and 126 Hz for the female speech (see Figure 13 in Methods). This coherence was
204 due to all pulse trains beginning and ending together at the onset and offset of voiced segments and was
205 the source of the low-frequency common component of each band’s response. To remove the common
206 component, there are two options. First, we could simply high-pass the response at 150 Hz to filter out the
207 regions of spectral coherence in the stimuli, as shown by the waveforms with dashed lines in Figure 4B.
208 However, this method reduces the amplitude of the responses, which in turn affects response SNR and
209 detectability. The second option is to calculate the common activity across the frequency band responses
210 and subtract this waveform from each of the frequency band responses. This common component was
211 calculated by regressing the EEG to multiband speech with 6 independent “fake” pulses trains – pulse
212 trains with slightly different fundamental frequencies that were not used to create the multiband peaky
213 speech stimuli that were presented during the experiment – and then averaging across these 6 responses.
214 This common component waveform is shown by the dot-dashed gray line, which is superimposed with
215 each response to the frequency bands in Figure 4A. The subtracted, frequency-specific waveforms to each
216 frequency band are shown by the solid lines in Figure 4B. Of course, the subtracted waveforms could also
217 be high-pass filtered at 150 Hz to highlight earlier waves of the brainstem responses, as shown by the
218 dashed lines in Figure 4B. Overall, the frequency-specific responses showed characteristic ABR and MLR
219 waves with longer latencies for lower frequency bands, as would be expected from responses arising from
220 different cochlear regions. Also, waves I and III of the ABR were visible in the group average waveforms of
221 the 2–4 kHz ($\geq 41\%$ of subjects) and 4–8 kHz ($\geq 86\%$ of subjects) bands, whereas the MLR waves were
222 more prominent in the 0–1 kHz ($\geq 95\%$ of subjects) and 1–2 kHz ($\geq 54\%$ of subjects) bands.

223 These frequency-dependent latency changes for the frequency-specific responses are highlighted further in
224 Figure 4C, which shows mean \pm SEM peak latencies and the number of subjects who had a clearly
225 identifiable wave. ICC3 ≥ 0.89 indicated good agreement in peak wave choices (lowest two 95%
226 confidence intervals were 0.82–0.93 for P_a and 0.88–0.95 for N_a). The nonlinear change in peak latency
227 with frequency band was modeled using mixed effects regression by including orthogonal linear and
228 quadratic terms for frequency band and their interactions with wave, as well as random effects of intercept
229 and each frequency band term for each subject. A model was completed for each filter cutoff of 30 and 150
230 Hz. There were insufficient numbers of subjects with identifiable waves I and III for the 0–1 kHz and 1–2
231 kHz bands, so these waves were not included in the full model. Details of each model are described in
232 Supplemental Table 1. As expected, there were significantly different latencies for each MLR wave P_0 , N_a
233 and P_a compared to the ABR wave V (all effects of wave on the intercept $p < 0.001$, for each high-pass
234 filter cutoff of 30 and 150 Hz). The significant decrease in latency with frequency band (linear term, slope: p
235 < 0.001 for 30 and 150 Hz) was steeper (i.e., more negative) for MLR waves compared to the ABR wave V
236 (all $p < 0.001$ for interactions between wave and the linear frequency band term for 30 and 150 Hz). The
237 rate of latency decrease also changed significantly (quadratic frequency band term: $p = 0.001$ for 30 Hz
238 and $p < 0.001$ for 150 Hz), but in a similar way for each component wave (all $p > 0.091$ for interactions
239 between wave and the quadratic frequency band term, for 30 and 150 Hz models).

240 Next, the frequency-specific responses (i.e., multiband responses with common component subtracted)
241 were summed and the common component added to derive the entire response to multiband peaky
242 speech. As shown in Figure 5, this summed multiband response was strikingly similar in morphology to the
243 broadband peaky speech. Both responses were high-passed filtered at 150 Hz and 30 Hz to highlight the
244 earlier ABR waves and later MLR waves, respectively. The median (interquartile range) correlation
245 coefficients from the 22 subjects were 0.90 (0.86–0.9) for 0–15 ms ABR lags, and 0.55 (0.46–0.74) for 0–
246 40 ms MLR lags. The similarity verifies that the frequency-dependent responses are truly independent from
247 each other, and that these responses are complementary to the common component. If there were overlap
248 in the cochlear regions, for example, the summed response would not resemble the broadband response
249 to such a degree. The similarity also verified that the additional changes we made to create re-synthesized
250 multiband peaky speech did not significantly affect responses compared to broadband peaky speech.

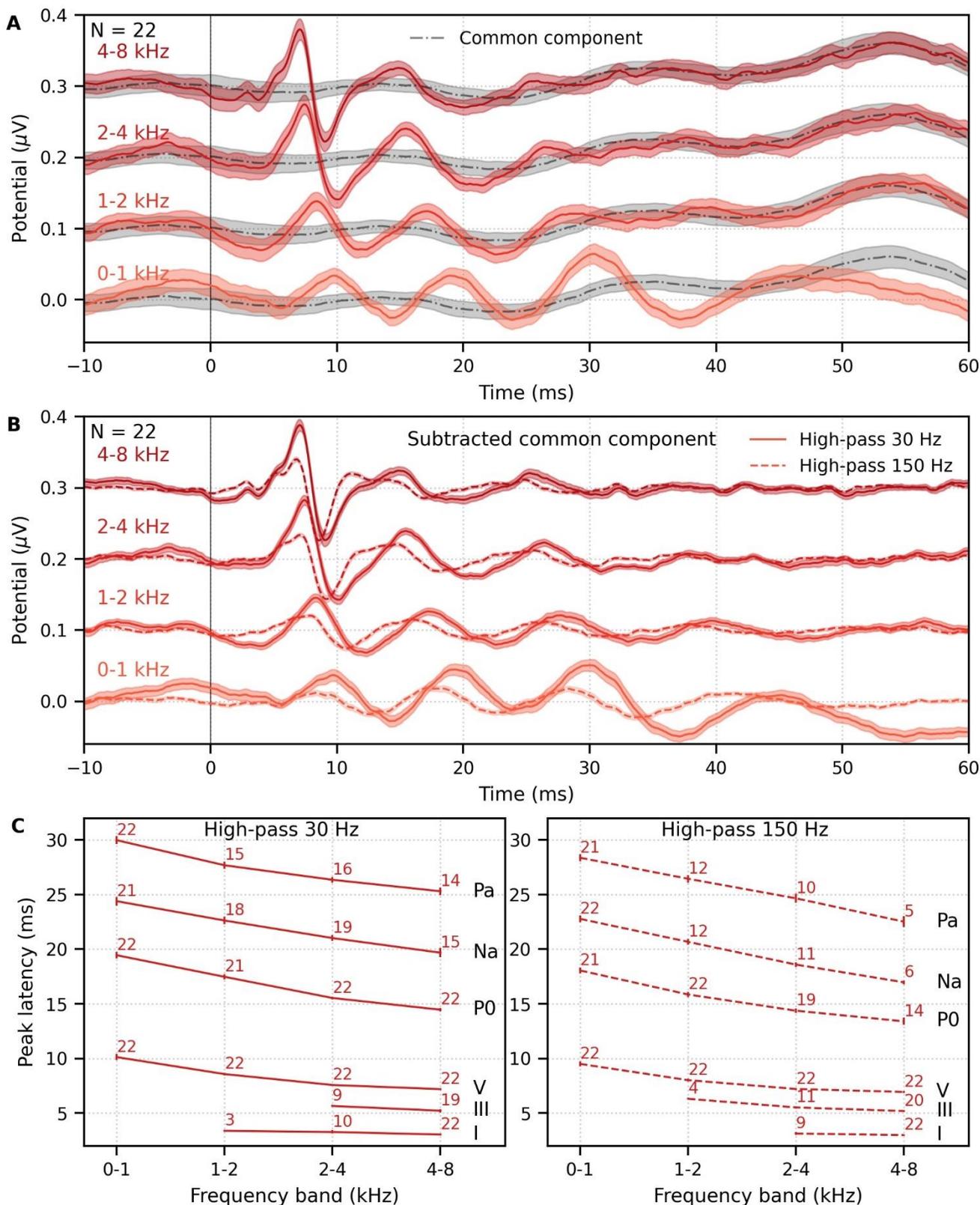


Figure 4. Comparison of responses to ~43 minutes of male-narrated multiband peaky speech. (A) Average waveforms across subjects (areas show ± 1 SEM) are shown for each band (colored solid lines) and for the common component (dot-dash gray line, same waveform replicated as a reference for each band), which was calculated using 6 false pulse trains. (B) The common component was subtracted from each band's response to give the frequency-specific waveforms (areas show ± 1 SEM), which are shown with high-pass filtering at 30 Hz (solid lines) and 150 Hz (dashed lines). (C) Mean \pm SEM peak latencies for each wave decreased with increasing band frequency. Numbers of subjects with an identifiable wave are given for each wave and band.

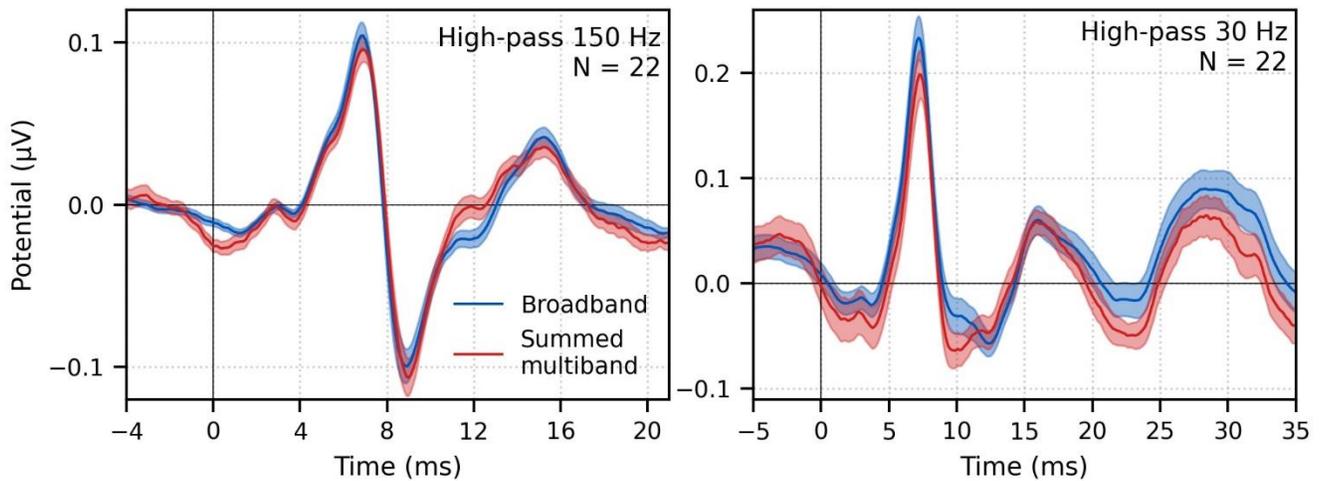


Figure 5. Comparison of responses to ~43 minutes of male-narrated peaky speech in the same subjects. Average waveforms across subjects (areas show ± 1 SEM) are shown for broadband peaky speech (blue) and for the summed frequency-specific responses to multiband peaky speech with the common component added (red), high-pass filtered at 150 Hz (left) and 30 Hz (right). Regressors in the deconvolution were pulse trains.

252

253 *Frequency-specific responses also differ by narrator*

254 We also investigated the effects of narrator on multiband peaky speech by deriving responses to 32
255 minutes (30 epochs of 64 s each) each of male- and female-narrated multiband peaky speech in the same
256 11 subjects. As with broadband peaky speech, responses to both narrators showed similar morphology,
257 but the responses were smaller and the MLR waves more variable for the female than male narrator
258 (Figure 6A). Figure 6B shows the male-female correlation coefficients for responses between 0–40 ms with
259 a high-pass filter of 30 Hz and between 0–15 ms with a high-pass filter of 150 Hz. The median (interquartile
260 range) male-female correlation coefficients were better for higher frequency bands, ranging from 0.20
261 (–0.04–0.29) for the 1–2 kHz band to 0.38 (0.27–0.48) for the 4–8 kHz band for MLR lags (Figure 6B, left
262 panel), and from 0.57 (0.29–0.68) for the 0–1 kHz band to 0.81 (0.70–0.84) for the 4–8 kHz band for ABR
263 lags (Figure 6B, right panel). These male-female correlation coefficients were significantly weaker than
264 those of the same EEG split into even and odd trials for all but the 2–4 kHz frequency band when
265 responses were high-pass filtered at 30 Hz and correlated across 0–40 ms lags (2–4 kHz: $W_{(10)} = 21.0$, $p =$
266 0.320; other bands: $W_{(10)} \leq 8.0$, $p \leq 0.024$), but were similar to the even/odd trials for responses from all
267 frequency bands high-pass filtered at 150 Hz ($W_{(10)} \geq 13.0$, $p \geq 0.083$). These results indicate that the
268 specific narrator can affect the robustness of frequency-specific responses, particularly for the MLR waves.

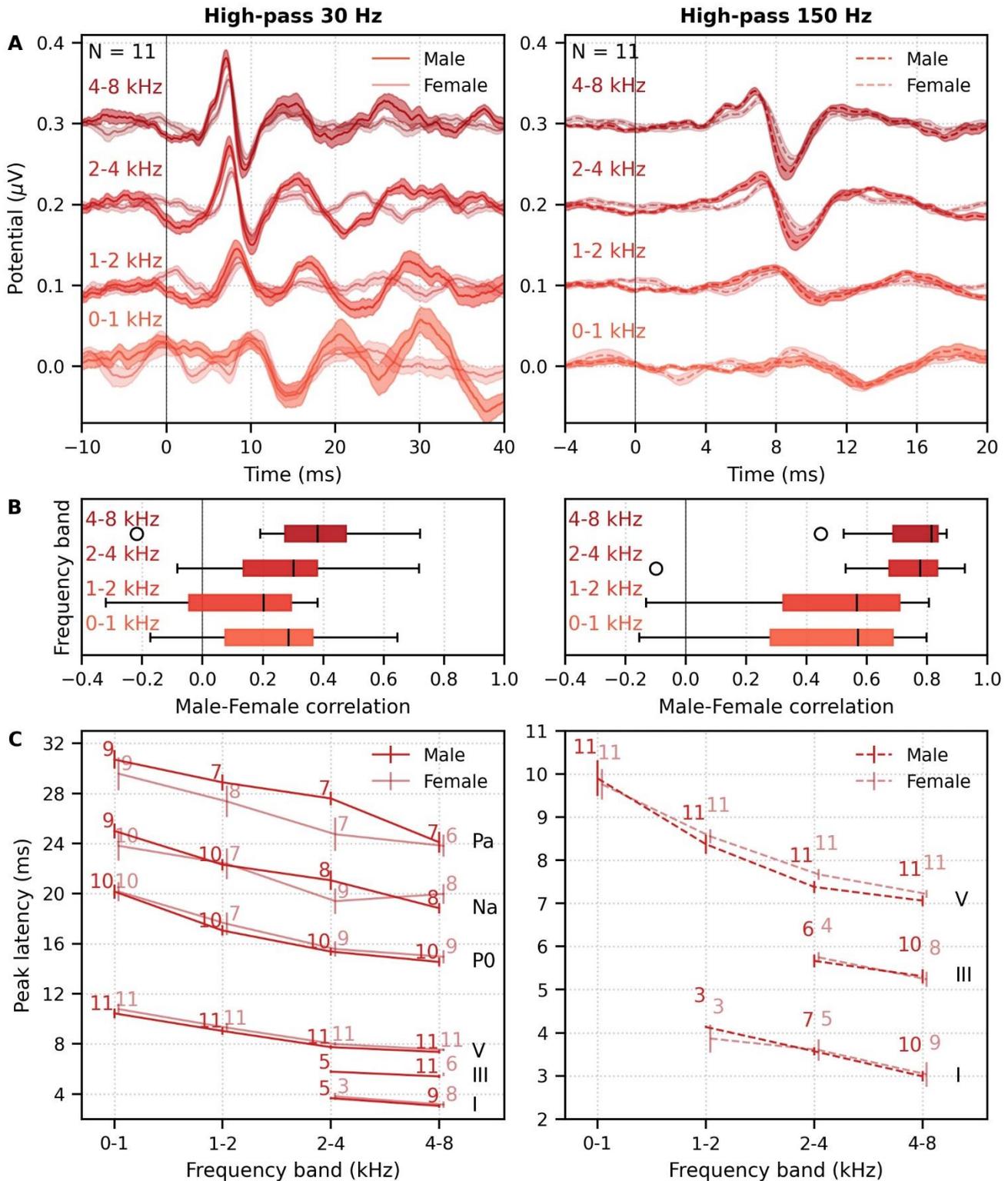


Figure 6. Comparison of responses to 32 minutes each of male- and female-narrated re-synthesized multiband peaky speech. (A) Average frequency-specific waveforms across subjects (areas show ± 1 SEM; common component removed) are shown for each band in response to male- (dark red lines) and female-narrated (light red lines) speech. Responses were high-pass filtered at 30 Hz (left) and 150 Hz (right) to highlight the MLR and ABR respectively. (B) Correlation coefficients between responses evoked by male- and female-narrated multiband peaky speech during ABR/MLR (left) and ABR (right) time lags for each frequency band. Black lines denote the median. (C) Mean \pm SEM peak latencies for male- (dark) and female- (light) narrated speech for each wave decreased with increasing frequency band. Numbers of subjects with an identifiable wave are given for each wave, band and narrator. Lines are given a slight horizontal offset to make the error bars easier to see.

269

270

271

As expected from the grand average waveforms and male-female correlations, there were fewer subjects who had identifiable waves across frequency bands for the female- than male-narrated speech. These

272 numbers are shown in Figure 6C, along with the mean \pm SEM peak latencies for each wave, frequency
273 band and narrator. ICC3 ≥ 0.93 indicated good agreement in peak latency choices (the two lowest 95%
274 confidence intervals were 0.83–0.97 for wave III and 0.98–0.99 for P_a, both with a high-pass filter cutoff of
275 30 Hz). Again, there were few numbers of subjects with identifiable waves I and III for the lower frequency
276 bands. Therefore, the mixed effects model was completed for waves V, P₀, N_a, and P_a of responses in the 4
277 frequency bands that were high-pass filtered at 30 Hz. The model included fixed effects of narrator, wave,
278 linear and quadratic terms for frequency band, the interaction between narrator and wave, and the
279 interactions between wave and frequency band terms, as well as a random intercept per subject and
280 random frequency band terms for per subject. Details of the model are described in Supplemental Table 2.
281 For those subjects with identifiable waves, peak latencies shown in Figure 6C differed by wave ($p < 0.001$
282 for effects of each wave on the intercept), and latency decreased with increasing frequency band ($p <$
283 0.001 for the linear term, slope). This change with frequency was greater (i.e., steeper slope) for each MLR
284 wave compared to wave V ($p < 0.013$ for all interactions between wave and the linear term for frequency
285 band). There was no change in slope with band ($p = 0.190$ for the quadratic term, $p > 0.318$ for the
286 interactions between wave and the quadratic term). There was also no main effect of narrator on peak
287 latencies (narrator $p = 0.481$), except that the latencies for P_a were faster for the female than male narrator
288 (P_a–narrator interaction $p = 0.003$; other wave–narrator interactions $p > 0.195$). Therefore, as with
289 broadband peaky speech, frequency-specific peaky responses were more robust with the male narrator,
290 but unlike the broadband responses, the frequency-specific responses did not peak earlier for a narrator
291 with a lower fundamental frequency.

292

293 *Frequency-specific responses can be measured simultaneously in each ear (dichotically)*

294 The focus so far has been on scientific use of peaky speech but there are also potential clinical
295 applications. Frequency-specific ABRs to tone pips are traditionally used to assess hearing function in each
296 ear across octave bands with center frequencies of 500–8000 Hz. Applying the same principles to generate
297 multiband peaky speech, we investigated whether ear-specific responses could be evoked across 5
298 standard, clinically-relevant (audiological) frequency bands using dichotic multiband speech. For peaky
299 dichotic (stereo) audiological multiband speech we created 10 independent pulse trains, 2 for each ear in
300 each of the 5 frequency bands (see Multiband peaky speech and Band filters in Methods).

301 We recorded responses to 64 minutes (60 epochs of 64 s each) each of male- and female-narrated
302 dichotic multiband peaky speech in 11 subjects. The frequency-specific (i.e., common component-
303 subtracted) group average waveforms for each ear and frequency band are shown in Figure 7A. The ten
304 waveforms were small, especially for female-narrated speech, but a wave V was identifiable for both
305 narrators. MLR waves were not clearly identifiable for responses to female-narrated speech. Therefore,
306 correlations between responses were performed for ABR lags between 0–15 ms. As shown in Figure 7B,
307 the median (interquartile range) left-right ear correlation coefficients (averaged across narrators) ranged
308 from 0.17 (–0.12–0.49) for the 0.5 kHz band to 0.63 (0.33–0.86) for the 8 kHz band. Male-female
309 correlation coefficients (averaged across ear) ranged from 0.08 (–0.25–0.22) for the 0.5 kHz band to 0.70
310 (0.26–0.80) for the 4 kHz band. Although the female-narrated responses were smaller than the male-
311 narrated responses, these male-female coefficients did not significantly differ from the left-right ear
312 coefficients ($W_{(10)} \geq 20.0$, $p \geq 0.278$), or from correlations of same EEG split into even-odd trials and
313 averaged across ear ($W_{(10)} \geq 20.0$, $p \geq 0.278$), likely reflecting the variability in such small responses.

314 Figure 7C shows the mean \pm SEM peak latencies of wave V for each ear and frequency band for the male-
315 and female-narrated dichotic multiband peaky speech. The ICC3 for wave V was 0.98 (95% confidence
316 interval 0.98–0.99), indicating reliable peak latency choices. The nonlinear change in wave V latency with
317 frequency was modeled using mixed effects regression with fixed effects of narrator, ear, linear and
318 quadratic terms for frequency band, and the interactions between narrator and frequency band terms.
319 Random effects included an intercept and both frequency band terms for each subject. Details of the model
320 are described in Supplemental Table 3. Wave V latency was significantly longer for female- than male-
321 narrated multiband peaky speech in the 0.5 kHz band (narrator effect on the intercept, $p = 0.001$),
322 decreased at a steeper rate with frequency band (interaction between narrator and linear frequency band
323 term $p < 0.001$), and had a significantly different rate of change with frequency (interaction between
324 narrator and quadratic frequency band term, $p < 0.001$). Overall, latency did not differ between ears ($p =$
325 0.116). Taken together, these results confirm that, while small in amplitude, frequency-specific responses

326
327

can be elicited in both ears across 5 different frequency bands and show characteristic latency changes across the different frequency bands.

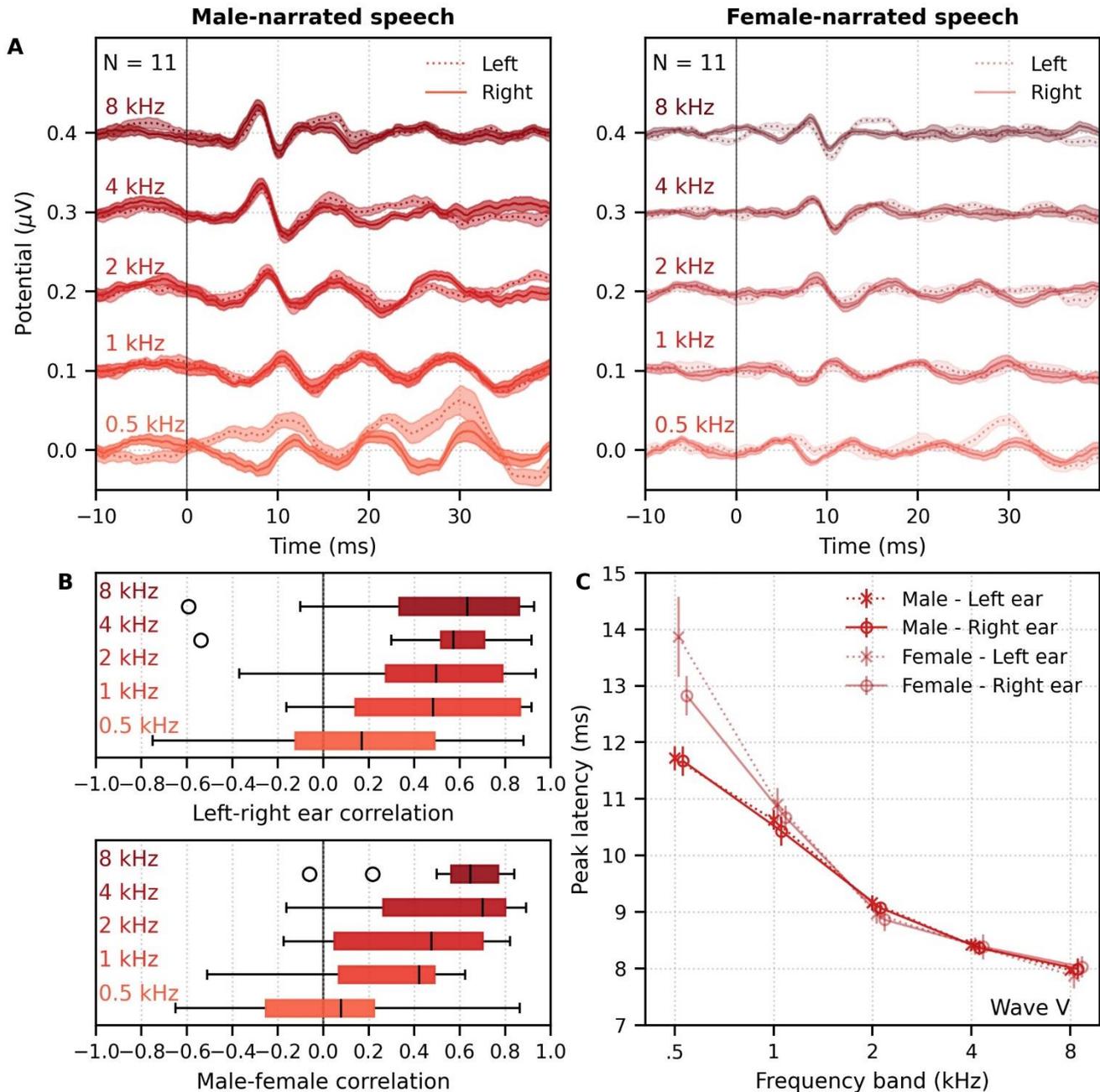


Figure 7. Comparison of responses to ~60 minutes each of male- and female-narrated dichotic multiband peaky speech with standard audiological frequency bands. (A) Average frequency-specific waveforms across subjects (areas show ± 1 SEM; common component removed) are shown for each band for the left ear (dotted lines) and right ear (solid lines). Responses were high-pass filtered at 30 Hz. (B) Left-right ear correlation coefficients (top, averaged across gender) and male-female correlation coefficients (bottom, averaged across ear) during ABR time lags (0–15 ms) for each frequency band. Black lines denote the median. (C) Mean \pm SEM wave V latencies for male- (dark red) and female-narrated (light red) speech for the left (dotted line, cross symbol) and right ear (solid line, circle symbol) decreased with increasing frequency band. Lines are given a slight horizontal offset to make the error bars easier to see.

328

329

Responses are obtained quickly for male-narrated broadband peaky speech but not multiband speech

330

331

332

333

334

Having demonstrated that peaky broadband and multiband speech provides canonical waveforms with characteristic changes in latency with frequency, we next evaluated the acquisition time required for waveforms to reach a decent SNR. We chose 0 dB SNR based on visual assessment of when waveforms were easily inspected and based on what we have done previously (Maddox and Lee, 2018; Polonenko and Maddox, 2019). SNR was calculated by comparing the variance in the MLR time interval 0–30 ms (for

335 responses high-pass filtered at 30 Hz) or ABR time interval 0–15 ms (for responses high-pass filtered at
336 150 Hz) to the variance in the pre-stimulus noise interval –480 to –20 ms (see Response SNR calculation
337 in Methods for details).

338 Figure 8 shows the cumulative proportion of subjects who had responses with ≥ 0 dB SNR to unaltered and
339 broadband peaky speech as a function of recording time. Acquisition times for 22 subjects were similar for
340 responses to both unaltered and broadband peaky male-narrated speech, with 0 dB SNR achieved by 8
341 minutes in 50% of subjects and by 18 and 20 minutes respectively in 100% of subjects. This time reduced
342 to 2 and 5 minutes for 50% and 100 % of subjects respectively for broadband peaky responses high-pass
343 filtered at 150 Hz to highlight the ABR (0–15 ms interval). These times for male-narrated broadband peaky
344 speech were confirmed in our second cohort of 11 subjects, who also all achieved 0 dB SNR within 26
345 minutes for the 0–30 ms MLR interval (10 / 11 subjects in 18 minutes; 50% by 10 minutes) and 4 minutes
346 for the 0–15 ms ABR interval (50% by 2 minutes). However, acquisition times were at least 3.6 times – but
347 up to over 10 times – longer for female-narrated broadband peaky speech, with 50% of subjects achieving
348 0 dB SNR by 36 minutes for the MLR interval and 8 minutes for ABR interval. In contrast to male-narrated
349 speech, not all subjects achieved this threshold for female-narrated speech by the end of the 32-minute
350 recording (45% and 63% for the MLR and ABR intervals respectively). Taken together, these acquisition
351 times confirm that responses with useful SNRs can be measured quickly for male-narrated broadband
352 peaky speech but longer recording sessions are necessary for narrators with higher fundamental
353 frequencies.

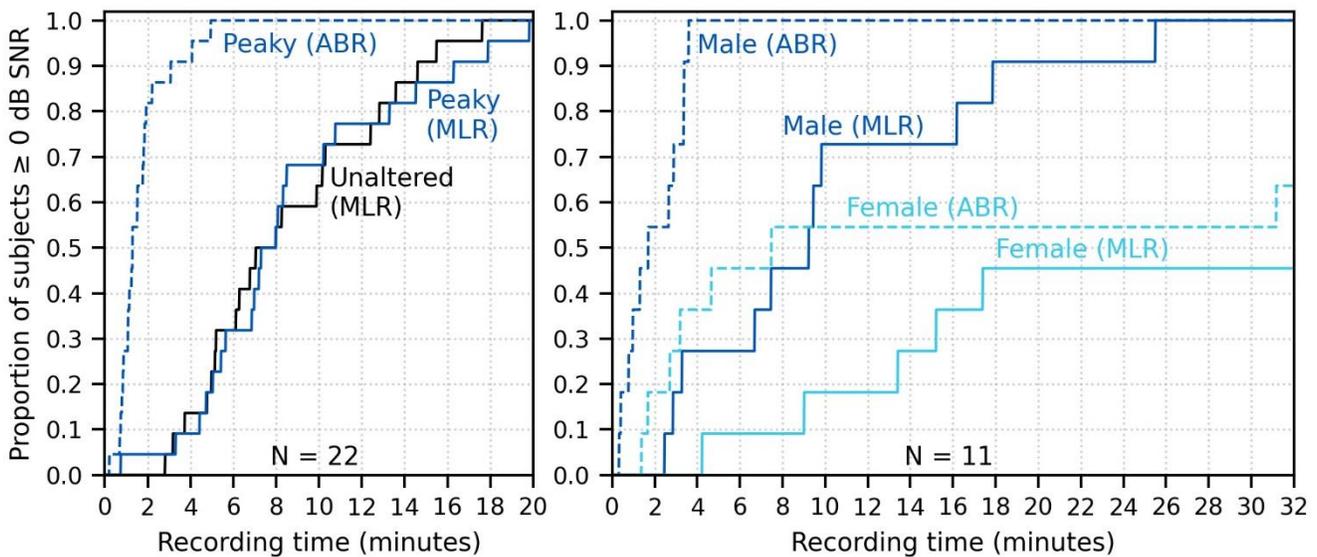


Figure 8. Cumulative proportion of subjects who have responses with ≥ 0 dB SNR as a function of recording time. Time required for unaltered (black) and broadband peaky speech (dark blue) of a male narrator are shown for 22 subjects in the left plot, and for male (dark blue) and female (light blue) broadband peaky speech is shown for 11 subjects in the right plot. Solid lines denote SNRs calculated using variance of the signal high-pass filtered at 30 Hz over the ABR/MLR interval 0–30 ms, and dashed lines denote SNR variances calculated on signals high-pass filtered at 150 Hz over the ABR interval 0–15 ms. Noise variance was calculated in the pre-stimulus interval –480 to –20 ms.

354

355 The longer recording times necessary for a female narrator became more pronounced for the multiband
356 peaky speech. Figure 9A shows the cumulative density function for responses high-pass filtered at 150 Hz
357 and the SNR estimated over the ABR interval. Most subjects (72%) had frequency-specific responses
358 (common component subtracted) with ≥ 0 dB SNR for all 4 frequency bands by the end of the 32-minute
359 recording for the male-narrated speech, but this was achieved in only 45% of subjects for the female-
360 narrated speech. Multiband peaky speech required significantly longer recording times than broadband
361 peaky speech, with 50% of subjects achieving 0 dB SNR by 22 minutes compared to 2 minutes for the
362 male-narrated responses across the ABR 0–15 ms interval and 23 minutes compared to 5 minutes for the
363 MLR 0–30 ms interval. For the MLR interval, 72% and 18% of subjects had reached 0 dB SNR by the 32-
364 minute recording for male- and female-narrated speech respectively. Even more time was required for
365 dichotic multiband speech, which was comprised of a larger number of frequency bands (Figure 9B). All 10
366 audiological band responses achieved ≥ 0 dB SNR in 36% of ears (8 / 22 ears from 11 subjects) by 64
367 minutes for male-narrated speech and in 22% of ears (5 / 22 ears) for female-narrated speech. The smaller

368 and broader responses in the low frequency bands limited this testing time – for male-narrated speech, at
369 least 90% of subjects had 2–4 and 4–8 kHz responses (diotic 4-band speech) with ≥ 0 dB SNR in 15
370 minutes, 70% of subjects had 6 frequency-specific responses (2, 4, 8 kHz bands in both ears for dichotic
371 speech) by the end of the recording, and 50% of subjects had the 6 higher frequency-specific responses
372 within 40 minutes. These significant recording times suggest that deriving multiple frequency-specific
373 responses will require at least more than 30 minutes per condition for < 5 bands, and more than an hour
374 session for one condition of peaky multiband speech with 10 bands.

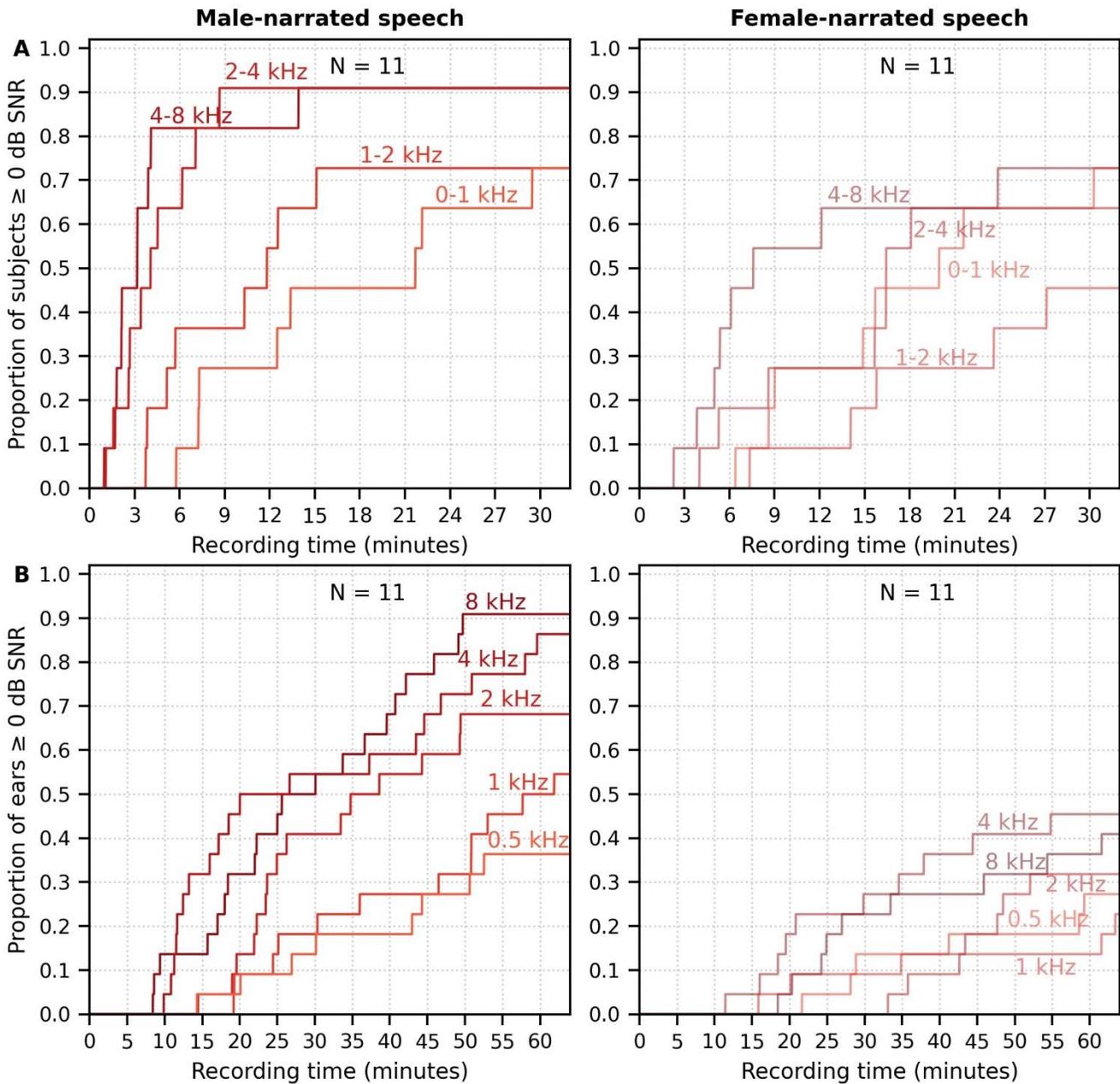


Figure 9. Cumulative proportion of subjects who have frequency-specific responses (common component subtracted) with ≥ 0 dB SNR as a function of recording time. Acquisition time was faster for male (left) than female (right) narrated multiband peaky speech with (A) 4 frequency bands presented diotically, and with (B) 5 frequency bands presented dichotically (total of 10 responses, 5 bands in each ear). SNR was calculated by comparing variance of signals high-pass filtered at 150 Hz across the ABR interval of 0–15 ms to variance of noise in the pre-stimulus interval –480 to –20 ms.

375
376

377 DISCUSSION

378 The major goal of this work was to develop a method to investigate early stages of naturalistic speech
379 processing. We re-synthesized continuous speech taken from audio books so that the phases of all
380 harmonics aligned at each glottal pulse during voiced segments, thereby making speech as impulse-like
381 (peaky) as possible to drive the auditory brainstem. Then we used the glottal pulse trains as the regressor
382 in deconvolution to derive the responses. Indeed, comparing waveforms to broadband peaky and unaltered
383 speech validated the superior ability of peaky speech to evoke additional waves of the canonical ABR and
384 MLR, reflecting neural activity from multiple subcortical structures. Robust ABR and MLR responses were
385 recorded in less than 5 and 20 minutes respectively for all subjects, with half of the subjects exhibiting a
386 strong ABR within 2 minute and MLR within 8 minutes. Longer recording times were required for the
387 smaller responses generated by a narrator with a higher fundamental frequency. We also demonstrated
388 the flexibility of this stimulus paradigm by simultaneously recording up to 10 frequency-specific responses
389 to multiband peaky speech that was presented either diotically or dichotically, although these responses
390 required much longer recording times. Taken together, our results show that peaky speech effectively
391 yields responses from distinct subcortical structures and from different frequency bands, paving the way for
392 new investigations of speech processing and new tools for clinical application.

393 For the purpose of investigating responses from different subcortical structures, we accomplished our goal
394 of creating a stimulus paradigm that overcame some of the limitations of current methods using natural
395 speech. Methods that do not use re-synthesized impulse-like speech generate responses characterized by
396 a broad peak between 6–9 ms (Forte et al., 2017; Maddox and Lee, 2018), with contributions predominantly
397 from the inferior colliculus (Saiz-Alia and Reichenbach, 2020). In contrast, for the majority of our subjects,
398 peaky speech evoked responses with canonical morphology comprised of waves I, III, V, P₀, N_a, P_a (Figure
399 1), reflecting neural activity from distinct stages of the auditory system from the auditory nerve to thalamus
400 and primary auditory cortex (e.g., Picton et al., 1974). Presence of these additional waves allows for new
401 investigations into the contributions of each of these neural generators to speech processing while using a
402 continuously dynamic and ecologically salient stimulus.

403 The same ABR waves evoked here were also evoked by a method using embedded chirps intermixed within
404 alternating octave bands of speech, particularly if presented monaurally over headphones instead of in free
405 field (Backer et al., 2019; Miller et al., 2017). Chirps are transients that compensate for the cochlear traveling
406 delay wave by introducing different phases across frequency, leading to a more synchronized response
407 across the cochlea and a larger brainstem response than for clicks (Dau et al., 2000; Elberling and Don,
408 2008; Shore and Nuttall, 1985). The responses to embedded chirps elicited waves with larger mean
409 amplitude than those to our broadband peaky speech (~0.4 versus ~0.2 μ V, respectively), although a similar
410 proportion of subjects had identifiable waves and several other factors may contribute to amplitude
411 differences. For example, higher click rates (e.g., Burkard et al., 1990; Burkard and Hecox, 1983; Chiappa
412 et al., 1979; Don et al., 1977; Jiang et al., 2009) and higher fundamental frequencies (Maddox and Lee,
413 2018; Saiz-Alía et al., 2019; Saiz-Alia and Reichenbach, 2020) reduce the brainstem response amplitude,
414 and dynamic changes in rate may create interactions across neural populations that lead to smaller
415 amplitudes. Our stimuli kept the dynamic changes in pitch across all frequencies (instead of alternate octave
416 bands of chirps and speech) and created impulses at every glottal pulse, with an average pitch of ~115 Hz
417 and ~198 Hz for the male and female narrators respectively. These presentation rates were much higher and
418 more variable than the flat 42 Hz rate at which the embedded chirps were presented (pitch flattened to 82
419 Hz and chirps presented every other glottal pulse). We could evaluate whether chirps would improve
420 response amplitude to our dynamic peaky speech by simply all-pass filtering the re-synthesized voiced
421 segments by convolving with a chirp prior to mixing the re-synthesized parts with the unvoiced segments.
422 While maintaining the amplitude spectrum of speech, the harmonics would then have the different phases
423 associated with chirps at each glottal pulse instead of all phases set to 0. Regardless, our peaky speech
424 generated robust canonical responses with good SNR while maintaining a natural-sounding, if very slightly
425 “buzzy”, quality to the speech. Overall, continuous speech re-synthesized to contain impulse-like
426 characteristics is an effective way to elicit responses that distinguish contributions from different subcortical
427 structures.

428 The latencies of the component waves of the responses to peaky speech are consistent with activity arising
429 from known subcortical structures. The inter-wave latencies between I-III, III-V and I-V fall within the
430 expected range for brainstem responses elicited by transients at 50–60 dB sensation level (SL) and 50–

431 100 Hz rates (Burkard and Hecox, 1983; Chiappa et al., 1979; Don et al., 1977), suggesting the
432 transmission times between auditory nerve, cochlear nucleus and rostral brainstem remain similar for
433 speech stimuli. However, these speech-evoked waves peak at later absolute latencies than responses to
434 transient stimuli at 60 dB SL and 90–100 Hz, but at latencies more similar to those presented at 50 dB SL
435 or 50 dB nHL in the presence of some masking noise (Backer et al., 2019; Burkard and Hecox, 1983;
436 Chiappa et al., 1979; Don et al., 1977; Maddox and Lee, 2018; Miller et al., 2017). There are a couple of
437 reasons why the speech-evoked latencies may be later. First, our level of 60 dB sound pressure level
438 (SPL) may be more similar to click levels of 50 dB SL. Second, although spectra of both speech and
439 transients are broad, clicks, chirps and even our previous speech stimuli (which was high-pass filtered at 1
440 kHz; Maddox and Lee, 2018) have relatively greater high-frequency energy than the unaltered and peaky
441 broadband speech used in the present work. Neurons with higher characteristic frequencies respond
442 earlier due to their basal cochlear location, and contribute relatively more to brainstem responses (e.g.,
443 Abdala and Folsom, 1995), leading to quicker latencies for stimuli that have greater high frequency energy.
444 Also consistent with having greater lower frequency energy, our unaltered and peaky speech responses
445 were later than the response from the same speech segments that were high-pass filtered at 1 kHz
446 (Maddox and Lee, 2018). In fact, the ABR to broadband peaky speech bore a close resemblance to the
447 summation of each frequency-specific response and the common component to peaky multiband speech
448 (Figure 5), with peak wave latencies representing the relative contribution of each frequency band. Third,
449 higher stimulation rates prolong latencies due to neural adaptation, and the 115–198 Hz average
450 fundamental frequencies of our speech were much higher than the 41 Hz embedded chirps and 50–100 Hz
451 click rates (e.g., Burkard et al., 1990; Burkard and Hecox, 1983; Chiappa et al., 1979; Don et al., 1977;
452 Jiang et al., 2009). The effect of stimulation rate was also demonstrated by the later ABR wave I, III, and V
453 peak latencies for the female narrator with the higher average fundamental frequency of 198 Hz (Figure
454 3A&C). Fourth, continuous speech is dynamic with much greater variability in pitch than the variability in
455 presentation rate of typical period stimulation with clicks, chirps or short syllables. Across all our 64 second
456 speech segments, the average standard deviation in pitch was 28.5 and 50.6 Hz for the male and female
457 narrators respectively, with pitch varying from a minimum of 61 and 91 Hz to a maximum of 288 and 456
458 Hz, respectively. This variability in rate over a continuous stimulus may create interactions across different
459 neural populations that may delay responses compared to the relatively regular transient stimuli. Therefore,
460 the differing characteristics of typical periodic transients (such as clicks and chirps) and continuous speech
461 may give rise to differences in brainstem responses, even though they share canonical waveforms arising
462 from similar contributing subcortical structures.

463 Latency of the peaky speech-evoked response also differed from the non-standard, broad responses to
464 unaltered speech. However, latencies from these waveforms are difficult to compare due to the differing
465 morphology and the different analyses that were used to derive the responses. Evidence for the effect of
466 analysis comes from the fact that the same EEG collected in response to peaky speech could be regressed
467 with pulse trains to give canonical ABRs (Figures 1, 2), or regressed with the half-wave rectified peaky
468 speech to give the different, broad waveform (Supplemental Figure 1). Furthermore, non-peaky continuous
469 speech stimuli with similar ranges of fundamental frequencies (between 100–300 Hz) evoke non-standard,
470 broad brainstem responses that also differ in morphology and latency depending on whether the EEG is
471 analyzed by deconvolution with the half-wave rectified speech (Figure 2, Maddox and Lee, 2018) or
472 complex cross-correlation with the fundamental frequency waveform (Forte et al., 2017). Therefore, again,
473 even though the inferior colliculus and lateral lemniscus may contribute to generating these different
474 responses (Møller and Jannetta, 1983; Saiz-Alia and Reichenbach, 2020; Starr and Hamilton, 1976), the
475 morphology and latency may differ (sometimes substantially) depending on the analysis technique used.

476 In addition to evoking canonical brainstem responses, peaky speech can be exploited for other traditional
477 uses of ABR, such as investigating subcortical responses across different frequencies. Frequency-specific
478 responses were measurable to two different types of multiband peaky speech: 4 frequency bands
479 presented diotically (Figures 4, 6), and 5 frequency bands presented dichotically (Figure 7). Peak wave
480 latencies of these responses decreased with increasing band frequency in a similar way to responses
481 evoked by tone pips (Gorga et al., 1988; Rasetshwane et al., 2013), thereby representing activity evoked
482 from different areas across the cochlea. Interestingly, the frequency-specific responses across frequency
483 band were similar in amplitude (Figures 4, 6, 7) even though the relative energy of each band decreased
484 with increasing frequency, resulting in a ~30 dB difference between the lowest and highest frequency
485 bands (Figure 12). A greater response elicited by higher frequency bands is consistent with the relatively

486 greater contribution of neurons with higher characteristic frequencies to ABRs (Abdala and Folsom, 1995),
487 as well as the need for higher levels to elicit low frequency responses to tone pips that are close to
488 threshold (Gorga et al., 2006, 1993; Hyde, 2008; Stapells and Oates, 1997). Also, canonical waveforms
489 were derived in the higher frequency bands of diotically presented speech, with waves I and III identifiable
490 in most subjects. Measuring waves I, III, and V of high frequency responses may have applications to
491 studying cochlear synaptopathy (Lieberman et al., 2016) using naturalistic speech in humans. Another
492 exciting application is the evaluation of supra-threshold hearing across frequency in toddlers and
493 individuals who do not provide reliable behavioral responses, as they may be more responsive to sitting for
494 longer periods of time while listening to a narrated story than to a series of tone pips. An extension of this
495 assessment would be to evaluate neural speech processing in the context of hearing loss, as well as
496 rehabilitation strategies such as hearing aids and cochlear implants. Therefore, the ability of peaky speech
497 to yield both canonical waveforms and frequency-specific responses makes this paradigm a flexible
498 method that assesses speech processing in new ways.

499 Having established that peaky speech is a flexible stimulus for investigating different aspects of speech
500 processing, there are several practical considerations for using the peaky speech paradigm. First, filtering
501 should be performed carefully. As recommended in Maddox and Lee (2018), causal filters – which have
502 impulse responses with non-zero values at positive lags – should be used to ensure cortical activity at later
503 peak latencies does not spuriously influence earlier peaks corresponding to subcortical origins. Applying
504 less aggressive, low-order filters (i.e., broadband with shallow roll-offs) will help reduce the effects of
505 causal filtering on delaying response latency. The choice of high-pass cutoff will also affect the response
506 amplitude and morphology. After evaluating several orders and cutoffs to the high-pass filters, we
507 determined that early waves of the peaky broadband ABRs were best visualized with a 150 Hz cutoff,
508 whereas a lower cutoff frequency of 30 Hz was necessary to view the ABR and MLR of the broadband
509 responses. For multiband responses, the 150 Hz high-pass filter significantly reduced the response but
510 also decreased the low-frequency noise in the pre-stimulus interval. For the 4-band multiband peaky
511 speech the 150 Hz and 30 Hz filters provided similar acquisition times for 0 dB SNR, but better SNRs were
512 obtained quicker with 150 Hz filtering for the 10-band multiband peaky speech.

513 Second, the choice of narrator impacts the responses to both broadband and multiband peaky speech.
514 Although overall morphology was similar, the male-narrated responses were larger, contained more clearly
515 identifiable component waves in a greater proportion of subjects, and achieved a 0 dB SNR at least 3.6 to
516 over 10 times faster than those evoked by a female narrator. These differences likely stemmed from the
517 ~77 Hz difference in average pitch, as higher stimulation rates evoke smaller responses due to adaptation
518 and refractoriness (e.g., Burkard et al., 1990; Burkard and Hecox, 1983; Chiappa et al., 1979; Don et al.,
519 1977; Jiang et al., 2009). Indeed, a 50 Hz change in fundamental frequency yields a 24% reduction in the
520 modelled auditory brainstem response that was derived as the complex cross-correlation with the
521 fundamental frequency (Saiz-Alia and Reichenbach, 2020). The narrator differences exhibited in the
522 present study may be larger than those in other studies with continuous speech (Forte et al., 2017; Maddox
523 and Lee, 2018; Saiz-Alia et al., 2019) as a result of the different regressors. These response differences do
524 not preclude using narrators with higher fundamental frequencies in future studies, but the time required for
525 usable responses from each narrator must be considered when planning experiments, and caution taken
526 when interpreting comparisons between conditions with differing narrators. The strongest results will come
527 from comparing responses to the same narrator (or even the same speech recordings) under different
528 experimental conditions.

529 Third, the necessary recording time depends on the chosen SNR threshold, experimental demands, and
530 stimulus. We chose a threshold SNR of 0 dB based on when waveforms became clearly identifiable, but of
531 course a different threshold would change our recording time estimates (though, notably, not the ratios
532 between them). With this SNR threshold, acquisition times were quick enough for broadband peaky
533 responses to allow multiple conditions in a reasonable recording session. With male-narrated broadband
534 peaky speech, all subjects achieved 0 dB SNR ABRs in < 5 minutes and MLRs in < 20 minutes, thereby
535 affording between 3 and 12 conditions in an hour recording session. These recording times are
536 comparable, if not faster, than the 8 minutes for the broad response to unaltered speech, 6–12 minutes for
537 the chirp-embedded speech (Backer et al., 2019), ~10 minutes for the broad complex-cross correlation
538 response to the fundamental waveform (Forte et al., 2017), and 33 minutes for the broad response to high-
539 passed continuous speech (Maddox and Lee, 2018). However, using a narrator with a higher fundamental
540 frequency could increase testing time by 3- to over 10-fold. In this experiment, at most 2 conditions per

541 hour could be tested with the female-narrated broadband peaky speech. Furthermore, longer testing times
542 are likely needed, even for male-narrated speech, in order to reliably compare differences in the smaller
543 amplitude component waves I and III of the ABR. Our 30- to 40-minute recording sessions provided robust
544 responses with very good SNRs to evaluate the earlier ABR waves, but this long may be unnecessary. The
545 cumulative density functions in Figure 8 suggest that between 12 to 20 minutes should constitute ample
546 time to generate comparable responses with highly positive SNRs. Unlike broadband peaky speech, the
547 testing times required for all frequency-specific responses to reach 0 dB SNR were significantly longer,
548 making only 1 condition feasible within a recording session. At least 30 minutes was necessary for the
549 dichotically presented multiband peaky speech with 4 frequency bands, but based on extrapolated testing
550 times, about 56–88 minutes is required for 90% of subjects to achieve this threshold for all 4 bands. For
551 dichotically presented multiband peaky speech with 5 frequency bands (for a total of 10 frequency-specific
552 waveforms), only 36% the responses achieved 0 dB SNR within an hour. Extrapolated testing times
553 suggest that over 2 hours is required for at least 75% of subjects, limiting the feasibility or utility of
554 multiband peaky speech with several frequency bands.

555 Fourth, as mentioned above, the number of frequency bands incorporated into multiband peaky speech
556 decreases SNR and increases testing time. Although it is possible to simultaneously record up to 10
557 frequency-specific responses, the significant time required to obtain decent SNRs reduces the feasibility of
558 testing multiple conditions or having recording sessions lasting less than 1–2 hours. However, pursuing
559 shorter testing times with multiband peaky speech is possible. Depending on the experimental question,
560 different multiband options could be considered. For male-narrated speech, the 2–4 and 4–8 kHz
561 responses had good SNRs and exhibited waves I, III, and V within 15 minutes for 90% of subjects.
562 Therefore, if researchers were more interested in comparing responses in these higher frequency bands,
563 they could stop recording once these bands reach threshold but before the lower frequency bands reach
564 criterion (i.e., within 15 minutes). Alternatively, the lower frequencies could be combined into a single
565 broader band in order to reduce the total number of bands, or the intensity could be increased to evoke
566 responses with larger amplitudes. Therefore, different band and parameter considerations could reduce
567 testing time and improve the feasibility, and thus utility, of multiband peaky speech.

568 Fifth, and finally, a major advantage of deconvolution analysis is that the analysis window for the response
569 can be extended arbitrarily in either direction to include a broader range of latencies (Maddox and Lee,
570 2018). Extending the pre-stimulus window leftward provides a better estimate of the SNR, and extending
571 the window rightward allows parts of the response that come after the ABR and MLR to be analyzed as
572 well, which are driven by the cortex. These later responses can be evaluated in response to broadband
573 peaky speech, but as shown in Figures 6 and 7, only ABR and early MLR waves are present in the
574 frequency-specific responses. The same broadband peaky speech data from Figure 3 are displayed with
575 an extended time window in Figure 10, which shows component waves of the ABR, MLR and late latency
576 responses (LLR). Thus, this method allows us to simultaneously investigate speech processing ranging
577 from the earliest level of the auditory nerve all the way through the cortex without requiring extra recording
578 time. Usually the LLR is larger than the ABR/MLR, but our subjects were encouraged to relax and rest,
579 yielding a passive LLR response. Awake and attentive subjects may improve the LLR; however, other
580 studies that present continuous speech to attentive subjects also report smaller and different LLR (Backer
581 et al., 2019; Maddox and Lee, 2018), possibly from cortical adaptation to a continuous stimulus. Here we
582 used a simple 2-channel montage that is optimized for recording ABRs, but a full multi-channel montage
583 could also be used to more fully explore the interactions between subcortical and cortical processing of
584 naturalistic speech. The potential for new knowledge about how the brain processes naturalistic and
585 engaging stimuli cannot be undersold.

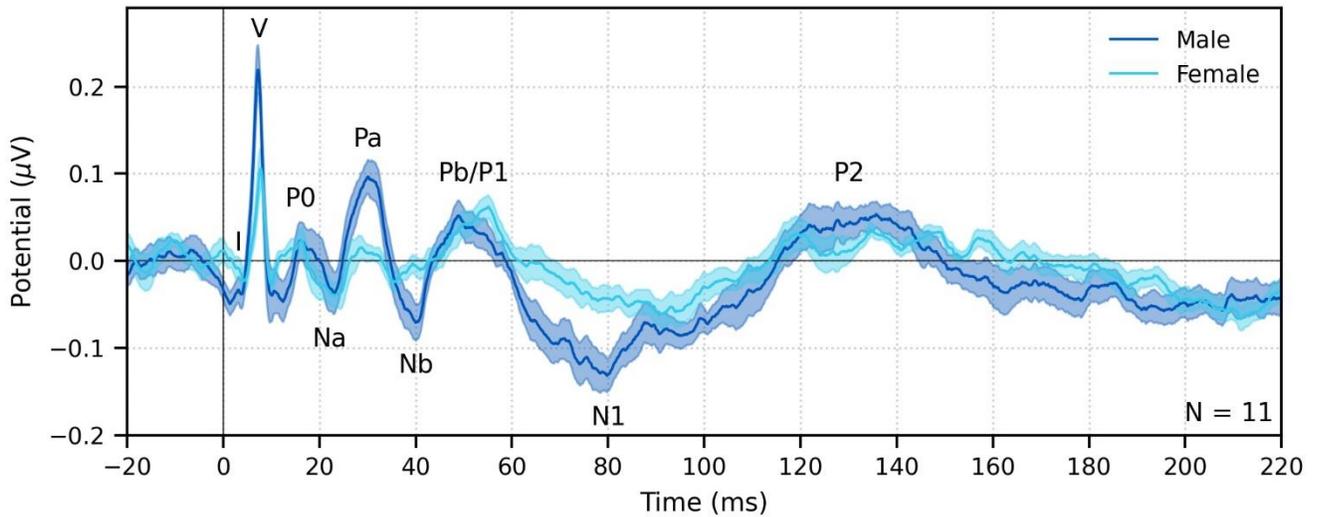


Figure 10. The range of lags can be extended to allow early, middle and late latency responses to be analyzed from the same recording to broadband peaky speech. Average waveforms across subjects (areas show ± 1 SEM) are shown for responses measured to 32 minutes of broadband peaky speech narrated by a male (dark blue) and female (light blue). Responses were high-pass filtered at 30 Hz using a first order Butterworth filter, but different filter parameters can be used to focus on each stage of processing. Canonical waves of the ABR, MLR and LLR are labeled for the male-narrated speech. Due to adaptation, amplitudes of the late potentials are smaller than typically seen with other stimuli that are shorter in duration with longer inter-stimulus intervals than our continuous speech. Waves I and III become more clearly visible by applying a 150 Hz high-pass cutoff.

586

587 The flexibility of peaky speech to evoke broadband and frequency-specific responses from distinct
588 subcortical structures facilitates new lines of query, both in neuroscientific and clinical domains. Speech
589 often occurs within a mixture of sounds, such as other speech sources, background noise, or music.
590 Furthermore, visual cues from a talker's face are often available to aid speech understanding, particularly
591 in environments with low SNR (e.g., Bernstein and Grant, 2009; Grant et al., 2007). Peaky speech allows
592 investigation into the complex subcortical processing that underpins successful listening in these scenarios
593 using naturalistic, engaging tasks. Indeed, previous methods have been quite successful in elucidating
594 cortical processing of speech under these conditions (O'Sullivan et al., 2019; Teoh and Lalor, 2019).
595 Finally, as aforementioned, the ability to customize peaky speech for measuring frequency-specific
596 responses provides potential applications to clinical research in the context of facilitating assessment of
597 supra-threshold hearing function and changes following intervention strategies and technologies.

598 In summary, the peaky speech paradigm is a viable method for recording canonical waveforms and
599 frequency-specific responses to an engaging, continuous speech stimulus. The customizability and
600 flexibility of peaky speech facilitates new ways of investigating subcortical contributions to speech
601 processing and holds great potential for future implementation into clinical assessment of hearing function.

602

603 METHODS

604

605 Participants

606 Data were collected over 3 experiments that were conducted under a protocol approved by the University
607 of Rochester Research Subjects Review Board. All subjects gave informed consent before the experiment
608 began and were compensated for their time. In each of experiments 1 and 2, there were equipment
609 problems during testing for one subject, rendering data unusable in the analyses. Therefore, there were a
610 total of 22, 11, and 11 subjects included in experiments 1, 2, and 3 respectively. Four subjects completed
611 both experiments 1 and 2, and 2 subjects completed both experiments 2 and 3. The 38 unique subjects (25
612 females, 66%) were aged 18–32 years with a mean \pm SD age of 23.0 ± 3.6 years. Audiometric screening
613 confirmed subjects had normal hearing in both ears, defined as thresholds ≤ 20 dB HL from 250 to 8000
614 Hz. All subjects identified English as their primary language.

615

616 Stimulus presentation and EEG measurement

617 In each experiment, subjects listened to 128 minutes of continuous speech stimuli while reclined in a
618 darkened sound booth. They were not required to attend to the speech and were encouraged to relax and
619 to sleep. Speech was presented at an average level of 65 dB SPL over ER-2 insert earphones (Etymotic
620 Research, Elk Grove, IL) plugged into an RME Babyface Pro digital soundcard (RME, Haimhausen,
621 Germany) via an HB7 headphone amplifier (Tucker Davis Technologies, Alachua, FL). Stimulus
622 presentation was controlled by a custom python script using publicly available software (available at
623 <https://github.com/LABSN/expyfun>; Larson et al., 2014). We interleaved conditions in order to prevent slow
624 impedance drifts or transient periods of higher EEG noise from unevenly affecting one condition over the
625 others. Physical measures to reduce stimulus artifact included: 1) hanging earphones from the ceiling so
626 that they were as far away from the EEG cap as possible; and 2) sending an inverted signal to a dummy
627 earphone (blocked tube) attached in the same physical orientation to the stimulus presentation earphones
628 in order to cancel electromagnetic fields away from transducers. The soundcard also produced a digital
629 signal at the start of each epoch, which was converted to trigger pulses through a custom trigger box
630 (modified from a design by the National Acoustic Laboratories, Sydney, NSW, Australia) and sent to the
631 EEG system so that audio and EEG data could be synchronized with sub-millisecond precision.

632 EEG was recorded using BrainVision's PyCorder software. Ag/AgCl electrodes were placed at the high
633 forehead (FCz, active non-inverting), left and right earlobes (A1, A2, inverting references), and the frontal
634 pole (Fpz, ground). These were plugged into an EP-Preamp system specifically for recording ABRs,
635 connected to an ActiCHamp recording system, both manufactured by BrainVision. Data were sampled at
636 10,000 Hz and high-pass filtered at 0.1 Hz. Offline, raw data were high-pass filtered at 1 Hz using a first-
637 order causal Butterworth filter to remove slow drift in the signal, and then notch filtered with 5 Hz wide
638 second-order infinite impulse response (IIR) notch filters to remove 60 Hz and its first 3 odd harmonics
639 (180, 300, 420 Hz). To optimize parameters for viewing the ABR and MLR components of peaky speech
640 responses, we evaluated several orders and high-pass cutoffs to the filters. Early waves of the broadband
641 peaky ABRs were best visualized with a 150 Hz cutoff, whereas a lower cutoff frequency of 30 Hz was
642 necessary to view the ABR and MLR of the broadband responses. Conservative filtering with a first order
643 filter was sufficient with these cutoff frequencies.

644

645 Speech stimuli and conditions

646 Speech stimuli were taken from two audiobooks. The first was *The Alchemyst* (Scott, 2007), read by a
647 male narrator and used in all 3 experiments. The second was *A Wrinkle in Time* (L'Engle, 2012), read by a
648 female narrator and used in experiments 2 and 3. These stimuli were used in Maddox and Lee (2018), but
649 in that study a gentle high-pass filter was applied which was not done for this study. Briefly, the audiobooks
650 were resampled to 44,100 Hz and then silent pauses were truncated to 0.5 s. Speech was segmented into
651 64 s epochs with 1 s raised cosine fade-in and fade-out. Because conditions were interleaved, the last 4 s
652 of a segment were repeated in the next segment so that subjects could pick up where they left off if they
653 were listening.

654 In experiment 1 subjects listened to 3 conditions of male speech (42.7 min each): unaltered speech, re-
655 synthesized broadband peaky speech, and re-synthesized multiband peaky speech (see below for a
656 description of re-synthesized speech). In experiment 2 subjects listened to 4 conditions of re-synthesized
657 peaky speech (32 minutes each): male and female narrators of both broadband and multiband peaky
658 speech. For these first 2 experiments, speech was presented diotically (same speech to both ears). In
659 experiment 3 subjects listened to both male and female dichotic (different speech in each ear) multiband
660 peaky speech designed for audiological applications (64 min each). The same 64 s of speech was
661 presented simultaneously to each ear, but the stimuli were dichotic due to how the re-synthesized
662 multiband speech was created (see below).

663

664 Stimulus design

665 The brainstem responds best to impulse-like stimuli, so we re-synthesized the speech segments from the
666 audiobooks (termed "unaltered") to create 3 types of "peaky" speech, with the objectives of 1) evoking
667 additional waves of the ABR reflecting other neural generators, and 2) measuring responses to different

668 frequency regions of the speech. The process is described in detail below, but is best read in tandem with
669 the code that will be publicly available (<https://github.com/maddoxlab>). Figure 11 compares the unaltered
670 speech and re-synthesized broadband and multiband peaky speech. Comparing the pressure waveforms
671 shows that the peaky speech is as click-like as possible, but comparing the spectrograms (how sound
672 varies in amplitude at every frequency and time point) shows that the overall spectrotemporal content that
673 defines speech is basically unchanged by the re-synthesis. See supplementary files for audio examples of
674 each stimulus type for both narrators.
675

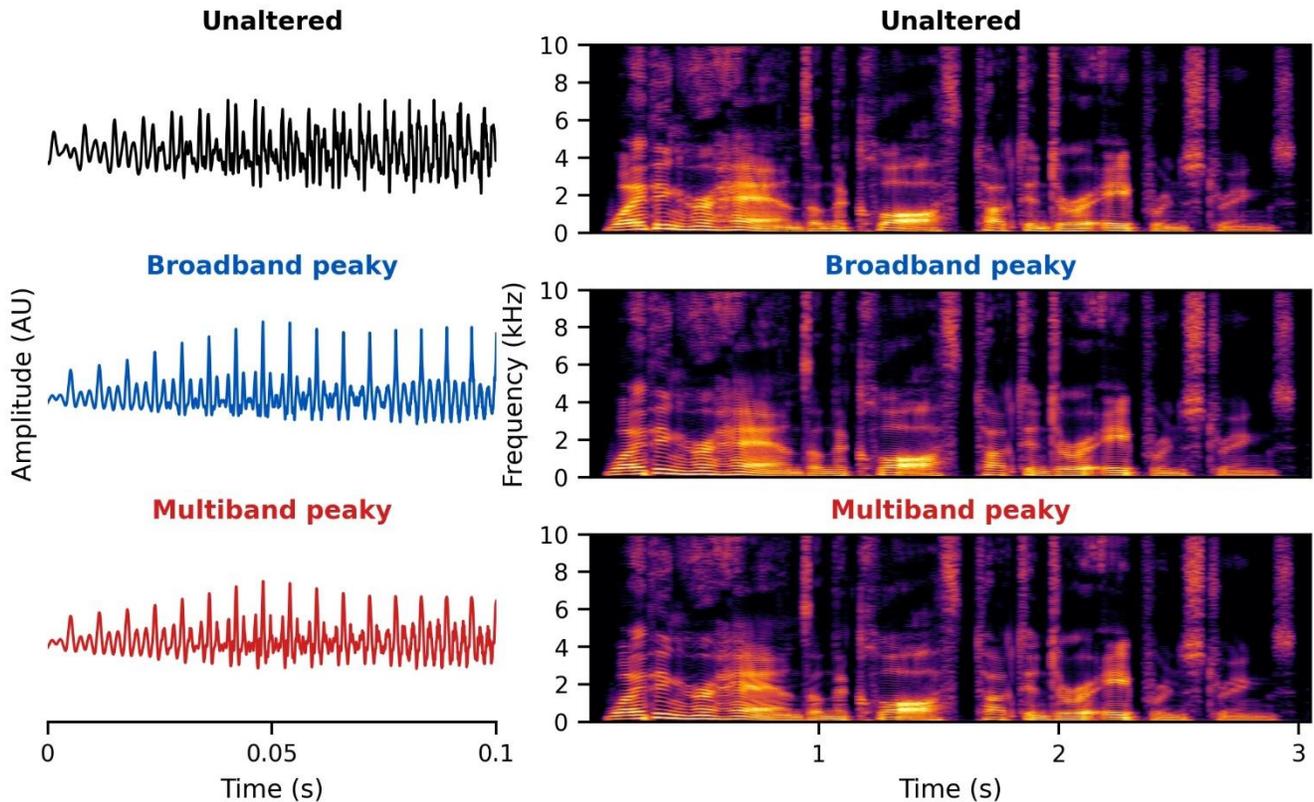


Figure 11. Unaltered speech waveform (top left) and spectrogram (top right) compared to re-synthesized broadband peaky speech (middle left and right) and multiband peaky speech (bottom left and right). Comparing waveforms shows that the peaky speech is as “click-like” as possible, while comparing the spectrograms shows that the overall spectrotemporal content that defines speech is basically unchanged by the re-synthesis. A naïve listener is unlikely to notice that any modification has been performed, and subjective listening confirms the similarity. Yellow/lighter colors represent larger amplitudes than purple/darker colors in the spectrogram. See supplementary files for audio examples of each stimulus type for both narrators.

676 Broadband peaky speech

677 Voiced speech comprises rapid openings and closings of the vocal folds which are then filtered by the
678 mouth and vocal tract to create different vowel and consonant sounds. The first processing step in creating
679 peaky speech was to use speech processing software (PRAAT; Boersma and Weenink, 2018) to extract
680 the times of these glottal pulses. Sections of speech where glottal pulses were within 17 ms of each other
681 were considered voiced (vowels and voiced consonants like /z/). 17 ms is the longest inter-pulse interval
682 one would expect in natural speech because it is the inverse of 60 Hz, the lowest pitch at which someone
683 with a deep voice would likely speak. A longer gap in pulse times was considered a break between voiced
684 sections. These segments were identified in a “mixer” function of time, with indices of 1 indicating unvoiced
685 and 0 indicating voiced segments (and would later be responsible for time-dependent blending of re-
686 synthesized and natural speech, hence its name). Transitions of the binary mixer function were smoothed
687 using a raised cosine envelope spanning the time between the first and second pulses, as well as the last
688 two pulses of each voiced segment. During voiced segments, the glottal pulses set the fundamental
689 frequency of speech (i.e., pitch), which were allowed to vary from a minimum to maximum of 60–350 Hz for
690 the male narrator and 90–500 Hz for the female narrator. For the male and female narrators, these pulses
691 gave a mean \pm SD fundamental frequency (i.e., pulse rate) in voiced segments of 115.1 ± 6.7 Hz and 198.1

692 ± 20 Hz respectively, and a mean \pm SD pulses per second over the entire 64 s, inclusive of unvoiced
693 periods and silences, of 69.1 ± 5.7 Hz and 110.8 ± 11.4 respectively. These pulse times were smoothed
694 using 10 iterations of replacing pulse time p_i with the mean of pulse times p_{i-1} to p_{i+1} if the \log_2 absolute
695 difference in the time between p_i and p_{i-1} and p_{i+1} was less than $\log_2(1.6)$.

696 The fundamental frequency of voiced speech is dynamic, but the signal always consists of a set of integer-
697 related frequencies (harmonics) with different amplitudes and phases. To create the waveform component
698 at the fundamental frequency, $f_0(t)$, we first created a phase function, $\varphi(t)$, which increased smoothly by
699 2π between glottal pulses within the voiced sections as a result of cubic interpolation. We then computed
700 the spectrogram of the unaltered speech waveform – which is a way of analyzing sound that shows its
701 amplitude at every time and frequency (Figure 11, top-right) – which we called $A[t, f_0(t)]$. We then created
702 the fundamental component of the peaky speech waveform as:

$$703 \quad h_0(t) = A[t, f_0(t)] \cos[\varphi(t)].$$

704 This waveform has an amplitude that changes according to the spectrogram but always peaks at the time
705 of the glottal pulses.

706 Next the harmonics of the speech were synthesized. The k^{th} harmonic of speech is at a frequency of
707 $(k + 1)f_0$ so we synthesized each harmonic waveform as:

$$708 \quad h_k(t) = A[t, (k + 1)f_0(t)] \cos[(k + 1)\varphi(t)].$$

709 Each of these harmonic waveforms has multiple peaks per period of the fundamental, but every harmonic
710 also has a peak at exactly the time of the glottal pulse. Because of these coincident peaks, when the
711 harmonics are summed to create the re-synthesized voiced speech, there is always a large peak at the
712 time of the glottal pulse. In other words, the phases of all the harmonics align at each glottal pulse, making
713 the pressure waveform of the speech appear “peaky” (left-middle panel of Figure 11).

714 The resultant re-synthesized speech contained only the voiced segments of speech and was missing
715 unvoiced sounds like /s/ and /k/. Thus the last step was to mix the re-synthesized voiced segments with the
716 original unvoiced parts. This was done by cross-fading back and forth between the unaltered speech and
717 re-synthesized speech during the unvoiced and voiced segments respectively, using the binary mixer
718 function created when determining where the voiced segments occurred. We also filtered the peaky
719 speech to an upper limit of 8 kHz, and used the unaltered speech above 8 kHz, to improve the quality of
720 voiced consonants such as /z/. Filter properties for the broadband peaky speech are further described
721 below in the “Band filters” subsection.

722

723 *Multiband peaky speech*

724 The same principles to generate broadband peaky speech were applied to create stimuli designed to
725 investigate the brainstem’s response to different frequency bands that comprise speech. This makes use of
726 the fact that over time, speech signals with slightly different f_0 are independent, or have (nearly) zero
727 cross-correlation, at the lags for the ABR. To make each frequency band of interest independent, we
728 shifted the fundamental frequency and created a fundamental waveform and its harmonics as:

$$729 \quad h_k(t) = A[t, (k + 1)f_0(t)] \cos[(k + 1)\varphi(t)],$$

730 where

$$731 \quad \varphi(t) = 2\pi \int_0^t (f_0(\tau) + f_\Delta) d\tau,$$

732 and where f_Δ is the small shift in fundamental frequency.

733 In these studies, we increased fundamentals for each frequency band by the square root of each
734 successive prime number and subtracting one, resulting in a few tenths of a hertz difference between
735 bands. The first, lowest frequency band contained the un-shifted f_0 . Responses to this lowest, un-shifted
736 frequency band showed some differences from the common component for latencies > 30 ms that were not
737 present in the other, higher frequency bands (Figure 4, 0–1 kHz band), suggesting some low-frequency
738 privilege/bias in this response. Therefore, we suggest that following studies create independent frequency
739 bands by synthesizing a new fundamental for each band. The static shifts described above could be used,
740 but we suggest an alternative method that introduces random dynamic frequency shifts of up to ± 1 Hz over

741 the duration of the stimulus. From this random frequency shift we can compute a dynamic random phase
 742 shift, to which we also add a random starting phase, θ_{Δ} , which is drawn from a uniform distribution between
 743 0 and 2π . The phase function from the above set of formulae would be replaced with this random dynamic
 744 phase function:

745
$$\varphi(t) = 2\pi \int_0^t [f_0(\tau) + f_{\Delta}(\tau)] d\tau + \theta_{\Delta}$$

746 This random f_0 shift method is described further in the supplementary material and validation data from
 747 one subject is provided in Supplemental Figure 2. Responses from all four bands show more consistent
 748 resemblance to the common component, indicating that this method is effective at reducing stimulus-
 749 related bias. However, low-frequency dependent differences remained, suggesting there is also unique
 750 neural-based low-frequency activity to the speech-evoked responses.

751 This re-synthesized speech was then band-pass filtered to the frequency band of interest (e.g. from 0–1
 752 kHz or 2–4 kHz). This process was repeated for each independent frequency band, then the bands were
 753 mixed together and then these re-synthesized voiced parts were mixed with the original unaltered voiceless
 754 speech. This peaky speech comprised octave bands with center frequencies of: 707, 1414, 2929, 5656 Hz
 755 for experiments 1 and 2, and of 500, 1000, 2000, 4000, 8000 Hz for experiment 3. Note that for the lowest
 756 band, the actual center frequency was slightly lower because the filters were set to pass all frequencies
 757 below the upper cutoff. Filter properties for these two types of multiband speech are shown in the middle
 758 and right panels of Figure 14 and further described below in the “Band filters” subsection. For the dichotic
 759 multiband peaky speech, we created 10 fundamental waveforms – 2 in each of the five filter bands for the
 760 two different ears, making the output audio file stereo (or dichotic). We also filtered this dichotic multiband
 761 peaky speech to an upper limit of 11.36 kHz to allow for the highest band to have a center frequency of 8
 762 kHz and octave width. The relative mean-squared magnitude in decibels for components of the multiband
 763 peaky speech (4 filter bands) and dichotic (audiological) multiband peaky speech (5 filter bands) are shown
 764 in Figure 12.

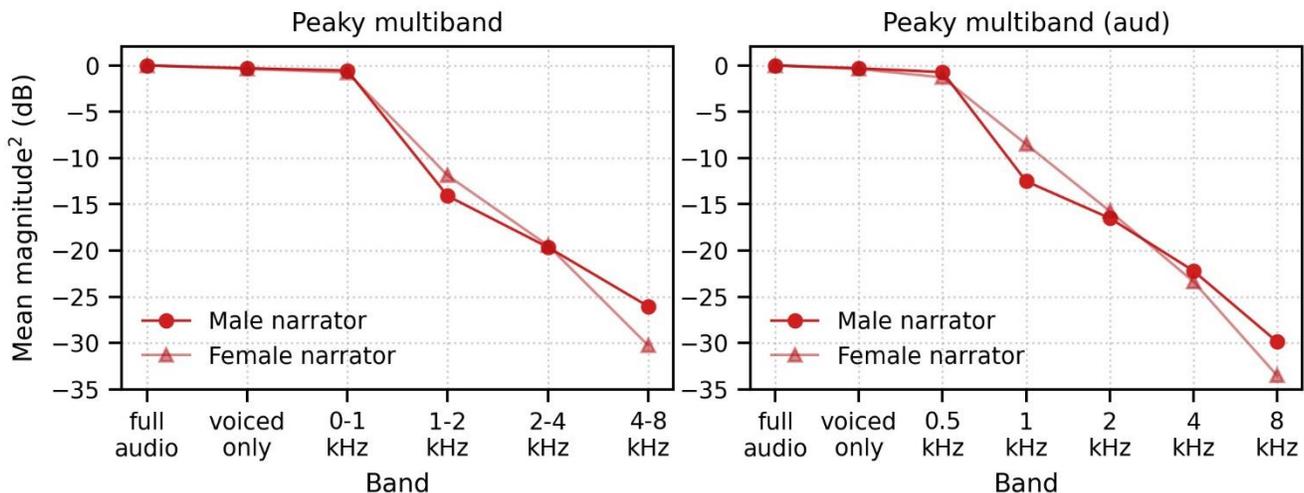


Figure 12. Relative mean-squared magnitude in decibels of multiband peaky speech with 4 filter bands (left) and 5 filter bands (right) for male-(blue) and female-(orange) narrated speech. The full audio comprises unvoiced and re-synthesized voiced sections, which was presented to the subjects during the experiments. The other bands reflect the relative magnitude of the voiced sections (voiced only), and each filtered frequency band.

765

766 For peaky speech, the re-synthesized speech waveform was presented during the experiment but the
 767 pulse trains were used as the input stimulus for calculating the response (i.e., the regressor, see Response
 768 derivation section below). These pulse trains all began and ended together in conjunction with the onset
 769 and offset of voiced sections of the speech. To verify which frequency ranges of the multiband pulse trains
 770 were independent across frequency bands, and would thus yield truly band-specific responses, we
 771 conducted a spectral coherence analyses on the pulse trains. All 60 unique 64 s sections of each male-
 772 and female-narrated multiband peaky speech used in the three experiments were sliced into 1 s segments
 773 for a total of 3,840 slices. Phase coherence across frequency was then computed across these slices for
 774 each combination of pulse trains according to the formula:

775

$$C_{xy} = \frac{|E[\mathcal{F}\{x_i\}^* \mathcal{F}\{y_i\}]|}{\sqrt{E[\mathcal{F}\{x_i\}^* \mathcal{F}\{x_i\}] E[\mathcal{F}\{y_i\}^* \mathcal{F}\{y_i\}]}}$$

776

where C_{xy} denotes coherence between bands x and y , $E[\]$ the average across slices, \mathcal{F} the fast Fourier transform, $*$ complex conjugation, x_i the pulse train for slice i in band x , and y_i the pulse train for slice i in band y .

777

778

779

Spectral coherence for each narrator is shown in Figure 13. For the 4-band multiband peaky speech used in experiments 1 and 2 there were 6 pulse train comparisons. For the audiological multiband peaky speech used in experiment 3, there were 5 bands for each of 2 ears, resulting in 10 pulse trains and 45 comparisons. All 45 comparisons are shown in Figure 13. Pulse trains were coherent (> 0.1) up to a maximum of 71 and 126 Hz for male- and female-narrated speech respectively, which roughly correspond to the mean \pm SD pulse rates (calculated as total pulses / 64 s) of 69.1 ± 5.7 Hz and 110.8 ± 11.4 Hz respectively. This means that above ~ 130 Hz the stimuli were no longer coherent and evoked frequency-specific responses. Importantly, responses would be to correlated stimuli, i.e., not frequency-specific, at frequencies below this cutoff and would result in a low-frequency response component that is present in (or common to) all band responses.

789

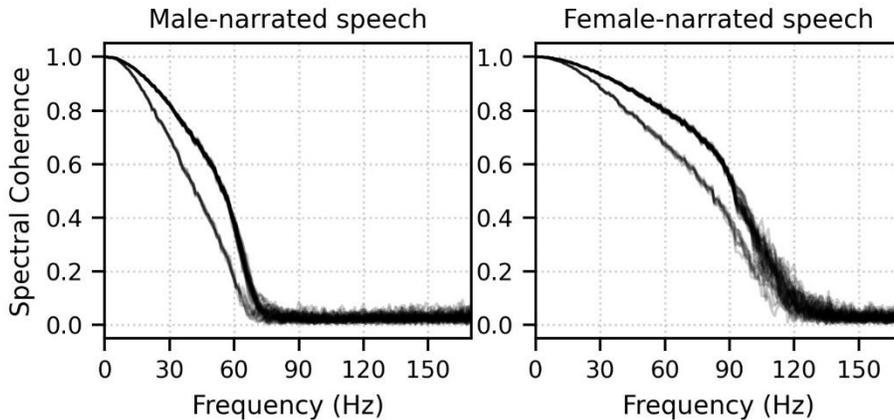


Figure 13. Spectral coherence of pulse trains for multiband peaky speech narrated by a male (left) and female (right). Spectral coherence was computed across 1 s slices from 60 unique 64 s multiband peaky speech segments (3,840 total slices) for each combination of bands. Each light gray line represents the coherence for one band comparison. There were 45 comparisons across the 10-band (audiological) speech used in experiment 3 (5 frequency bands x 2 ears). Pulse trains (i.e., the input stimuli, or regressors, for the deconvolution) were frequency-dependent (coherent) below 72 Hz for the male multiband speech and 126 Hz for the female multiband speech.

790

To identify the effect of the low-frequency stimulus coherence in the responses, we computed the common component across pulse trains by creating an averaged response to 6 additional “fake” pulse trains that were created during stimulus design but were not used during creation of the multiband peaky speech wav files. The common component was assessed for both “fake” pulse trains taken from shifts lower than the original fundamental frequency and those taken from shifts higher than the highest “true” re-synthesized fundamental frequency. To assess frequency-specific responses to multiband speech, we subtracted this common component from the band responses. Alternatively, one could simply high-pass the stimuli at 150 Hz using a first-order causal Butterworth filter (being mindful of edge artifacts). However, this high-pass filtering reduces response amplitude and may affect response detection (see Results for more details).

798

799

We also verified the independence of the stimulus bands by treating the regressor pulse train as the input to a system whose output was the rectified stimulus audio and performed deconvolution (see Deconvolution and Response derivation section below). Further details are provided in the supplementary material. The responses are given in Supplemental Figure 5, and showed that the non-zero responses only occurred when the correct pulse train was paired with the correct audio.

804

805

806 Band filters

807 Because the fundamental frequencies for each frequency band were designed to be independent over
808 time, the band filters for the speech were designed to cross over in frequency at half power. To make the
809 filter, the amplitude was set by taking the square root of the specified power at each frequency. Octave
810 band filters were constructed in the frequency domain by applying trapezoids – with center bandwidth and
811 roll-off widths of 0.5 octaves. For the first (lowest frequency) band, all frequencies below the high-pass
812 cutoff were set to 1, and likewise for all frequencies above the low-pass cutoff for the last (highest
813 frequency) band were set to 1 (Figure 14 top row). The impulse response of the filters was assessed by
814 shifting the inverse FFT (IFFT) of the bands so that time zero was in the center, and then applied a Nuttall
815 window, thereby truncating the impulse response to length of 5 ms (Fig 14 middle row). The actual
816 frequency response of the filter bands was assessed by taking the FFT of the impulse response and
817 plotting the magnitude (Figure 14 bottom row).

818 As mentioned above, broadband peaky speech was filtered to an upper limit of 8 kHz for diotic peaky
819 speech and 11.36 kHz for dichotic peaky speech. This band filter was constructed from the second last
820 octave band filter from the multiband filters (i.e., the 4–8 kHz band from the top-middle of Figure 14, dark
821 red line) by setting the amplitude of all frequencies less than the high-pass cutoff frequency to 1 (Figure 14
822 top-left panel, blue line). As mentioned above, unaltered (unvoiced) speech above 8 kHz (diotic) or 11.36
823 kHz (dichotic) was mixed with the broadband and multiband peaky speech, which was accomplished by
824 applying the last (highest) octave band filter (8+ or 11.36+ kHz band, black line) to the unaltered speech
825 and mixing this band with the re-synthesized speech using the other bands.

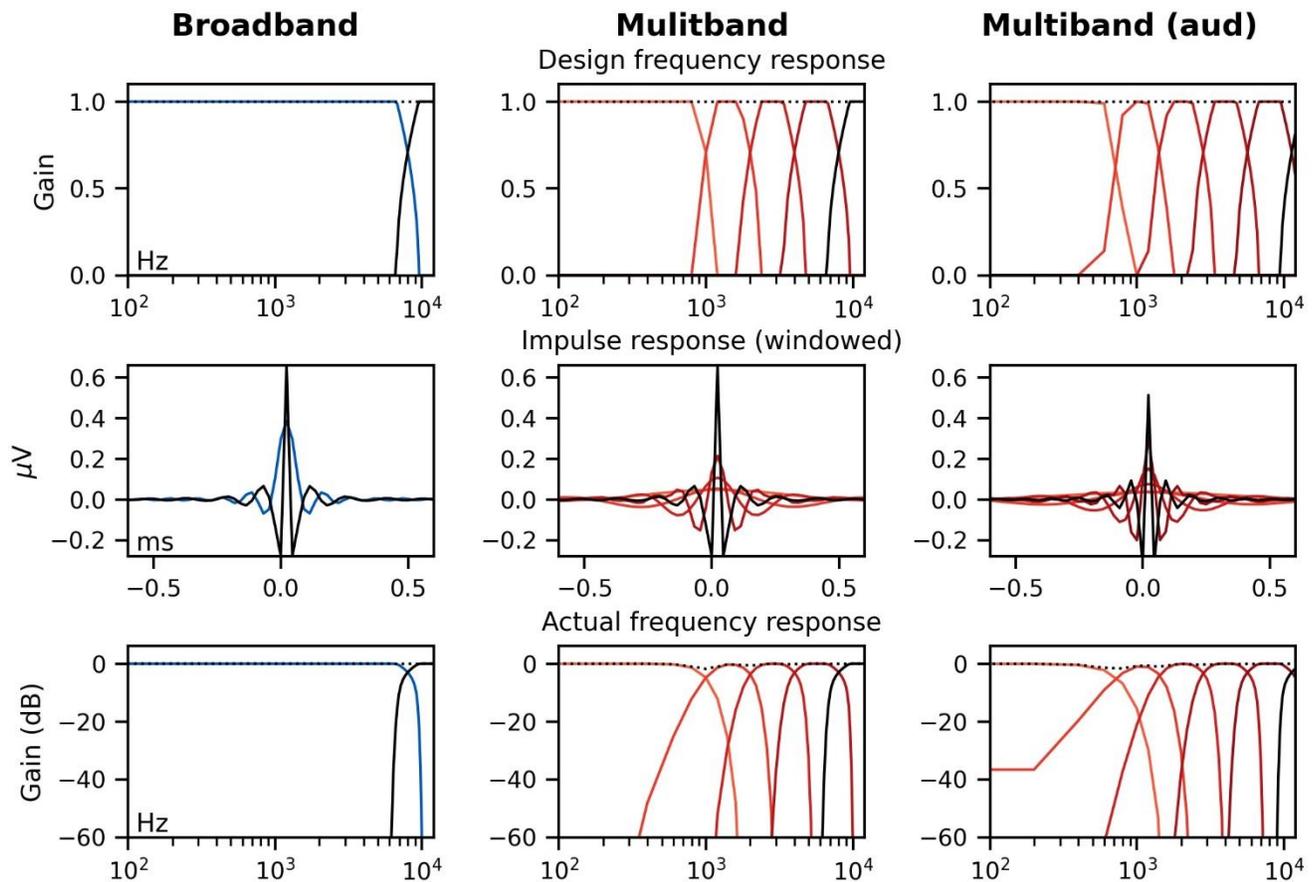


Figure 14. Octave band filters used to create re-synthesized broadband peaky speech (left, blue), diotic multiband peaky speech with 4 bands (middle, red), and dichotic multiband peaky speech using 5 bands with audiological center frequencies (right, red). The last band (2nd, 5th, 6th respectively, black line) was used to filter the high-frequencies of unaltered speech during mixing to improve the quality of voiced consonants. The designed frequency response using trapezoids (top) were converted into the time-domain using IFFT, shifted and Nuttall windowed to create impulse responses (middle), which were then used to assess the actual frequency response by converting into the frequency domain using FFT (bottom).

827 *Alternating polarity*

828 To limit stimulus artifact, we also alternated polarity between segments of speech. To identify regions to flip
 829 polarity, the envelope of speech was extracted using a first-order causal Butterworth low-pass filter with a
 830 cutoff frequency of 6 Hz applied to the absolute value of the waveform. Then flip indices were identified
 831 where the envelope became less than 1 percent of the median envelope value, and then a function that
 832 changed back and forth between 1 and -1 at each flip index was created. This function of spikes was
 833 smoothed using another first-order causal Butterworth low-pass filter with a cutoff frequency of 10,000 Hz,
 834 which was then multiplied with the re-synthesized speech before saving to a wav file.

835

836 Response derivation

837 *Deconvolution*

838 The peaky-speech ABR was derived by using deconvolution, as in previous work (Maddox and Lee, 2018),
 839 though the computation was performed in the frequency domain for efficiency. The speech was considered
 840 the input to a linear system whose output was the recorded EEG signal, with the ABR computed as the
 841 system's impulse response. As in Maddox and Lee (2018), for the unaltered speech, we used the half-
 842 wave rectified audio as the input waveform. Half-wave rectification was accomplished by separately
 843 calculating the response to all positive and all negative values of the input waveform for each epoch and
 844 then combining the responses together during averaging. For our new re-synthesized peaky speech, the
 845 input waveform was the sequence of impulses that occurred at the glottal pulse times and corresponded to
 846 the peaks in the waveform. Figure 15 shows a section of stimulus and the corresponding input signal of
 847 glottal pulses used in the deconvolution.

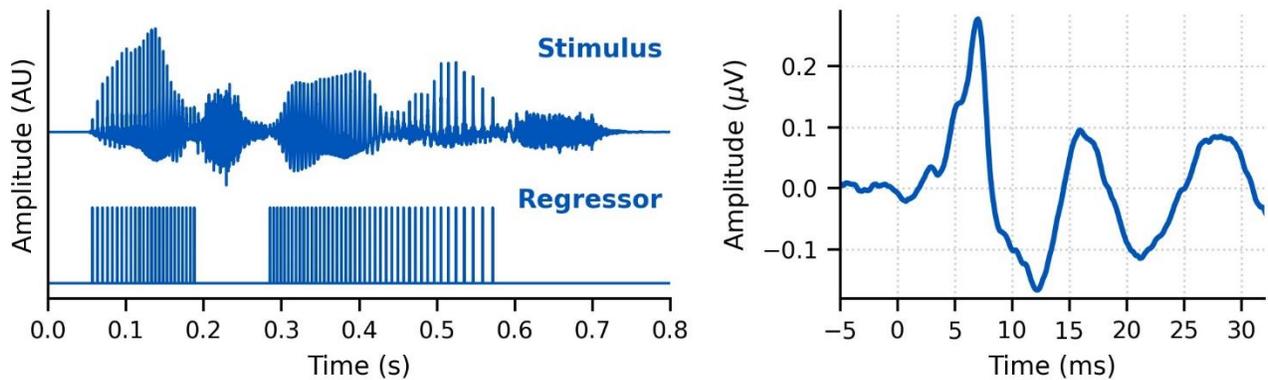


Figure 15. Left: A segment of broadband peaky speech stimulus (top) and the corresponding glottal pulse train (bottom) used in calculating the broadband peaky speech response. Right: An example broadband peaky speech response from a single subject. The response shows ABR waves I, III, and V at ~3, 5, 7 ms respectively. It also shows later peaks corresponding to thalamic and cortical activity at ~17 and 27 ms respectively.

848 The half-wave rectified waveforms and glottal pulse sequences were down-sampled to the EEG sampling
 849 frequency prior to deconvolution. To avoid temporal splatter due to standard downsampling, the pulse
 850 sequences were resampled by placing unit impulses at sample indices closest to each pulse time.
 851 Regularization was not necessary because the amplitude spectra of these regressors were sufficiently
 852 broadband. For efficiency, the time-domain response waveform, w , for a given 64 s epoch was calculated
 853 using frequency-domain division for the deconvolution, with the numerator the cross-spectral density
 854 (corresponding to the cross-correlation in the time domain) of the stimulus regressor and EEG response,
 855 and the denominator the power spectral density of the stimulus regressor (corresponding to its
 856 autocorrelation in the time domain). For a single epoch, that would be:

857
$$w = \mathcal{F}^{-1} \left\{ \frac{\mathcal{F}\{x\}^* \mathcal{F}\{y\}}{\mathcal{F}\{x\}^* \mathcal{F}\{x\}} \right\}$$

858 where \mathcal{F} denotes the fast Fourier transform, \mathcal{F}^{-1} the inverse fast Fourier Transform, * complex conjugation,
 859 x the input stimulus regressor (half-wave rectified waveform or glottal pulse sequence), and y the EEG

860 data for each epoch. We used methods incorporated into the mne-python package (Gramfort et al., 2013).
861 In practice, we made adaptations to this formula to improve SNR with Bayesian-like averaging (see below).
862 For multiband peaky speech the same EEG was deconvolved with the pulse train of each band separately,
863 and then with an additional 6 “fake” pulse trains to derive the common component across bands due to the
864 pulse train coherence at low frequencies (shown in Figure 13). The averaged response across these 6 fake
865 pulse trains, or common component, was then subtracted from the multiband responses to identify the
866 frequency-specific band responses.

867

868 *Response averaging*

869 The quality of the ABR waveforms as a function of each type of stimulus was of interest, so we calculated
870 the averaged response after each 64 s epoch. We followed a Bayesian-like process (Elberling and
871 Wahlgreen, 1985) to account for variations in noise level across the recording time (such as slow drifts or
872 movement artifacts) and to avoid rejecting data based on thresholds. Each epoch was weighted by its
873 inverse variance, $1/\sigma_i^2$, to the sum of the inverse variances of all epochs. Thus, epoch weights, b_i , were
874 calculated as follows:

$$875 \quad b_i = \frac{1/\sigma_i^2}{\sum_{i=1}^n 1/\sigma_i^2}$$

876 where i is the epoch number and n is the number of epochs collected. For efficiency, weighted averaging
877 was completed during deconvolution. Because auto-correlation of the input stimulus (denominator of the
878 frequency domain division) was similar across epochs it was averaged with equal weighting. Therefore, the
879 numerator of the frequency domain division was summed across weighted epochs and the denominator
880 averaged across epochs, according to the following formula:

$$881 \quad w = \mathcal{F}^{-1} \left\{ \frac{\sum_{i=1}^n b_i \mathcal{F}\{x_i\}^* \mathcal{F}\{y_i\}}{\sum_{i=1}^n \frac{1}{n} \mathcal{F}\{x_i\}^* \mathcal{F}\{x_i\}} \right\}$$

882

883 where w is the average response waveform, i is again the epoch number, n is the number of epochs
884 collected.

885 Due to the circular nature of the discrete frequency domain deconvolution, the resulting response has an
886 effective time interval of $[0, 32]$ s at the beginning and $[-32, 0]$ s at the end, so that concatenating the two –
887 with the end first – yields the response from $[-32, 32]$ s. Consequently, to avoid edge artifacts, all filtering
888 was performed after the response was shifted to the middle of the 64 s time window. To remove high-
889 frequency noise and some low-frequency noise, the average waveform was band-pass filtered between
890 30–2000 Hz using a first-order causal Butterworth filter. An example of this weighted average response to
891 broadband peaky speech is shown in the right panel of Figure 15. This bandwidth of 30 to 2000 Hz is
892 sufficient to identify additional waves in the brainstem and middle latency responses (ABR and MLR
893 respectively). To further identify earlier waves of the auditory brainstem responses (i.e., waves I and III),
894 responses were high-pass filtered at 150 Hz using a first-order causal Butterworth filter. This filter was
895 determined to provide the best morphology without compromising the response by comparing responses
896 filtered with common high-pass cutoffs of 1, 30, 50, 100 and 150 Hz each combined with first, second and
897 fourth order causal Butterworth filters.

898

899 *Response normalization*

900 An advantage of this method over our previous one (Maddox and Lee, 2018) is that because the regressor
901 comprises unit impulses, the deconvolved response is given in meaningful units which are the same as the
902 EEG recording, namely microvolts. With a continuous regressor, like the half-wave rectified speech
903 waveform, this is not the case. Therefore, to compare responses to half-wave rectified speech versus
904 glottal pulses, we calculated a normalization factor, g , based on data from all subjects:

905

$$g = \frac{1/n \sum_{i=1}^n \sigma_{u,i}}{1/n \sum_{i=1}^n \sigma_{p,i}}$$

906

907

908

909

910

911

912

913

914

915

916

where n is the number of subjects, $\sigma_{u,i}$ is the SD of subject i 's response to unaltered speech between 0–20 ms, and $\sigma_{p,i}$ is the same for the broadband peaky speech. Each subject's responses to unaltered speech were multiplied by this normalization factor to bring these responses within a comparable amplitude range as those to broadband peaky speech. Consequently, amplitudes were not compared between responses to unaltered and peaky speech. This was not our prime interest, rather we were interested in latency and presence of canonical component waves. In this study the normalization factor was 0.26, which cannot be applied to other studies because this number also depends on the scale when storing the digital audio. In our study, this unitless scale was based on a root-mean-square amplitude of 0.01. The same normalization factor was used when the half-wave rectified speech as used as the regressor with EEG collected in response to unaltered speech, broadband peaky speech and multiband peaky speech (Figure 2, Supplemental Figure 1).

917

918 *Response SNR calculation*

919

920

921

We were also interested in the recording time required to obtain robust responses to re-synthesized peaky speech. Therefore, we calculated the time it took for the ABR and MLR to reach a 0-dB SNR. The SNR of each waveform in dB, SNR_w , was estimated as:

922

$$SNR_w = 10 \log_{10} \left[\frac{\sigma_{S+N}^2 - \sigma_N^2}{\sigma_N^2} \right],$$

923

924

925

926

927

where σ_{S+N}^2 represents the variance (i.e., mean-subtracted energy) of the waveform between 0 and 15 ms or 30 ms for the ABR and MLR respectively (contains both component signals as well as noise, $S + N$), and σ_N^2 represents the variance of the noise, N , estimated by averaging the variances of 15 ms (ABR) to 30 ms (MLR) segments of the pre-stimulus baseline between -480 and -20 ms. Then the SNR for 1 min of recording, SNR_{60} , was computed from the SNR_w as:

928

$$SNR_{60} = SNR_w + 10 \log_{10}[60/t_w],$$

929

930

931

where t_w is the duration of the recording in seconds, as specified in the "Speech stimuli and conditions" subsection. For example, in experiment 3, the average waveform resulted from 64 min of recording, or a t_w of 3,840 s. The time to reach 0 dB SNR for each subject, $t_{0dB SNR}$, was estimated from this SNR_{60} by:

932

$$t_{0dB SNR} = 60 \times 10^{-SNR_{60}/10}.$$

933

934

Cumulative density functions were used to show the proportion of subjects that reached an SNR ≥ 0 dB and to determine the necessary acquisition times that can be expected for each stimulus on a group level.

935

936 Statistical Analyses

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

Data were checked for normality using the Shapiro-Wilk test. Waveform morphology of responses to different narrators was compared using Pearson correlations of the responses between 0 and 15 ms for the ABR waveforms or 0 and 40 ms for both ABR and MLR waveforms. The Wilcoxon signed rank test was used to determine whether narrator differences (waveform correlations) were significantly different than the correlations of the same EEG split into even and odd epochs with equal numbers of epochs from each narrator. The intraclass correlation coefficient type 3 (absolute agreement) was used to verify good agreement in peak latencies chosen by an experienced Audiologist and neuroscientist (MJP) at two different time points, 3 months apart. Independent t-tests with $\mu = 0$ were conducted on the peak latency differences of ABR/MLR waves for unaltered and broadband peaky speech. For multiband peaky speech, the component wave peak latency changes across frequency band were assessed with linear mixed effects regression using the lme4 and lmerTest packages in R (Bates et al., 2015; Kuznetsova et al., 2017; R Core Team, 2020). A likelihood ratio test was used to determine that the model significantly improved upon adding an orthogonal 2nd order polynomial of frequency band (transformed so that each coefficient is independent), with linear and quadratic components to estimate the decrease in latency (slope) and change in the rate of latency decrease with increasing frequency band, respectively. Random effects of

952 subject and each frequency band term were included to account for individual variability that is not
953 generalizable to the fixed effects. A power analysis was completed using the simR package (Green and
954 MacLeod, 2016), which uses a likelihood ratio test on 1000 Monte Carlo permutations of the response
955 variables based on the fitted model.

956

957 **ACKNOWLEDGMENTS**

958 The authors wish to thank Sara Fiscella for assistance with recruitment.

959

960 **FUNDING**

961 This work was supported by National Institute for Deafness and Other Communication Disorders
962 [R00DC014288] awarded to RKM.

963

964 **DATA AVAILABILITY**

965 Data and python code will be made available on the lab GitHub account (<https://github.com/maddoxlab>).

966

967 **REFERENCES**

- 968 Abdala C, Folsom RC. 1995. Frequency contribution to the click-evoked auditory brain-stem response in
969 human adults and infants. *J Acoust Soc Am* **97**:2394–2404. doi:10.1121/1.411961
- 970 Backer KC, Kessler AS, Lawyer LA, Corina DP, Miller LM. 2019. A novel EEG paradigm to simultaneously
971 and rapidly assess the functioning of auditory and visual pathways. *J Neurophysiol* **122**:1312–1329.
972 doi:10.1152/jn.00868.2018
- 973 Bajo VM, King AJ. 2012. Cortical modulation of auditory processing in the midbrain. *Front Neural Circuits*
974 **6**:114. doi:10.3389/fncir.2012.00114
- 975 Bajo VM, Nodal FR, Moore DR, King AJ. 2010. The descending corticocollicular pathway mediates
976 learning-induced auditory plasticity. *Nat Neurosci* **13**:253–260. doi:10.1038/nn.2466
- 977 Bates D, Mächler M, Bolker B, Walker S. 2015. Fitting Linear Mixed-Effects Models Using lme4. *J Stat*
978 *Softw* **67**:1–48.
- 979 Bernstein JGW, Grant KW. 2009. Auditory and auditory-visual intelligibility of speech in fluctuating maskers
980 for normal-hearing and hearing-impaired listeners. *J Acoust Soc Am* **125**:3358–3372.
981 doi:10.1121/1.3110132
- 982 Boersma P, Weenink D. 2018. Praat: doing phonetics by computer.
- 983 Bramhall N, Beach EF, Epp B, Le Prell CG, Lopez-Poveda EA, Plack CJ, Schaette R, Verhulst S, Carlon
984 B. 2019. The search for noise-induced cochlear synaptopathy in humans: Mission impossible? *Hear*
985 *Res* **377**:88–103. doi:10.1016/j.heares.2019.02.016
- 986 Burkard R, Hecox K. 1983. The effect of broadband noise on the human brainstem auditory evoked
987 response. I. Rate and intensity effects. *J Acoust Soc Am* **74**:1204–1213. doi:10.1121/1.390024
- 988 Burkard R, Shi Y, Hecox KE. 1990. A comparison of maximum length and Legendre sequences for the
989 derivation of brain-stem auditory-evoked responses at rapid rates of stimulation. *J Acoust Soc Am*
990 **87**:1656–1664. doi:10.1121/1.399413
- 991 Carney LH, Li T, McDonough JM. 2015. Speech Coding in the Brain: Representation of Vowel Formants by
992 Midbrain Neurons Tuned to Sound Fluctuations,.. *eNeuro* **2**. doi:10.1523/ENEURO.0004-15.2015
- 993 Chiappa KH, Gladstone KJ, Young RR. 1979. Brain Stem Auditory Evoked Responses: Studies of
994 Waveform Variations in 50 Normal Human Subjects. *Arch Neurol* **36**:81–87.
995 doi:10.1001/archneur.1979.00500380051005
- 996 Dau T, Wegner O, Mellert V, Kollmeier B. 2000. Auditory brainstem responses with optimized chirp signals
997 compensating basilar-membrane dispersion. *J Acoust Soc Am* **107**:1530–1540.
998 doi:10.1121/1.428438
- 999 Don M, Allen AR, Starr A. 1977. Effect of Click Rate on the Latency of Auditory Brain Stem Responses in
1000 Humans. *Ann Otol Rhinol Laryngol* **86**:186–195. doi:10.1177/000348947708600209
- 1001 Elberling C, Don M. 2008. Auditory brainstem responses to a chirp stimulus designed from derived-band
1002 latencies in normal-hearing subjects. *J Acoust Soc Am* **124**:3022–3037. doi:10.1121/1.2990709

- 1003 Elberling C, Wahlgreen O. 1985. Estimation of Auditory Brainstem Response, Abr, by Means of Bayesian
1004 Inference. *Scand Audiol* **14**:89–96. doi:10.3109/01050398509045928
- 1005 Forte AE, Etard O, Reichenbach T. 2017. The human auditory brainstem response to running speech
1006 reveals a subcortical mechanism for selective attention. *eLife* **6**:e27203. doi:10.7554/eLife.27203
- 1007 Geisler CD, Frishkopf LS, Rosenblith WA. 1958. Extracranial Responses to Acoustic Clicks in Man.
1008 *Science* **128**:1210–1211. doi:10.1126/science.128.3333.1210
- 1009 Goldstein R, Rodman LB. 1967. Early components of averaged evoked responses to rapidly repeated
1010 auditory stimuli. *J Speech Hear Res* **10**:697–705. doi:10.1044/jshr.1004.697
- 1011 Gorga MP, Johnson TA, Kaminski JR, Beauchaine KL, Garner CA, Neely ST. 2006. Using a Combination
1012 of Click- and Tone Burst-Evoked Auditory Brain Stem Response Measurements to Estimate Pure-
1013 Tone Thresholds. *Ear Hear* **27**:60–74. doi:10.1097/01.aud.0000194511.14740.9c
- 1014 Gorga MP, Kaminski JR, Beauchaine KA, Jesteadt W. 1988. Auditory brainstem responses to tone bursts
1015 in normally hearing subjects. *J Speech Hear Res* **31**:87–97. doi:10.1044/jshr.3101.87
- 1016 Gorga MP, Kaminski JR, Beauchaine KL, Bergman BM. 1993. A Comparison of Auditory Brain Stem
1017 Response Thresholds and latencies Elicited by Air- and Bone-Conducted Stimuli. *Ear Hear* **14**:85–
1018 94.
- 1019 Gramfort A, Luessi M, Larson E, Engemann DA, Strohmeier D, Brodbeck C, Goj R, Jas M, Brooks T,
1020 Parkkonen L, Hämäläinen M. 2013. MEG and EEG data analysis with MNE-Python. *Front Neurosci*
1021 **7**. doi:10.3389/fnins.2013.00267
- 1022 Grant KW, Tufts JB, Greenberg S. 2007. Integration efficiency for speech perception within and across
1023 sensory modalities by normal-hearing and hearing-impaired individuals. *J Acoust Soc Am*
1024 **121**:1164–1176. doi:10.1121/1.2405859
- 1025 Green P, MacLeod CJ. 2016. SIMR: an R package for power analysis of generalized linear mixed models
1026 by simulation. *Methods Ecol Evol* **7**:493–498. doi:10.1111/2041-210X.12504
- 1027 Grothe B, Pecka M. 2014. The natural history of sound localization in mammals--a story of neuronal
1028 inhibition. *Front Neural Circuits* **8**:116. doi:10.3389/fncir.2014.00116
- 1029 Hashimoto I. 1982. Auditory evoked potentials from the human midbrain: slow brain stem responses.
1030 *Electroencephalogr Clin Neurophysiol* **53**:652–657. doi:10.1016/0013-4694(82)90141-9
- 1031 Hyde M. 2008. Ontario Infant Hearing Program Audiologic Assessment Protocol Version 3.1.
- 1032 Jiang ZD, Wu YY, Wilkinson AR. 2009. Age-related changes in BAER at different click rates from neonates
1033 to adults. *Acta Paediatr Oslo Nor* **1992** **98**:1284–1287. doi:10.1111/j.1651-2227.2009.01312.x
- 1034 Kileny P, Paccioletti D, Wilson AF. 1987. Effects of cortical lesions on middle-latency auditory evoked
1035 responses (MLR). *Electroencephalogr Clin Neurophysiol* **66**:108–120. doi:10.1016/0013-
1036 4694(87)90180-5
- 1037 Kuznetsova A, Brockhoff PB, Christensen RHB. 2017. lmerTest Package: Tests in Linear Mixed Effects
1038 Models. *J Stat Softw* **82**:1–26. doi:10.18637/jss.v082.i13
- 1039 Larson E, McCloy D, Maddox R, Pospisil D. 2014. expyfun: Python experimental paradigm functions,
1040 version 2.0.0. Zenodo. doi:10.5281/zenodo.11640
- 1041 L'Engle M. 2012. A Wrinkle in Time: 50th Anniversary Commemorative Edition. Macmillan.
- 1042 Liberman MC, Epstein MJ, Cleveland SS, Wang H, Maison SF. 2016. Toward a Differential Diagnosis of
1043 Hidden Hearing Loss in Humans. *PLOS ONE* **11**:e0162726. doi:10.1371/journal.pone.0162726
- 1044 Maddox RK, Lee AKC. 2018. Auditory Brainstem Responses to Continuous Natural Speech in Human
1045 Listeners. *eNeuro* **5**. doi:10.1523/ENEURO.0441-17.2018
- 1046 Mesgarani N, David SV, Fritz JB, Shamma SA. 2009. Influence of Context and Behavior on Stimulus
1047 Reconstruction From Neural Activity in Primary Auditory Cortex. *J Neurophysiol* **102**:3329–3339.
1048 doi:10.1152/jn.91128.2008
- 1049 Miller L, IV BM, Bishop C. 2017. Frequency-multiplexed speech-sound stimuli for hierarchical neural
1050 characterization of speech processing. US20170196519A1.
- 1051 Møller AR, Jannetta PJ. 1983. Interpretation of brainstem auditory evoked potentials: results from
1052 intracranial recordings in humans. *Scand Audiol* **12**:125–133.
- 1053 Moore JK. 1987. The human auditory brain stem as a generator of auditory evoked potentials. *Hear Res*
1054 **29**:33–43. doi:10.1016/0378-5955(87)90203-6
- 1055 O'Sullivan AE, Lim CY, Lalor EC. 2019. Look at me when I'm talking to you: Selective attention at a
1056 multisensory cocktail party can be decoded using stimulus reconstruction and alpha power
1057 modulations. *Eur J Neurosci* **50**:3282–3295. doi:10.1111/ejn.14425

- 1058 Picton TW, Hillyard SA, Krausz HI, Galambos R. 1974. Human auditory evoked potentials. I. Evaluation of
1059 components. *Electroencephalogr Clin Neurophysiol* **36**:179–190. doi:10.1016/0013-4694(74)90155-
1060 2
- 1061 Polonenko MJ, Maddox RK. 2019. The Parallel Auditory Brainstem Response. *Trends Hear*
1062 **23**:2331216519871395. doi:10.1177/2331216519871395
- 1063 Prendergast G, Guest H, Munro KJ, Kluk K, Léger A, Hall DA, Heinz MG, Plack CJ. 2017. Effects of noise
1064 exposure on young adults with normal audiograms I: Electrophysiology. *Hear Res* **344**:68–81.
1065 doi:10.1016/j.heares.2016.10.028
- 1066 R Core Team. 2020. R: A language and environment for statistical computing. Vienna, Austria: R
1067 Foundation for Statistical Computing.
- 1068 Rasetshwane DM, Argenyi M, Neely ST, Kopun JG, Gorga MP. 2013. Latency of tone-burst-evoked
1069 auditory brain stem responses and otoacoustic emissions: Level, frequency, and rise-time effects. *J*
1070 *Acoust Soc Am* **133**:2803–2817. doi:10.1121/1.4798666
- 1071 Saiz-Alía M, Forte AE, Reichenbach T. 2019. Individual differences in the attentional modulation of the
1072 human auditory brainstem response to speech inform on speech-in-noise deficits. *Sci Rep* **9**:1–10.
1073 doi:10.1038/s41598-019-50773-1
- 1074 Saiz-Alia M, Reichenbach T. 2020. Computational modeling of the auditory brainstem response to
1075 continuous speech. *J Neural Eng*. doi:10.1088/1741-2552/ab970d
- 1076 Scott M. 2007. *The Alchemyst: the secrets of the immortal Nicholas Flamel, Book 1*. New York: Listening
1077 Library.
- 1078 Shore SE, Nuttall AL. 1985. High-synchrony cochlear compound action potentials evoked by rising
1079 frequency-swept tone bursts. *J Acoust Soc Am* **78**:1286–1295. doi:10.1121/1.392898
- 1080 Stapells DR, Oates P. 1997. Estimation of the Pure-Tone Audiogram by the Auditory Brainstem Response:
1081 A Review. *Audiol Neurotol* **2**:257–280. doi:10.1159/000259252
- 1082 Starr A, Hamilton AE. 1976. Correlation between confirmed sites of neurological lesions and abnormalities
1083 of far-field auditory brainstem responses. *Electroencephalogr Clin Neurophysiol* **41**:595–608.
1084 doi:10.1016/0013-4694(76)90005-5
- 1085 Teoh ES, Lalor EC. 2019. EEG decoding of the target speaker in a cocktail party scenario: considerations
1086 regarding dynamic switching of talker location. *J Neural Eng* **16**:036017. doi:10.1088/1741-
1087 2552/ab0cf1
- 1088 Winer JA. 2005. Decoding the auditory corticofugal systems. *Hear Res* **207**:1–9.
1089 doi:10.1016/j.heares.2005.06.007
- 1090

1091 FIGURE CAPTIONS

1092 **Figure 1.** Single subject and group average (bottom right) weighted-average auditory brainstem responses
1093 (ABR) to ~43 minutes of broadband peaky speech. Area for the group average shows ± 1 SEM.
1094 Responses were high-pass filtered at 150 Hz using a first order Butterworth filter. Waves I, III, and V of the
1095 canonical ABR are evident in most of the single subject responses (N = 22, 16, 22 respectively), and are
1096 marked by the average peak latencies on the average response.

1097 **Figure 2.** Comparison of auditory brainstem (ABR) and middle latency responses (MLR) to ~43 minutes
1098 each of unaltered speech and broadband peaky speech. (A) The average waveform to broadband peaky
1099 speech (blue) shows additional, and sharper, waves of the canonical ABR and MLR than the broader
1100 average waveform to unaltered speech (black). Responses were high-pass filtered at 30 Hz with a first
1101 order Butterworth filter. Areas show ± 1 SEM. (B) Comparison of peak latencies for ABR wave V (circles)
1102 and MLR waves N_a (downward triangles) and P_a (upward triangles) that were common between responses
1103 to broadband peaky and unaltered speech. Blue symbols depict individual subjects and black symbols
1104 depict the mean.

1105 **Figure 3.** Comparison of responses to 32 minutes each of male (dark blue) and female (light blue) narrated
1106 re-synthesized broadband peaky speech. (A) Average waveforms across subjects (areas show ± 1 SEM)
1107 are shown for auditory brainstem response (ABR) time lags with high-pass filtering at 150 Hz (top), and
1108 both ABR and middle latency response (MLR) time lags with a lower high-pass filtering cutoff of 30 Hz
1109 (bottom). (B) Histograms of the correlation coefficients between responses evoked by male- and female-
1110 narrated broadband peaky speech during ABR (top) and ABR/MLR (bottom) time lags. Solid lines denote
1111 the median and dotted lines the inter-quartile range. (C) Comparison of ABR (top) and MLR (bottom) wave

1112 peak latencies for individual subjects (gray) and the group mean (black). ABR and MLR responses were
1113 similar to both types of input but are smaller for female-narrated speech, which has a higher glottal pulse
1114 rate. Peak latencies for female-evoked speech were delayed during ABR time lags but faster for early MLR
1115 time lags.

1116 **Figure 4.** Comparison of responses to ~43 minutes of male-narrated multiband peaky speech. (A) Average
1117 waveforms across subjects (areas show ± 1 SEM) are shown for each band (colored solid lines) and for the
1118 common component (dot-dash gray line, same waveform replicated as a reference for each band), which
1119 was calculated using 6 false pulse trains. (B) The common component was subtracted from each band's
1120 response to give the frequency-specific waveforms (areas show ± 1 SEM), which are shown with high-pass
1121 filtering at 30 Hz (solid lines) and 150 Hz (dashed lines). (C) Mean \pm SEM peak latencies for each wave
1122 decreased with increasing band frequency. Numbers of subjects with an identifiable wave are given for
1123 each wave and band.

1124 **Figure 5.** Comparison of responses to ~43 minutes of male-narrated peaky speech in the same subjects.
1125 Average waveforms across subjects (areas show ± 1 SEM) are shown for broadband peaky speech (blue)
1126 and for the summed frequency-specific responses to multiband peaky speech with the common component
1127 added (red), high-pass filtered at 150 Hz (left) and 30 Hz (right). Regressors in the deconvolution were
1128 pulse trains.

1129 **Figure 6.** Comparison of responses to 32 minutes each of male- and female-narrated re-synthesized
1130 multiband peaky speech. (A) Average frequency-specific waveforms across subjects (areas show ± 1 SEM;
1131 common component removed) are shown for each band in response to male- (dark red lines) and female-
1132 narrated (light red lines) speech. Responses were high-pass filtered at 30 Hz (left) and 150 Hz (right) to
1133 highlight the MLR and ABR respectively. (B) Correlation coefficients between responses evoked by male-
1134 and female-narrated multiband peaky speech during ABR/MLR (left) and ABR (right) time lags for each
1135 frequency band. Black lines denote the median. (C) Mean \pm SEM peak latencies for male- (dark) and
1136 female- (light) narrated speech for each wave decreased with increasing frequency band. Numbers of
1137 subjects with an identifiable wave are given for each wave, band and narrator. Lines are given a slight
1138 horizontal offset to make the error bars easier to see.

1139 **Figure 7.** Comparison of responses to ~60 minutes each of male- and female-narrated dichotic multiband
1140 peaky speech with standard audiological frequency bands. (A) Average frequency-specific waveforms
1141 across subjects (areas show ± 1 SEM; common component removed) are shown for each band for the left
1142 ear (dotted lines) and right ear (solid lines). Responses were high-pass filtered at 30 Hz. (B) Left-right ear
1143 correlation coefficients (top, averaged across gender) and male-female correlation coefficients (bottom,
1144 averaged across ear) during ABR time lags (0–15 ms) for each frequency band. Black lines denote the
1145 median. (C) Mean \pm SEM wave V latencies for male- (dark red) and female-narrated (light red) speech for
1146 the left (dotted line, cross symbol) and right ear (solid line, circle symbol) decreased with increasing
1147 frequency band. Lines are given a slight horizontal offset to make the error bars easier to see.

1148 **Figure 8.** Cumulative proportion of subjects who have responses with ≥ 0 dB SNR as a function of
1149 recording time. Time required for unaltered (black) and broadband peaky speech (dark blue) of a male
1150 narrator are shown for 22 subjects in the left plot, and for male (dark blue) and female (light blue)
1151 broadband peaky speech is shown for 11 subjects in the right plot. Solid lines denote SNRs calculated
1152 using variance of the signal high-pass filtered at 30 Hz over the ABR/MLR interval 0–30 ms, and dashed
1153 lines denote SNR variances calculated on signals high-pass filtered at 150 Hz over the ABR interval 0–15
1154 ms. Noise variance was calculated in the pre-stimulus interval –480 to –20 ms.

1155 **Figure 9.** Cumulative proportion of subjects who have frequency-specific responses (common component
1156 subtracted) with ≥ 0 dB SNR as a function of recording time. Acquisition time was faster for male (left) than
1157 female (right) narrated multiband peaky speech with (A) 4 frequency bands presented diotically, and with
1158 (B) 5 frequency bands presented dichotically (total of 10 responses, 5 bands in each ear). SNR was
1159 calculated by comparing variance of signals high-pass filtered at 150 Hz across the ABR interval of 0–15
1160 ms to variance of noise in the pre-stimulus interval –480 to –20 ms.

1161 **Figure 10.** The range of lags can be extended to allow early, middle and late latency responses to be
1162 analyzed from the same recording to broadband peaky speech. Average waveforms across subjects (areas
1163 show ± 1 SEM) are shown for responses measured to 32 minutes of broadband peaky speech narrated by
1164 a male (dark blue) and female (light blue). Responses were high-pass filtered at 30 Hz using a first order

1165 Butterworth filter, but different filter parameters can be used to focus on each stage of processing.
1166 Canonical waves of the ABR, MLR and LLR are labeled for the male-narrated speech. Due to adaptation,
1167 amplitudes of the late potentials are smaller than typically seen with other stimuli that are shorter in
1168 duration with longer inter-stimulus intervals than our continuous speech. Waves I and III become more
1169 clearly visible by applying a 150 Hz high-pass cutoff.

1170 **Figure 11.** Unaltered speech waveform (top left) and spectrogram (top right) compared to re-synthesized
1171 broadband peaky speech (middle left and right) and multiband peaky speech (bottom left and right).
1172 Comparing waveforms shows that the peaky speech is as “click-like” as possible, while comparing the
1173 spectrograms shows that the overall spectrotemporal content that defines speech is basically unchanged
1174 by the re-synthesis. A naïve listener is unlikely to notice that any modification has been performed, and
1175 subjective listening confirms the similarity. Yellow/lighter colors represent larger amplitudes than
1176 purple/darker colors in the spectrogram. See supplementary files for audio examples of each stimulus type
1177 for both narrators.

1178 **Figure 12.** Relative mean-squared magnitude in decibels of multiband peaky speech with 4 filter bands
1179 (left) and 5 filter bands (right) for male-(blue) and female-(orange) narrated speech. The full audio
1180 comprises unvoiced and re-synthesized voiced sections, which was presented to the subjects during the
1181 experiments. The other bands reflect the relative magnitude of the voiced sections (voiced only), and each
1182 filtered frequency band.

1183 **Figure 13.** Spectral coherence of pulse trains for multiband peaky speech narrated by a male (left) and
1184 female (right). Spectral coherence was computed across 1 s slices from 60 unique 64 s multiband peaky
1185 speech segments (3,840 total slices) for each combination of bands. Each light gray line represents the
1186 coherence for one band comparison. There were 45 comparisons across the 10-band (audiological)
1187 speech used in experiment 3 (5 frequency bands x 2 ears). Pulse trains (i.e., the input stimuli, or
1188 regressors, for the deconvolution) were frequency-dependent (coherent) below 72 Hz for the male
1189 multiband speech and 126 Hz for the female multiband speech.

1190 **Figure 14.** Octave band filters used to create re-synthesized broadband peaky speech (left, blue), diotic
1191 multiband peaky speech with 4 bands (middle, red), and dichotic multiband peaky speech using 5 bands
1192 with audiological center frequencies (right, red). The last band (2nd, 5th, 6th respectively, black line) was
1193 used to filter the high-frequencies of unaltered speech during mixing to improve the quality of voiced
1194 consonants. The designed frequency response using trapezoids (top) were converted into the time-domain
1195 using IFFT, shifted and Nuttall windowed to create impulse responses (middle), which were then used to
1196 assess the actual frequency response by converting into the frequency domain using FFT (bottom).

1197 **Figure 15.** Left: A segment of broadband peaky speech stimulus (top) and the corresponding glottal pulse
1198 train (bottom) used in calculating the broadband peaky speech response. Right: An example broadband
1199 peaky speech response from a single subject. The response shows ABR waves I, III, and V at ~3, 5, 7 ms
1200 respectively. It also shows later peaks corresponding to thalamic and cortical activity at ~17 and 27 ms
1201 respectively.

1202