

1 **Gaussian-Enveloped Tones (GET): a vocoder that can simulate pulsatile**
2 **stimulation in cochlear implants**

3
4 Qinglin Meng^{a)} and Huali Zhou^{b)}

5 *Acoustics Laboratory, School of Physics and Optoelectronics, South China University of*
6 *Technology, Guangzhou, Guangdong 510641, China*

7 Thomas Lu and Fan-Gang Zeng^{c)}

8 *Center for Hearing Research, Department of Otolaryngology –Head and Neck Surgery,*
9 *University of California, Irvine, California 92697, USA*

10

11 Running title: Acoustic simulations of cochlear implants

12

13 This paper is submitted for consideration for a special issue on Reconsidering Classic Ideas in
14 Speech Communication on the Journal of the Acoustical Society of America.

15

16 Submission date: Dec. 19, 2021

17 Revision date: May 5, 2022

18

^{a)} Electronic mail: mengqinglin@scut.edu.cn. ORCID: 0000-0003-0544-1967.

^{b)} Also at: College of Electronics and Information Engineering, Shenzhen University, Shenzhen, Guangdong, 518060, China

^{c)} Electronic mail: fzeng@uci.edu. ORCID: 0000-0002-4325-2780

19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

ABSTRACT

Acoustic simulations of cochlear implants (CIs) allow for studies of perceptual performance with minimized effects of large CI individual variability. Different from conventional simulations using continuous sinusoidal or noise carriers, the present study employs Gaussian-enveloped tones (GETs) to simulate pulsatile stimulation in modern CIs. Subject to the time-frequency uncertainty principle, the GET has a well-defined tradeoff between its duration and bandwidth. Two types of GET vocoders were implemented and evaluated in normal-hearing listeners. In the first implementation, constant 100-Hz GETs were used to minimize within-channel temporal overlap while different GET durations were used to simulate electric channel interaction. This GET vocoder could produce vowel and consonant recognition similar to actual CI performance. In the second implementation, 900-Hz/channel pulse trains were directly mapped to 900-Hz GET trains to simulate a widely-used *n-of-m* processing strategy, or the Advanced Combination Encoder. The simulated and actual implant performance of speech in noise recognition was similar in terms of the overall trend, absolute mean scores, and standard deviations. **The present results suggest that the pulsatile GETs can be used as alternative vocoders to simulate speech perception with modern CIs.**

JASA abstract word limit: 200.

Current word number: 184.

38 I. INTRODUCTION

39 Vocoders as a means of speech synthesis have a long and rich history. At the 1939 New
40 York World's Fair, Homer Dudley of Bell Labs demonstrated his vocoder invention that could
41 "remake speech" automatically and instantaneously (18-ms delay) by controlling energy in 10
42 frequency bands (from 0 to 3000 Hz) that contained either buzz-like tone or hiss-like noise
43 carriers (Dudley, 1939). He later realized that the vocoder could be used in synthesizing speech,
44 and transformed in various ways to study the relative contributions of fundamental parameters
45 in speech synthesis and recognition. He found that good intelligibility can be achieved by
46 controlling "only low syllabic frequencies of the order of 10 cycles per second", whereas the
47 emotional content of speech can be controlled by altering the frequency of the buzzing tones.

48 The early multi-channel CIs followed Dudley's original vocoder idea closely by
49 extracting and delivering speech fundamental frequency (F0) in the form of electric pulse rate
50 and one or two formants (F2 or F1/F2) in the form of electrode position (Tong et al., 1980;
51 Skinner et al., 1991). The speech understanding of the early CIs was relatively low (<50%
52 correct for sentence recognition in quiet), due not only to crude F0 and formant extraction
53 methods (i.e., zero-crossing) at that time, but, more importantly, to complicated interactions
54 between sound frequency and electric pitch, for example, individual variability in electrode
55 insertion angle or depth, cochlear vs. ganglion cell tonotopic organization, current spread, and
56 nerve survival. These interactions make accurate F0 and formant representation difficult if not
57 impossible even if both F0 and formants can be exactly extracted by today's algorithms. As a
58 result, contemporary CIs have abandoned the F0 and formant extraction method but adopted
59 speech processing strategies that extract band-specific temporal envelopes from 8-24 frequency
60 bands. The envelopes are used to amplitude modulate a continuous, but fixed, high-rate (at least

61 two to four times the highest envelope frequency) pulse train, which is then delivered to a
62 corresponding electrode in an interleaved fashion in which no two electrodes fire simultaneously
63 ([Wilson et al., 1991](#); [Skinner et al., 2002](#)). These advances in multi-channel CIs have produced
64 70-80% correct sentence recognition in quiet, which is sufficient for an average user to carry on
65 a conversation without lipreading ([Zeng et al., 2008](#)).

66 Acoustic simulations of CIs have been developed and widely used ([Svirsky et al., 2021](#))
67 for at least three reasons. First, acoustic simulations minimize the effect of large CI individual
68 variability (e.g., cognitive differences, demographic variables, and electrode-neuron interface),
69 which may confound or mask the relative importance of speech processing parameters, e.g.,
70 [Skinner et al. \(2002\)](#). Second, acoustic simulations allow the evaluation of relative contributions
71 of different cues to auditory and speech perception, e.g., [Xu et al. \(2005\)](#); [Singh et al. \(2009\)](#).
72 Third, acoustic simulations allow a normal-hearing listener to appreciate the quality of CI
73 processing and the degree of difficulty facing a typical CI user.

74 Traditionally, acoustic simulations of CIs have used either noise- ([Shannon et al., 1995](#))
75 or sinusoid-excited ([Dorman et al., 1997](#)) vocoders. In these vocoders, the noise or sinusoid
76 simulates the electric pulse train, while the number of frequency bands and their overlaps
77 simulate the limited number of electrodes and their current spread, e.g., [Shannon et al. \(1998\)](#).
78 A significant drawback of these traditional vocoder models is the lack of simulation of the
79 pulsatile nature of CI electric stimulation. Several studies have attempted to develop acoustic
80 models that simulate pulsatile electric stimulation, such as filtered noise bursts ([Blamey et al.,](#)
81 [1984a](#); [Blamey et al., 1984b](#)), filtered harmonic complex tones ([Deeks and Carlyon, 2004](#)), and
82 pulse-spread harmonic complexes ([Hilkhuyzen and Macherey, 2014](#); [Mesnildrey et al., 2016](#)).
83 However, there are limitations to those methods in simulating some important features in

84 modern CIs. First, these vocoders cannot simulate the discrete nature of pulsatile stimulation on
85 a pulse-by-pulse basis. Second, they do not allow independent manipulation of the overlap
86 between spectral and temporal representation. Third, it is difficult for vocoders using continuous
87 carriers to simulate some CI speech processing strategies, e.g., *n-of-m*, in which the low-energy
88 bands are abandoned to produce temporally separated envelopes.

89 Here we identified the Gabor atom ([Gabor, 1947](#)), also known as the Gaussian-enveloped
90 tone (GET), as a means of simulating the essential features of modern CI processing as discussed
91 above. The GET has been used to study a wide range of auditory phenomena in normal hearing
92 or hearing-impaired listeners, e.g., temporal gap detection ([Schneider et al., 1994](#); [Trehub et al.,](#)
93 [1995](#)), intensity discrimination ([Baer et al., 1999](#); [van Schijndel et al., 1999](#); [Baer et al., 2001](#);
94 [Nizami et al., 2001](#)), simultaneous and non-simultaneous masking ([Laback et al., 2011](#); [Laback](#)
95 [et al., 2013](#)), interaural timing difference (ITD) ([Buell and Hafter, 1988](#)), and cortical encoding
96 of pulsatile stimulation ([Lu and Wang, 2000](#); [Lu et al., 2001](#); [Johnson et al., 2017](#)). More recently,
97 GET train has been used to simulate some basic tasks on binaural hearing with CIs, e.g., sound
98 localization ([Goupell et al., 2010](#); [Jones et al., 2014](#)), lateralization ([Ehlers et al., 2016](#)), binaural
99 masking level differences ([Lu et al., 2010](#)), temporal weighting of ITD and interaural level
100 difference (ILD) ([Brown and Stecker, 2010](#)), effects of electrode place mismatch on binaural
101 cues ([Goupell et al., 2013](#); [Kan et al., 2013](#)), and effects of temporal quantization on ITD
102 discrimination ([Dieudonne et al., 2020](#)).

103 In signal processing, due to the time-frequency uncertainty principle (also referred to as
104 the Gabor limit), the duration and bandwidth of a signal cannot be independently controlled, and
105 their product is no lower than a limit, which is reachable only by GETs (or say Gabor atoms)

106 (Gabor, 1947; Feichtinger and Strohmer, 1998; Gardner and Magnasco, 2006). This is an
107 important reason why most of the above-mentioned psychoacoustic studies use GETs as stimuli.

108 **However, the performance of GET-based vocoders in simulating speech perception with**
109 **CIs has not been investigated.** In much of the existing literature, conventional channel-vocoders
110 with eight channels using continuous noise or sine-wave carriers were used to replicate the sound
111 of 12-24 channel CIs. The main reason is the performance of eight-channel vocoders in normal-
112 hearing listeners usually matches the better performance of actual CI users (Winn and Nelson,
113 2021).

114 This study introduces a novel GET vocoder and demonstrates its potential for simulating
115 CI speech perception. In the following sections, the implementation and theory of the proposed
116 GET vocoders are introduced in detail; then two separate experiments of speech perception, each
117 with a different type of GET vocoder, are used to demonstrate the potential of the novel pulsatile
118 vocoders on CI speech perception simulation. Specifically, the first GET (Lu et al., 2007; Goupell
119 et al., 2010) is a naïve type using non-interleaved 100-pps (pulse per second) GET trains as
120 carriers to study the effect of current interaction among channels. The second GET (Meng et al.,
121 2018; Kong et al., 2019) is an advanced type that can directly map individual electric pulses from
122 a clinical n -of- m strategy with 900-pps pulse rate into an acoustic GET. In this way, any CI
123 electrodiagram (not limited to the selected strategy) can be directly transformed into a vocoded
124 sound. Such direct transformation can simulate not only pulsatile timing cues but also many other
125 features of CI electric stimuli (e.g., amplitude compression and maxima selection).

126 The pulsatile GET vocoder can replicate the temporal (pulsatile), intensity (compressed
127 and quantized), and spectral (maxima-selected) features of an actual CI strategy. Furthermore,
128 current spread at individual electrodes can be simulated by changing the GET bandwidth through

129 the pulse duration parameter. We hypothesized that the GET vocoder could be an alternative
130 vocoder model to simulate speech perception with CIs. Nevertheless, the uncertainty principle
131 imposes unavoidable physical constraints on the time-frequency tradeoff, which might limit the
132 performance of the pulsatile simulation and should be carefully controlled.

133 II. GET THEORY AND VOCODER ALGORITHMS

134 A. GET Theory

135 A Gaussian function is symmetrical in the time domain:

$$136 \quad g_{env}(t) = ae^{-\frac{\pi(t-t_0)^2}{2\sigma^2}} \quad (1)$$

137 where a determines the function's maximum amplitude, t_0 the maximum amplitude's
138 temporal position, and σ the effective duration or $D = \sqrt{2}\sigma$, at which the amplitude is 6.82-dB
139 down from the maximum amplitude (Baer et al., 1999). Its Fourier transform is:

$$140 \quad G_{env}(f) = \sqrt{2}a\sigma \cdot e^{-2\pi(\sigma f)^2} \cdot e^{-j2\pi f t_0} \quad (2)$$

141 The shape of its amplitude spectrum, $\sqrt{2}a\sigma \cdot e^{-2\pi(\sigma f)^2}$, is also a Gaussian function with an
142 effective bandwidth being $B = \frac{1}{\sqrt{2}\sigma}$ between the 6.82-dB down cutoff frequencies.

143 The effective duration (D) and the effective bandwidth (B) can be traded:

$$144 \quad D \cdot B = 1 \quad (3)$$

145 meaning that increasing the duration will narrow the bandwidth and vice versa.

146 Acoustic simulation of a single electric pulse in a frequency channel can be generated by
147 multiplying the above Gaussian function by a sinusoidal carrier:

$$148 \quad s(t) = g_{env}(t) \cdot \sin(2\pi f_c t + \varphi_0) = ae^{-\frac{\pi(t-t_0)^2}{2\sigma^2}} \cdot \sin(2\pi f_c t + \varphi_0) \quad (4)$$

149 where $s(t)$ has the same effective duration and effective bandwidth as $g_{env}(t)$ except for
 150 changing the center frequency from 0 to f_c , and φ_0 is an initial phase.

151 Fig. 1 illustrates both waveform (a) and spectrum (b) of a unit-amplitude Gaussian-
 152 enveloped single pulse (i.e., $a = 1$ in Eq. 4). The carrier frequency f_c is 5 kHz. The 6.82-dB
 153 cutoff point (corresponding to $D = \sqrt{2}\sigma$) with an amplitude of 0.456 in Fig. 1 was derived by
 154 substituting $t = t_1 = \frac{D}{2} + t_0 = \frac{\sqrt{2}}{2}\sigma + t_0$ into Eq. (1), i.e.,

$$155 \quad g_{env}(t_1) = e^{-\frac{\pi(\frac{\sqrt{2}}{2}\sigma)^2}{2\sigma^2}} = e^{-\frac{\pi}{4}} \approx 0.456 \quad (5)$$

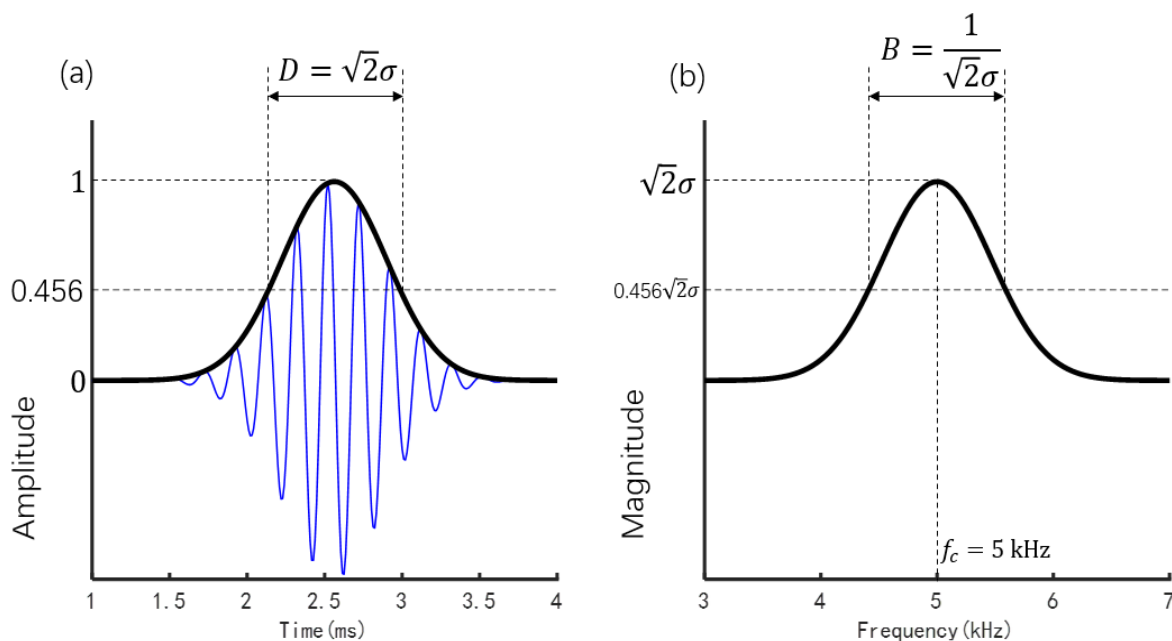


FIG. 1. (Color online) A unit-amplitude single pulse with Gaussian-shaped envelope (black line) in both the time (a) and frequency (b) domains. The carrier frequency is 5 kHz (the blue waveform in the left panel and the frequency with maximum amplitude in the right panel). The σ equals to $3/f_c = 0.6$ ms in Eq. (1), producing an effective duration of 0.85 ms and an effective bandwidth of 1.2 kHz.

156 Using the GET defined by Eq. 4, the change of amplitude and timing of an electric pulse can be
 157 simulated by manipulating a and t_0 respectively. Acoustic simulation of a continuous electric
 158 pulse train can be constructed by periodically repeating $s(t)$ or convolution of the electric pulse
 159 train and a GET.

160 Different from the CI electric pulses with constant duration at the order of tens of
 161 microseconds, the GET duration should be much longer to contain at least several (l) periods
 162 (e.g., $l = 2, 3, \text{ or } 4$) of the tone carrier. Therefore, the carrier period or frequency will determine
 163 the lower limits of the GET duration. The three lines in the two panels of Fig. 2 illustrate the
 164 dependent relationship between the GET duration (bandwidth), pulse rate, and carrier frequency,
 165 when $\sigma = \frac{l}{f_c} = \frac{2}{f_c}, \frac{3}{f_c}$, and $\frac{4}{f_c}$, respectively. The GET effective bandwidth equals in value to the
 166 maximum pulse rate that can be transmitted without obvious temporal interaction between
 167 neighboring GETs. Here the GET duration threshold for the “obvious temporal interaction” was
 168 defined as the effective duration of GET, i.e., $D = \sqrt{2}\sigma$. Increasing the duration (i.e., larger σ)
 169 can decrease the bandwidth with the maximum rate decreasing correspondingly.

170 At frequency bands with high carrier frequencies above ~ 2.5 kHz ($f_c = \frac{l}{\sigma} = l\sqrt{2}B = \sqrt{2} \cdot$
 171 $900l \approx 2546, 3818, \text{ and } 5091$ Hz for $l = 2, 3, \text{ and } 4$, respectively), a conventional pulse rate of
 172 900 pps could be simulated without obvious temporal interaction between neighboring GETs.

173 For carrier frequencies within the middle-frequency range around 2 kHz, the 900 pps is still
 174 possible to simulate, but neighboring GETs have moderate temporal interaction. The amplitude
 175 of the crossing point of neighboring GETs at a 2 kHz carrier would be

$$176 \quad 20 \lg e^{-\frac{\pi(t-t_0)^2}{2\sigma^2}} = 20 \lg e^{-\frac{\pi\left(\frac{1}{2 \times 900}\right)^2}{2\left(\frac{l}{f_c}\right)^2}} = 20 \lg e^{-\frac{\pi}{2}\left(\frac{10}{9l}\right)^2} \quad (6)$$

177 whose values are -4.21 , -1.87 , and -1.05 dB (relative to the maximum amplitude) for $l = 2, 3,$
 178 and 4, respectively. For a low-frequency carrier, the pulsatile feature for simulation of individual
 179 electric pulses cannot be guaranteed due to temporal interactions between neighboring GETs.

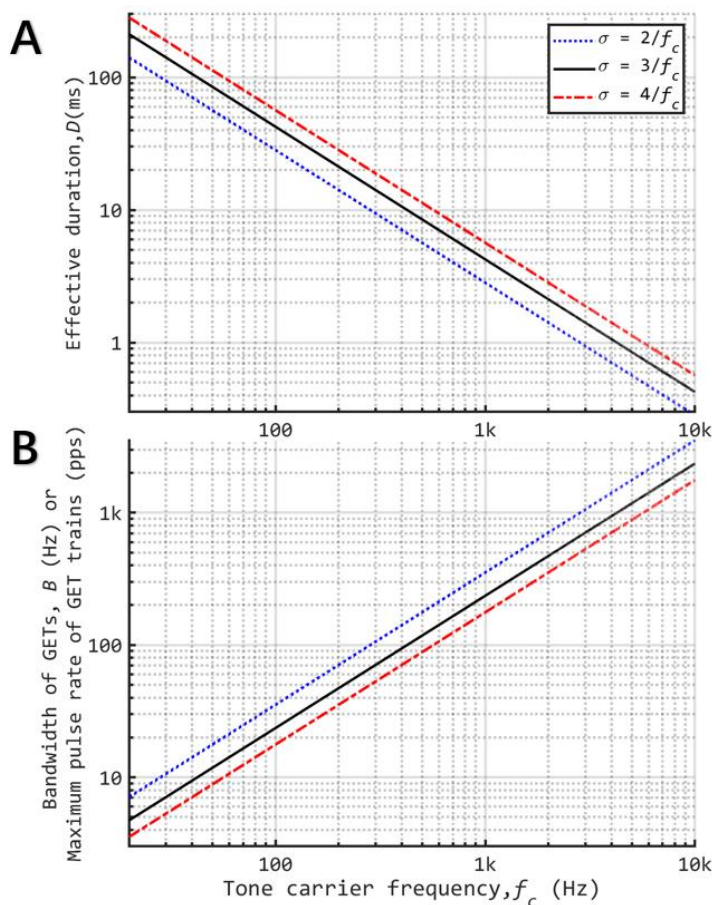


FIG 2. (Color online) The relationship between the tone carrier frequency and the effective duration $D = \sqrt{2}\sigma$ (see Panel A) or effective bandwidth $B = 1/D$ (see Panel B) of Gaussian-enveloped tones (GETs). All axes are logarithmically scaled. The σ was assumed to be $2/f_c$, $3/f_c$, or $4/f_c$ to demonstrate the effects of different duration of GETs. For certain combinations of f_c and σ , the maximum GET rate that can be transmitted with no temporal interaction between neighboring GETs is $1/D$, which equals in value to the effective bandwidth in Panel B.

180 **The temporal envelopes delivered in electric speech stimuli are often temporally separated**
 181 **across channels in many CI strategies, as** nature speech contains natural gaps within each
 182 channel of signal between syllables, and frame-wise low power bands are temporarily
 183 abandoned resulting from the maxima selection for n -of- m strategies. Additionally, envelope
 184 energies lower than the compression threshold level (or T level) are not represented in electric
 185 stimuli (i.e., no stimulation) in some strategies. For the temporally separated electric stimuli
 186 within each channel, GET carriers can better represent temporal separation features as well as
 187 CI compression (limited electric dynamic range), both of which are often omitted in
 188

189 conventional noise and sine-wave vocoders. The temporal separation features may be simulated
190 in all channels, and the low carrier frequency limit f_{c_low} is mainly determined by the duration
191 d_{gap} of each gap in the pulse trains:

$$192 \quad f_{c_low} = \frac{l}{\sigma} = \frac{\sqrt{2}l}{D_{max}} = \frac{\sqrt{2}l}{d_{gap}} \quad (7)$$

193 where D_{max} is the maximum possible GET duration, which equals the gap duration.

194 Current (or spectral) spread was acknowledged to be an important issue influencing the
195 frequency resolution of CIs (Mehta et al., 2020). For a single GET (defined by Eq. 4), its
196 bandwidth is determined by its duration due to the time-frequency uncertainty principle.
197 Therefore, it is possible to simulate CI current spread by manipulating the GET duration,
198 meaning the pulsatile timing feature and the current spread cannot be independently manipulated.

199 In short, the GETs can simulate and manipulate five important parameters of CI processing
200 or stimulation: (1) pulse rate by changing the period of pulse generation, (2) temporal envelope
201 (including its compression and quantization) by changing the amplitude of individual GETs in a
202 pulse train within a channel, (3) spectral envelope by changing the GET amplitude across
203 channels, (4) place of excitation by changing the carrier tone frequency, and (5) spread of
204 excitation by changing the effective bandwidth in GETs. The precise manipulation of these five
205 important parameters allows acoustic simulation of modern CIs using pulsatile electric
206 stimulation. The limitations from the dependent relationships between duration, bandwidth, and
207 carrier frequency of GETs are discussed above and should be taken into consideration during
208 algorithm design and experiments of CI simulations with GETs.

209 B. Vocoder Algorithm Frameworks

210 Fig. 3A shows the conventional acoustic simulation of CI using either noise (Shannon et al.,
211 1995) or sine-wave vocoders (Dorman et al., 1997). The output filters can be used to control the

212 current spread, but no temporal separation feature (e.g., pulsatile timing and temporally separated
 213 envelope) can be simulated.

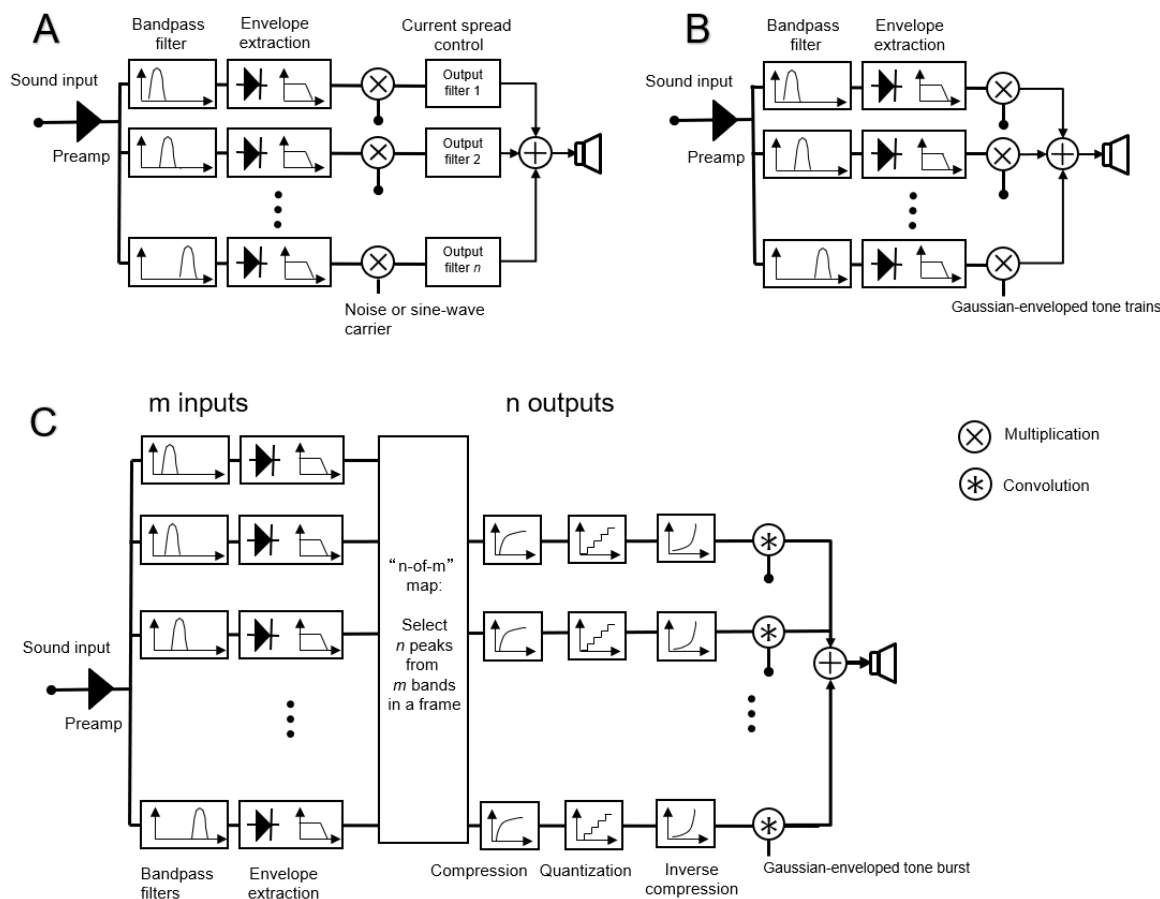


FIG. 3. Block diagrams of conventional channel vocoder (A), the first (B) and second (C) types of GET vocoders. The pulsatile vocoders are using GETs as carriers (the first type; used in Exp. 1) or using a single GET as an impulse response (the second type; used in Exp. 2). The front-end pre-emphasis, bandpass filter, and envelope extraction can be implemented either in the temporal or spectral domain.

214
 215 The first GET vocoder was proposed by [Lu et al. \(2007\)](#) (see Fig. 3B) and subsequently used
 216 in a sound localization study ([Goupell et al., 2010](#)). As a naïve implementation, this approach
 217 replaces the conventional continuous carriers with pulsatile GET carriers. To demonstrate the
 218 effects of current interaction realized by different GET durations, vowel and consonant perception
 219 with non-interleaved 100-pps GET carriers was measured in Experiment 1 (Section III).

220 The second GET vocoder was proposed by [Meng et al. \(2018\)](#) (see Fig. 3C). Compared to
221 the naïve implementation of the first type, the second GET vocoder hypothesized that a direct
222 mapping from individual CI electric pulses to individual GET acoustic pulses could transmit
223 similar speech information in both modes of CI and GET simulation. The implementation
224 framework of the second GET vocoder considers a common feature of temporal-frame-based n -
225 of- m selection in some CI processing strategies. The n -of- m selection means n maximum envelope
226 values are selected out of the envelope values from the m input channels within a given time
227 window. In this framework, the amplitude compression and quantization widely used in modern
228 CIs can also be simulated. In Experiment 2 (Section IV), sentence intelligibility tests were carried
229 out to demonstrate the feasibility of GET simulation on speech perception with the advanced
230 combination encoder (ACE) strategy, which is a typical n -of- m strategy and has a default pulse
231 rate of 900 pps.

232 The front-end processing stages of the three methods in Fig.3 share the same blocks of band-
233 pass filters and envelope extraction, e.g., in a traditional temporal envelope-based continuous
234 interleaved sampling (CIS) ([Wilson et al., 1991](#)) or ACE strategy ([Vandali et al., 2000](#)). Details
235 about the implementations of the two types of GET vocoders are provided in the following two
236 experiment sections.

237 **III. EXPERIMENT 1: SIMULATION OF CURRENT SPREAD**

238 **A. Rationale**

239 Experiment 1 was designed to study vowel and consonant speech perception with the first
240 type of GET vocoder ([Lu et al., 2007](#); [Goupell et al., 2010](#)) using non-interleaved GET carriers
241 (where the GET centers for all channels are in alignment with each other in each frame). The

242 interleaved sampling feature of modern CI strategies was not considered. A low pulse rate of
243 100 pps, which is much lower than the standard clinical rate (e.g., 900 pps or faster), was used
244 in this experiment to minimize the within-channel inter-pulse temporal interaction. The primary
245 purpose of this experiment is to examine the effects of current spread stimulated by manipulating
246 the GET duration based on the uncertainty principle.

247 There is a significant difference in simulating the spread of excitation between the
248 conventional vocoder (Shannon et al., 1995; Dorman et al., 1997) and the GET implementation
249 (Lu et al., 2007). In the conventional simulation, the spread of excitation is manipulated by
250 changing the filter type and the bandwidth of the synthesis band-pass filters at the vocoder output
251 stage (Croghan and Smith, 2018). For the GETs, the spread of excitation is manipulated by
252 increasing or decreasing the Gaussian tone duration, which produces a corresponding change in
253 narrowing or widening the spectral bandwidth for each pulse.

254 **B. Methods**

255 Five vocoders were used: three conventional vocoders - sine-wave, noise-separate, and
256 noise-spread (Fig. 3A) - and two proposed vocoders incorporating the GET simulation -GET-
257 separate and GET-spread (Fig. 3B).

258 *Analysis processing of all five vocoders:* The analysis filter banks consist of N band-pass
259 filters (4th order Butterworth). The frequency spacing for cutoffs for the filter bank was defined
260 in the range of [80, 7999] Hz according to a Greenwood map (Greenwood, 1990) (See Tab. I).
261 The filtered signals were half-wave rectified and low-pass filtered (50 Hz 4th order Butterworth)
262 to extract the envelope for each channel. This 50-Hz cutoff requires, in theory, at least a 100-Hz
263 carrier to avoid aliasing.

264

265 **TABLE I. Cutoff frequencies** of the band-pass filters in Exp. 1 according to a Greenwood map

Band number	Cutoff frequencies (Hz)
2	80, 1250, and 7999
4	80, 424, 1250, 3234, and 7999
8	80, 215, 424, 748, 1250, 2028, 3234, 5103, and 7999
16	80, 140, 215, 308, 424, 568, 748, 972, 1250, 1597, 2028, 2565, 3234, 4067, 5103, 6393, and 7999
32	80, 108, 140, 176, 215, 259, 308, 363, 424, 492, 568, 653, 748, 854, 972, 1103, 1250, 1414, 1597, 1801, 2028, 2282, 2565, 2881, 3234, 3628, 4067, 4556, 5103, 5713, 6393, 7152, and 7999

266 *Synthesis processing for the conventional vocoders:* For the sine-wave vocoder, a sine
267 wave with a frequency centered at the corresponding analysis filtering band was used as the
268 carrier. For the noise-separate vocoder, band-pass noise carriers were generated by passing
269 white noise through filters that were the same as the analysis filters. The noise-separate vocoder
270 provides upper-bound performance with a minimum of simulated electrode interaction. For the
271 noise-spread vocoder, low-pass filters (4th order Butterworth) were used to pass white noise for
272 generating low-pass noise carriers. The cutoff frequencies of the low-pass filters were the same
273 as the upper cutoff frequencies of the analysis filters. The signal carriers in each band were
274 corresponding low-pass noises. Low-pass filters were chosen to represent severe interactions
275 between channels (especially on the low-frequency side), and provide a lower bound of
276 performance with simple manipulation. For the two noise vocoders, after modulating each
277 channel of filtered noise with the channel envelope, the output was filtered again to band-limit
278 each channel. The band-limiting filters are the same as those used for the noise carrier
279 generation. The final vocoded signal was synthesized by summing all channels.

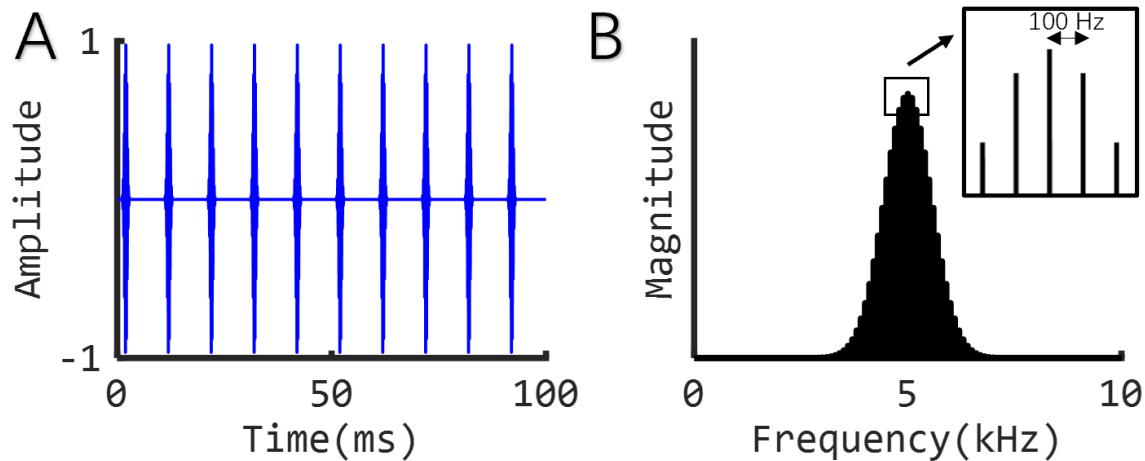


FIG. 4. (Color online) A 100-Hz pulse train, repeating a single pulse every 10 ms, in both the time (left panel) and frequency (right panel) domains. The parameters of the individual pulses are the same as those in Fig. 1.

280 *Synthesis processing for the GET vocoders:* For the GET vocoders, instead of modulating
281 a filtered noise signal at the synthesis stage, the envelope in each channel modulates the
282 amplitude of a GET train. Fig. 4 shows a 100-Hz pulse train, repeating the single pulse every 10
283 ms. The pulse train’s spectral envelope remains the same as the single pulse but its spectral fine
284 structure becomes discrete with 100-Hz spacing (in this case, the maximum-amplitude
285 frequency is 5 kHz with symmetrically decreasing-amplitude components at 4.9, 4.8, 4.7... and
286 5.1, 5.2, 5.3... kHz, respectively, see inset in the right panel). For the GET-separate vocoder,
287 $D = \sqrt{2}\sigma = 7.0$ ms, while for the GET-spread vocoder, $D = \sqrt{2}\sigma = 1.2$ ms. Because the first
288 experiment focused on the spread of excitation, the pulses among all channels were
289 synchronized, meaning that the “interleaved sampling” feature was not simulated.

290 CI stimulation was simulated using the above five different vocoders, i.e., sine-wave,
291 noise-separate, noise-spread, GET-separate, and GET-spread. The numbers of channels tested
292 were 2, 4, 8, 16, and 32. There were 12 medial vowels and 14 medial consonants in the vowel
293 and consonant tests, respectively. Fig. 5 provides an example of 16-channel vocoded stimuli for

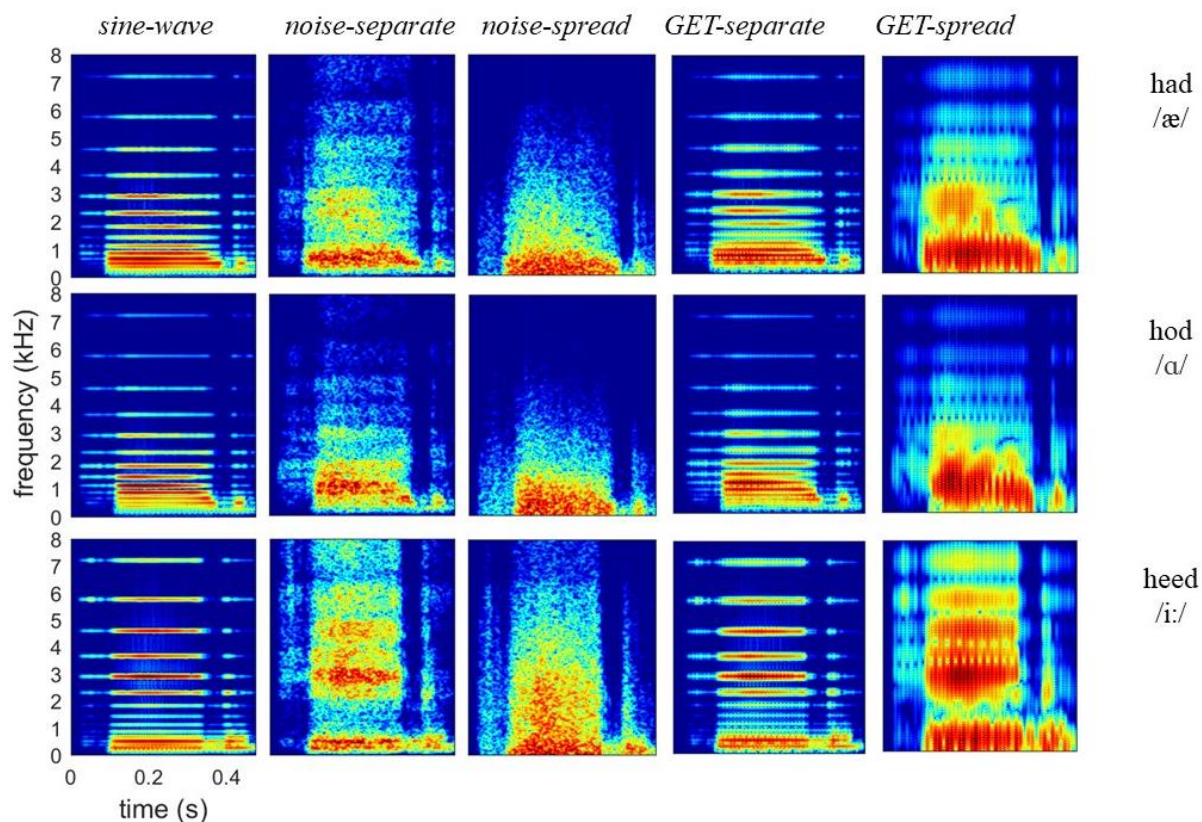


FIG. 5. (Color online) Spectrograms of three vowel stimuli encoded by the sine-wave, noise-separate, noise-spread, GET-separate, and GET-spread vocoders with 16 channels.

294 vowel tests. Each stimulus was presented 10 times. Stimuli were presented through headphones
295 (HDA 200, Sennheiser), and the sound level was calibrated to 70 dB SPL. This procedure was
296 conducted following procedures approved by the University of California Irvine Institutional
297 Review Board.

298 Seven normal hearing (NH) participants, ages 18-21, were tested in an anechoic chamber
299 (IAC) using the English vowel and consonant recognition tests adopted from [Friesen et al.](#)
300 [\(2001\)](#).

301 C. Results

302 Results are shown in Fig. 6. For the vowel test, the seven NH participants scored
303 approximately 20% under all simulation conditions with two channels. Increasing the number of
304 channels also improved performance. With eight channels, performance under the different
305 conditions began to separate. The sine-wave vocoder outperformed actual CI data, adapted from
306 [Friesen et al. \(2001\)](#), which showed no improvement beyond 8 channels. The noise-separate
307 vocoder and GET-separate vocoder showed similar performance trends. When electrode
308 interaction was simulated with overlapping filters, the subject performance showed a plateau near

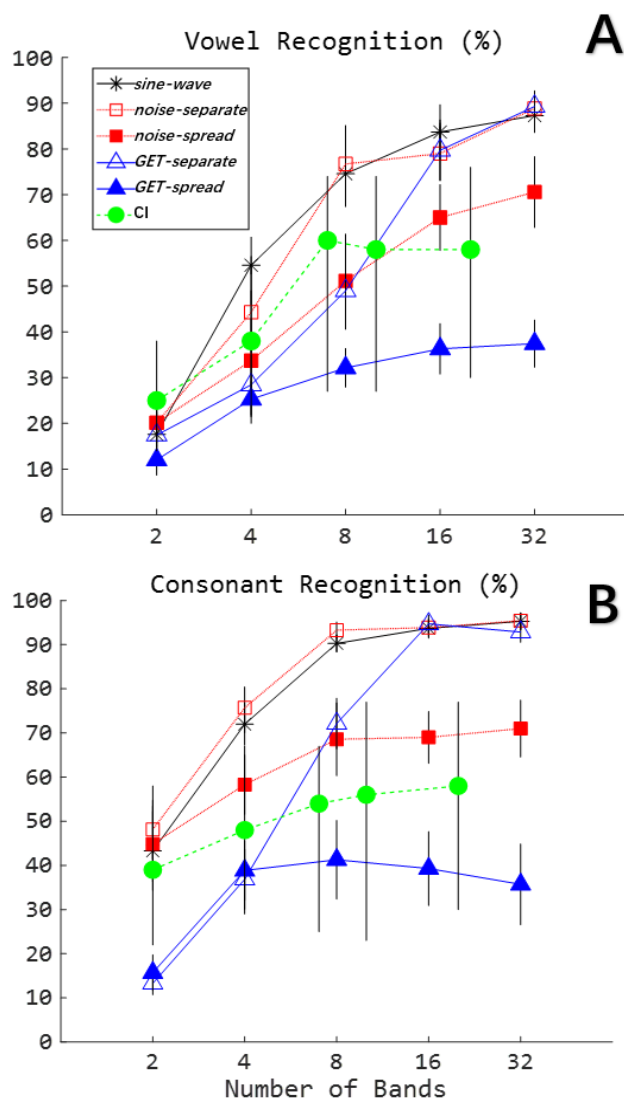


FIG. 6. (Color online) Vowel (A) and consonant (B) recognition as a function of number of bands (channels). Cochlear implant data is adapted from [Friesen et al. \(2001\)](#). Simulation data are averaged from seven normal hearing subjects listening to vocoded speech. For the simulation data, standard errors are indicated by the vertical bars. For the CI data, the bars show the entire ranges of performance across all their 19 participants.

309 60% with noise-spread, similar to actual CIs. The GET-spread condition underperformed CI data
310 in this case, saturating near 35% with eight channels.

311 Further, a two-way repeated-measures ANOVA with Geisser-Greenhouse correction was
312 used to analyze the vowel simulation results with vocoder and number of bands as the main
313 factors. The effect of vocoder ($F_{1.987, 11.92} = 49.87, p < 0.0001$), number of bands ($F_{2.018, 12.11} =$
314 $90.66, p < 0.0001$), and their interaction ($F_{3.890, 23.34} = 9.842, p < 0.0001$) were all significant. To
315 further analyze these effects, multiple comparisons with Bonferroni corrections were
316 implemented for each vocoder (to compare the five band numbers) and for each band number (to
317 compare the five vocoders). Table II shows the results of multiple comparisons between different
318 numbers of bands for each vocoder. Generally, there was a trend of better performance with more
319 bands. Still, the mean scores were not significantly different for 8, 16, and 32 bands (the only
320 exception was 8 vs. 32 with GET-separate). Table III shows the results of multiple comparisons
321 between vocoders for each number of bands. Because at 2 and 4 bands most vocoder pairs showed
322 no significant mean difference (the only exception was sine-wave vs. noise-spread at 4 number
323 of bands with $p = 0.009$), the comparison results bands were not listed. GET-spread derived the
324 lowest scores among the five vocoders at 16 and 32 bands, while GET-separate did not show
325 significantly different mean scores from the other three vocoders. The sine-wave, noise-separate,
326 and GET-separate vocoders did not show significantly different mean scores.

327 Consonant recognition showed similar performance trends across the simulation types, with
328 sine-wave, noise-separate, and GET-separate outperforming CIs (adapted from [Friesen et al.](#)
329 [\(2001\)](#)) when there were eight or more channels simulated. Noise-spread brought the performance
330 closer to actual CI data, while again GET-spread underperformed CIs. With only two channels,
331 both GET-separate and GET-spread showed much lower performance than actual CIs. For the

332 simulation results, consonant recognition scores were analyzed using the same statistical method
 333 as the above vowel data analysis. The effects of vocoder ($F_{1.404, 8.427} = 62.55, p < 0.0001$), number
 334 of bands ($F_{2.234, 13.40} = 379.0, p < 0.0001$), and their interaction ($F_{3.080, 18.48} = 10.88, p = 0.0002$)
 335 were all significant on consonant recognition. Results of multiple comparisons are shown in Table
 336 IV and V. The relative scores show similar trends as the results of multiple comparisons for
 337 vowel recognition (see Table II and III).

338 **TABLE II.** Results of multiple comparisons between vowel recognition scores with five band
 339 numbers for each of the five vocoders.

vocoder→ number of bands pair↓	sine-wave		noise-separate		noise-spread		GET-separate		GET-spread	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
2 vs. 4	*37.0	0.014	24.1	0.140	13.6	1.000	11.0	0.827	*13.3	0.019
2 vs. 8	**57.0	0.001	**56.6	0.002	31.0	0.085	31.6	0.082	*20.1	0.048
2 vs. 16	***66.1	<0.001	***58.9	<0.001	**44.9	0.001	***62.3	<0.001	*24.3	0.005
2 vs. 32	***69.7	<0.001	***68.7	<0.001	**50.4	0.001	***71.9	<0.001	*25.4	0.012
4 vs. 8	20	0.076	*32.4	0.019	17.4	0.057	20.6	0.209	6.9	0.564
4 vs. 16	*29.1	0.049	*34.7	0.023	***31.3	<0.001	**51.3	0.002	*11.0	0.034
4 vs. 32	*32.7	0.037	*44.6	0.011	***36.9	<0.001	***60.9	<0.001	12.1	0.296
8 vs. 16	9.1	0.199	2.29	1.000	13.9	0.147	30.7	0.058	4.1	1.000
8 vs. 32	12.7	0.288	12.1	0.696	19.4	0.059	*40.3	0.011	5.3	1.000
16 vs. 32	3.6	1.000	9.86	0.131	5.6	0.107	9.6	0.455	1.1	1.000

340
 341 **TABLE III.** Results of multiple comparisons between vowel recognition scores with five
 342 vocoders for each of the three band numbers (8, 16, and 32).

number of bands→ vocoder pair↓	8		16		32	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
sine-wave vs. noise-separate	-2.1	1.000	4.7	1.000	-1.6	1.000
sine-wave vs. noise-spread	23.4	0.054	**18.7	0.009	16.7	0.120
sine-wave vs. GET-separate	25.6	0.065	4.0	1.000	-2.0	1.000
sine-wave vs. GET-spread	**42.4	0.002	***47.4	<0.001	***49.9	<0.001
noise-separate vs. noise-spread	*25.6	0.030	*14.0	0.025	*18.3	0.048
noise-separate vs. GET-separate	27.7	0.080	-0.7	1.000	-0.4	1.000
noise-separate vs. GET-spread	**44.6	0.004	***42.7	<0.001	***51.4	<0.001
noise-spread vs. GET-separate	2.1	1.000	-14.7	0.130	-18.7	0.102
noise-spread vs. GET-spread	19.0	0.459	**28.7	0.002	**33.1	0.004
GET-separate vs. GET-spread	16.9	0.156	**43.4	0.001	***51.9	<0.001

343 **TABLE IV.** Results of multiple comparisons between consonant recognition scores with five
 344 number of bands for each of the five vocoders

vocoder→ number of bands pair↓	sine-wave		noise-separate		noise-spread		GET-separate		GET-spread	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
2 vs. 4	*37.0	0.014	24.1	0.140	13.6	1.000	11.0	0.827	*13.3	0.019
2 vs. 8	**57.0	0.001	**56.6	0.002	31.0	0.085	31.6	0.082	*20.1	0.048
2 vs. 16	***66.1	<0.001	***58.9	<0.001	**44.9	0.001	***62.3	<0.001	**24.3	0.005
2 vs. 32	***69.7	<0.001	***68.7	<0.001	**50.4	0.001	***71.9	<0.001	*25.4	0.012
4 vs. 8	20.0	0.076	*32.4	0.019	17.4	0.057	20.6	0.209	6.9	0.564
4 vs. 16	*29.1	0.049	*34.7	0.023	***31.3	<0.001	**51.3	0.002	*11.0	0.034
4 vs. 32	*32.7	0.037	*44.6	0.011	***36.9	<0.001	***60.9	<0.001	12.1	0.296
8 vs. 16	9.1	0.199	2.3	1.000	13.9	0.147	30.7	0.058	4.1	1.000
8 vs. 32	12.7	0.288	12.1	0.696	19.4	0.059	40.3	0.011	5.3	1.000
16 vs. 32	3.6	1.000	9.9	0.131	5.6	0.107	9.6	0.455	1.1	1.000

345 **TABLE V.** Results of multiple comparisons between consonant recognition scores with five
 346 vocoders for each of the five band numbers

number of bands→ vocoder pair↓	2		4		8		16		32	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
sine vs. noi-sep	-4.9	1.000	-3.7	1.000	-3.0	1.000	-0.1	1.000	-0.1	1.000
sine vs. noi-spr	-1.6	1.000	13.7	0.267	21.7	0.292	*24.7	0.016	*24.3	0.022
sine vs. GETsep	30.0	0.065	**35.1	0.002	18.1	0.096	-1.0	1.000	2.4	1.000
sine vs. GETspr	27.6	0.108	**33.1	0.010	**49.0	0.006	**54.4	0.002	**59.6	0.002
noi-sep vs. noi-spr	3.3	1.000	17.4	0.157	24.7	0.111	*24.9	0.021	*24.4	0.038
noi-sep vs. GETsep	**34.9	0.007	**38.9	0.007	21.1	0.115	-0.9	1.000	2.6	1.000
noi-sep vs. GETspr	*32.4	0.018	**36.9	0.002	**52.0	0.007	**54.6	0.003	**59.7	0.004
noi-spr vs. GETsep	*31.6	0.043	**21.4	0.003	-3.6	1.000	**25.7	0.010	*21.9	0.040
noi-spr vs. GETspr	29.1	0.078	**19.4	0.003	27.3	0.065	**29.7	0.006	***35.3	<0.001
GETsep vs. GETspr	-2.4	1.000	-2.0	1.000	*30.9	0.034	**55.4	0.002	**57.1	0.003

347 The current results suggest that the first type GET vocoder is feasible to simulate speech
 348 perception with CIs, and the CI current spread also could be simulated by manipulating durations
 349 of GETs. In both noise vocoder and GET vocoder, performance was substantially degraded by
 350 the increased current spread in both tasks. With eight or more bands, GET vocoders showed
 351 good simulation performance in that the actual CI data fell in the range between the separate
 352 and spread versions of the GETs.

353 **IV. EXPERIMENT 2. SIMULATION OF THE N-OF-M STRATEGY ACE**

354 **A. Rationale**

355 Some essential features of modern CI processing, including interleaved sampling, maxima
356 selection, amplitude compression and quantization, are omitted in not only conventional
357 continuous-carrier vocoders but also in the first type GET vocoder as used in Experiment 1. All
358 of these features may influence speech perception. According to the analysis in Section II, GETs
359 could be used to simulate them. The second type of GET vocoder ([Meng et al., 2018](#); [Kong et al.,](#)
360 [2019](#)) is introduced here in detail, and a battery of speech recognition tasks was carried out to
361 demonstrate its performance in Experiment 2. The experiment objective was to demonstrate the
362 potential of CI speech perception simulation with a GET vocoder involving all of the above-
363 mentioned essential features. The ACE strategy with 900-pps pulse rate was simulated by this
364 advanced GET vocoder.

365 **B. Vocoder Theory: Direct mapping from electric pulses to GETs**

366 In theory, the GETs are applicable for directly transferring any pulsatile CI electrodiagram
367 to a pulsatile vocoded sound. To be more illustrative, Fig. 7A demonstrates a 10-channel
368 electrodiagram (note: single vertical lines were used to represent electric pulses so that the
369 amplitude and timing of the electric pulse can be represented, while the phase and gap durations
370 in the common bi-phasic electric pulses were not considered in this study). To generate a GET
371 vocoder, the 10 channels were converted into frequency bands spanning over 10 equally divided
372 parts of the basilar membrane between characteristic frequencies of 150 and 8000 Hz
373 ([Greenwood, 1990](#)). The cutoff frequencies are 150, 271, 439, 672, 994, 1439, 2057, 2911, 4094,

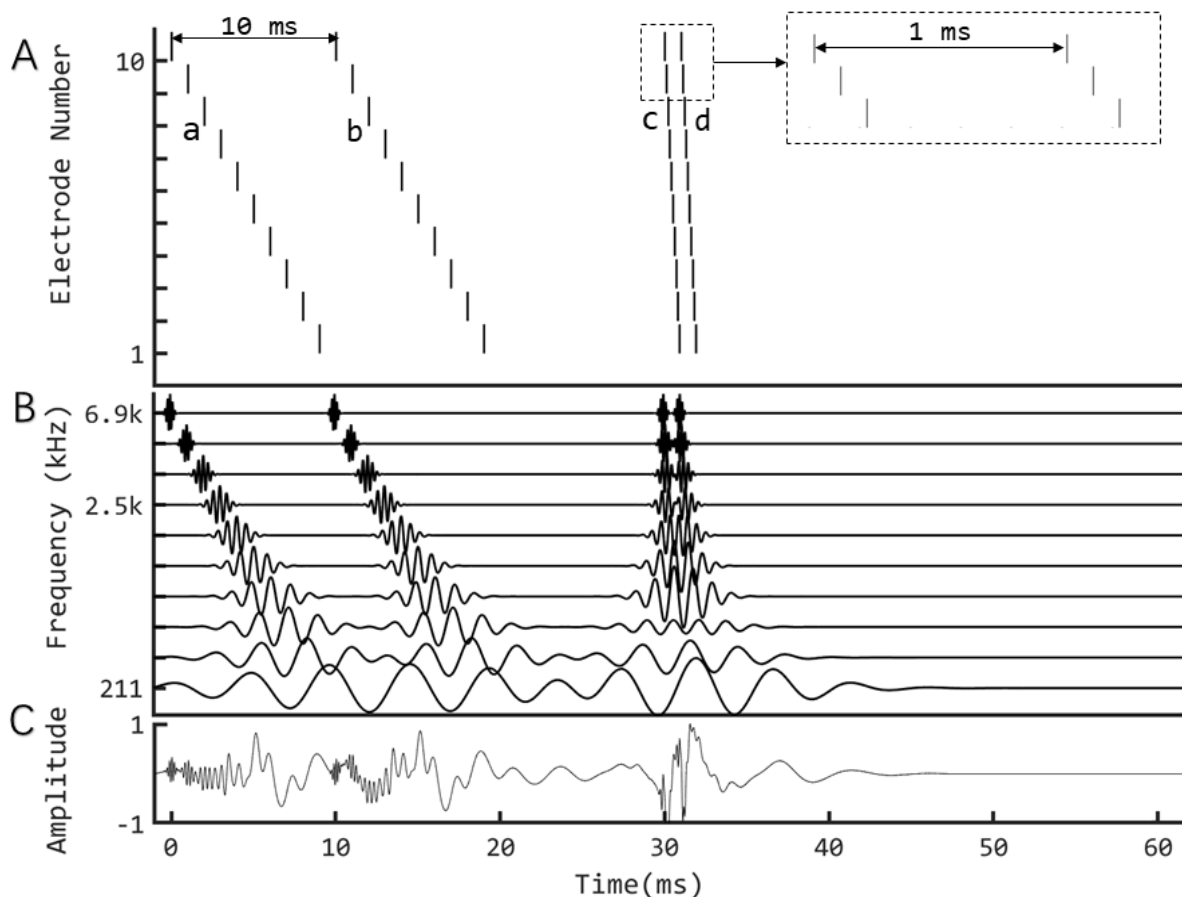


FIG. 7. Mapping a CI electrodiagram to a sound using the second type GET vocoder. **A.** An artificial 10-channel CI electrodiagram, including two pulse sweeps with a 10-ms difference between **a** and **b**, as well as two additional sweeps with a 1-ms difference between **c** and **d**, corresponding to stimulation rates of 100 pps and 1000 pps, respectively. **B.** GETs mimicking the electric pulse trains. **C.** The final GET waveform resulting from the sum of ten band-specific GET trains in **B**.

374 5732, and 8000 Hz. Then, a band-specific GET was generated in this demonstration by setting
 375 the parameters in Eq. 1 as $a = 1$, $t_0 = 0$, and

376
$$\sigma = \frac{2}{f_c} \quad (8)$$

377 where f_c denotes the center frequency of the specific band. As a result, the band-specific GET
 378 had a 6.82-dB duration of

379
$$D = \sqrt{2}\sigma = \frac{2\sqrt{2}}{f_c} \quad (9)$$

380 and a 6.82-dB bandwidth of

381
$$B = \frac{1}{D} = \frac{\sqrt{2}}{4} f_c \quad (10)$$

382 Then the acoustic GET train at the k^{th} channel in Fig. 7B is derived by

383
$$p_{a,k}(t) = (p_{e,k}(t) * e^{-\frac{\pi t^2}{2\sigma^2}}) \cdot \sin(2\pi f_c t + \varphi_0) \quad (11)$$

384 where $p_{e,k}(t)$ and $p_{a,k}(t)$ denotes the electric and acoustic pulse trains in Fig. 7A and 7B,
385 respectively, “*” denotes a convolution calculation, σ and f_c are band-dependent parameters as
386 defined above, and φ_0 is an initial phase that could be arbitrarily defined and was uniformly
387 randomized between 0 and 2π here.

388 Fig. 7B shows the 10-channel GET trains, which have temporally separated waveforms
389 for high-frequency channels, but overlapping waveforms for low-frequency channels. Fig. 7C
390 shows the overall waveform summed from the 10 bands.

391 According to the theoretical analysis of GET simulation, pulsatile features for individual
392 electric pulses cannot be guaranteed in the low-frequency channels, but the temporal-
393 separation feature between groups of pulses may be simulated to some extent. For example,
394 in Fig. 7B, at the lowest frequency channel, the 12-ms gap between b and c sweeps could have
395 a counterpart, i.e., a shallow amplitude-modulation dip, in the waveform.

396 C. Experiment method: Simulation of the n-of-m strategy ACE

397 Using the above method, any electrodiagrams, including the widely used *n-of-m* strategy
398 like ACE strategy which is the current default strategy in Nucleus cochlear implants (Vandali et

399 [al., 2000](#)), can be converted to vocoded sounds. The specific vocoder is named ACE-GET.
400 Following the preliminary results which showed comparable acute data between the ACE-GET
401 vocoder and actual CI users ([Kong et al., 2019](#)), in this paper a battery of speech recognition
402 tasks was carried out to further explore the potential of ACE-GET vocoder on simulation of
403 speech perception with CIs.

404 In the clinical fitting of ACE strategy, the intensity dynamic range should be measured
405 behaviorally electrode-by-electrode and is also limited and variable among users. In the ACE-
406 GET vocoders, the dynamic range could be easily manipulated either in the compression stage
407 of the ACE encoding or in the inverse compression stage of the GET synthesizing. The latter
408 method was used in this study, and two dynamic ranges corresponding to two ACE-GET
409 vocoders were tested. It was hypothesized that the vocoder with a higher dynamic range would
410 simulate the top CI participants while the vocoder with a lower dynamic range would simulate
411 the average performance of CI participants. The combination of $n = 8$ and $m = 22$ is one default
412 option in the clinical fitting of ACE and was simulated in this experiment.

413 In detail, two 22-channel ACE-GET vocoders (denoted by GETlargeDR and
414 GETsmallDR) were compared with two 22-channel sine-carrier conventional vocoders (125 Hz
415 and 250 Hz envelope cutoffs, denoted by Sin250 and Sin125, respectively) with minimum
416 channel overlapping as shown in Fig. 3A. The hypotheses for the parameter selection of the four
417 vocoders are discussed later.

418 *Detailed implementation methods of the vocoders:* First, the default setting of the ACE
419 software integrated in the CCI-Mobile software ([Ghosh et al., 2022](#)) was used to convert input
420 sounds into electrograms. An inverse-mapping function was used to transfer the electric
421 current value of each electric pulse in the electrogram to an envelope power value. Single-

422 sample pulse trains from each band were “convolved” with a Gaussian function with $\sigma = 3/f_c$.
423 In the specific implementation of the experiment, the convolution step was replaced by simply
424 comparing any overlapping sampling points from two GETs and preserving the larger point as
425 the final sample value. In the theory and framework analysis in Section II, a convolution
426 calculation was recommended, but in our experiment, we only preserved the largest point to
427 show better pulsatile waveform than the cumulative effect of a convolution. The output was
428 used to multiply a sinusoidal carrier with a frequency of f_c at the center of the corresponding
429 band and an arbitrary initial phase (a random initial phase in this study). The average power of
430 each band was kept unchanged. Finally, the modulated signals were summed to produce the
431 vocoded stimulus.

432 The difference between GETlargeDR and GETsmallDR was only between their inverse
433 (i.e., electric-to-acoustic) mapping functions, which are Eqs. 12 and 13, respectively:

$$434 \quad L_a = \frac{1}{\alpha} ((1 + \alpha)^{L_e} - 1) \quad (12)$$

435 and

$$436 \quad L_a = \frac{1}{2.72\alpha} (e^{L_e}(1 + \alpha) - 1) \quad (13)$$

437 in which, the L_a denotes the recovered acoustic level, L_e denotes the electric current level
438 defined by the electrodiagram from the ACE strategy based on a specific patient’s fitting map,
439 and α is a constant 416.0. In the present study, the threshold levels and most comfortable levels
440 are constantly defined as 100 and 255 CU (current unit), i.e., $100 \text{ CU} < L_e < 255 \text{ CU}$. In this
441 case, based on Eqs. 12 and 13, the recovered acoustic level ranges were 32.7 dB and 5.3 dB for
442 GETlargeDR and GETsmallDR, respectively. The output stimuli level was controlled at a
443 comfortable level around 65 dBA. Equation 12 is directly based on the default setting of the
444 acoustic-to-electric compression function in ACE. It was hypothesized that GETlargeDR could

445 simulate the best performance of CI listeners with the corresponding ACE strategy and
446 GETsmallDR would significantly degrade the performance because of the much narrower range.
447 Otherwise, the implementation details of the vocoder were the same as in [Meng et al. \(2018\)](#).

448 In the two sine vocoders, the frequency spacing for cutoffs for the analysis filters was
449 defined in the range of [80, 7999] Hz according to a Greenwood map ([Greenwood, 1990](#)).
450 Specifically, the cutoff frequencies were 80, 122, 172, 230, 298, 379, 473, 583, 712, 864, 1042,
451 1250, 1494, 1781, 2117, 2512, 2974, 3516, 4152, 4898, 5772, 6797, and 7999 Hz. The filtered
452 signals were full-wave rectified and low-pass filtered (6th order Butterworth; 125 Hz for Sin125
453 and 250 Hz for Sin250) to extract the envelope for each channel. A sine wave with a frequency
454 centered at the corresponding analysis band was used as the carrier, which was then multiplied
455 by the corresponding envelope. The final vocoded stimuli were generated by a summation of
456 the modulated carriers. In previous studies, it was found that speech intelligibility was better
457 with a higher cutoff frequency in the envelope extraction ([Souza and Rosen, 2009](#)). Therefore,
458 Sin250 was expected to be better than Sin125.

459 In Fig. 8, a Mandarin sentence was used to demonstrate the vocoded speech using the four
460 vocoders, i.e., GETlargeDR, GETsmallDR, Sin250, and Sin125. It shows that the GET vocoders
461 resemble the ACE-electrodiagram more than the sine vocoders. The temporal separation between
462 groups of pulses can also be found in the band signals of GET vocoded speech. Because the
463 GET vocoders directly use the information of the ACE electrodiagram, it was hypothesized that
464 speech intelligibility would be worse, but closer to actual CI results, with the GET vocoders than
465 with the sine vocoders.

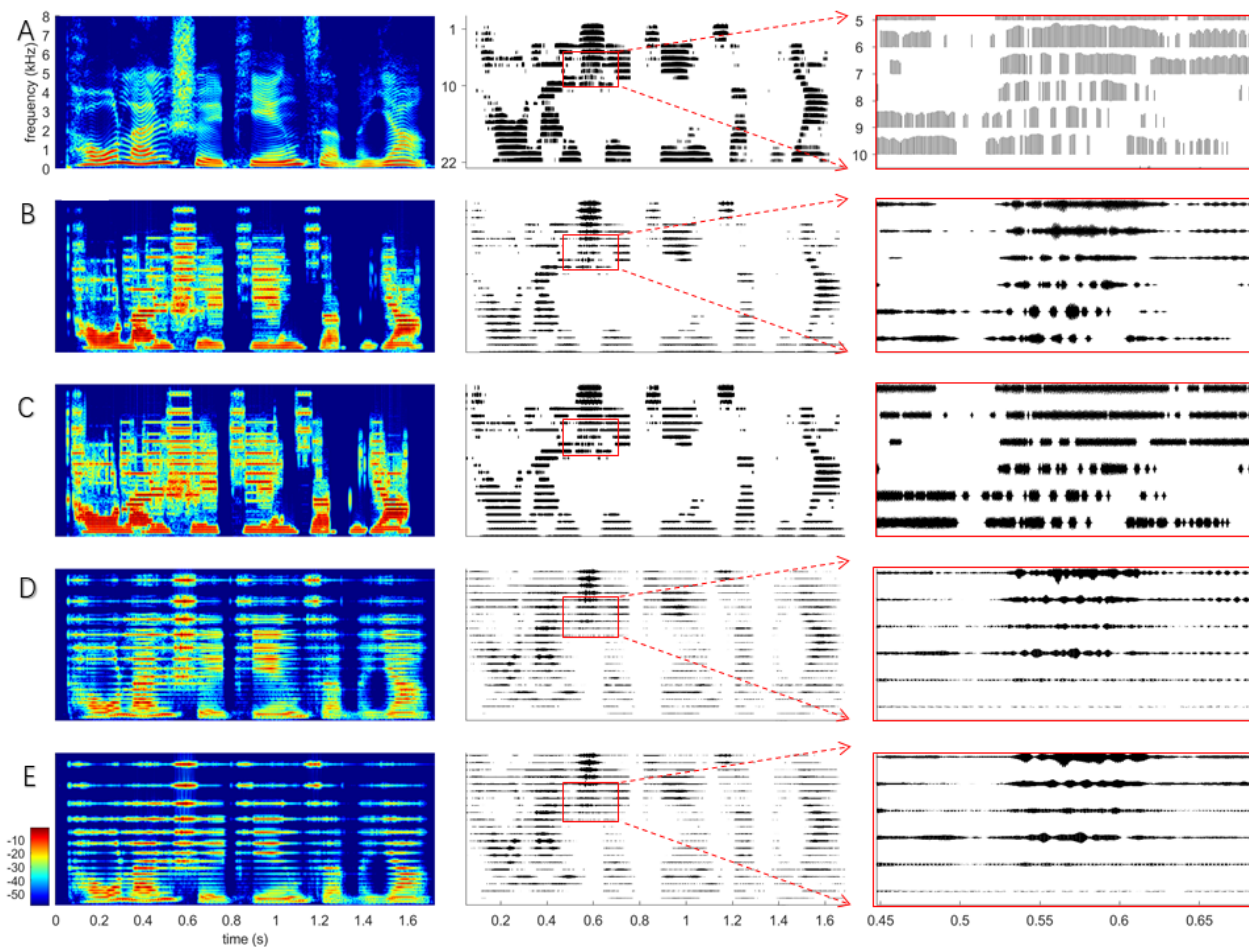


FIG. 8. (Color online) Speech stimulus demonstrations for the ACE-GET simulation experiment. Left: Spectrogram; middle: band-specific signal; right: zoom in of the boxed signals. **A.** Spectrogram and ACE electrodiagram of a clear sentence of speech. **B-E.** Spectrogram and band-specific waveforms of vocoded speech using two GET vocoders (GETlargeDR, and GETsmallDR) and two conventional sine-wave vocoders (Sin250 and Sin125), respectively.

466

467 **D. Experiment method: Participants and Tasks**

468 Two groups of NH participants (ten in each group, ages 18-29, and native Mandarin
469 speakers) were tested in a soundproof room. Group 1 used Sin250 and GETlargeDR, and Group
470 2 used Sin125 and GETsmallDR. Three open-set Mandarin Chinese recognition tasks were

471 tested, i.e., time-compression threshold, sentence-in-noise recognition, sentence-in-
472 reverberation recognition. The results for the four tasks with the two vocoders in these NH
473 participants were compared with actual CI results from our previous experiments (Meng et al.,
474 2019) as well as newly collected data in this work. These experiments were conducted
475 following procedures approved by the Medical Ethics Committee of Shenzhen University,
476 China. Detailed information about the three experiments is as follows:

477 1) Time-compression thresholds (TCTs), i.e., accelerated sentence speeds at which 50%
478 of words could be recognized correctly, were measured using the Mandarin speech perception
479 corpus (Fu et al., 2011).

480 2) Speech reception thresholds (SRTs) in speech-shaped noise (SSN) and babble noise,
481 i.e., signal-to-noise ratio (SNR) at which 50% of words could be recognized correctly, were
482 measured using the Mandarin hearing in noise test (MHINT) corpus (Wong et al., 2007). The
483 TCT and SRT test procedures followed Experiment 2 of Meng et al. (2019) strictly, in which
484 ten CI subjects (9/10 adults) with various hearing histories were tested.

485 3) Recognition of speech in reverberation was measured using a Mandarin BKB-like
486 sentence corpus (Xi et al., 2012), whose quiet sentences were convolved with simulated room
487 impulse responses (RIRs). The RIRs were generated using a MATLAB function
488 (<https://www.audiolabs-erlangen.de/fau/professor/habets/software/rir-generator>) with its default
489 setting, except the reverberation times (T60) were set as 0, 0.3, 0.6, and 0.9 s. For each T60, one
490 sentence list was used. Seven CI participants with various hearing histories were also tested for
491 comparison (See Table VI).

492 We had three subject groups, two of which were NH listeners each using two different
493 vocoders. A mixed model was used to assess the repeated measures within subjects as well as

494 independent measures between subjects. The paired-sample *t*-test and two-sample *t*-test were used
495 to examine the statistical significance of the means' difference for within-subject comparisons
496 and between-subject comparisons, respectively. For each task, the five CI processing conditions,
497 i.e., Sin250, Sin125, GETlargeDR, GETsmallDR, and CI, were pair-wisely examined to yield 10
498 pairs of comparison. Bonferroni corrections were used to adjust the *p* values, and the final
499 significance was examined using the criterion of 0.05.

500 **TABLE VI.** Detailed information of the 7 CI participants in the speech in reverberation test

Subject	Gender	Age (yr)	CI Experience (yr)	CI Processor	Etiology
C14	F	41	12	CP810	Drug induced
C23	M	31	11	CP810	Sudden deafness
C30	M	13	10	Freedom	LVAS
M5	M	18	15	OPUS-2	Virus infection
C16	F	25	2	Freedom	Unknown
M17	F	18	5	OPUS-2	Genetic
M16	F	17	7	OPUS-2	Unknown

501 E. Results

502 The results with the four 22-channel vocoders, i.e., GETlargeDR, GETsmallDR, Sin250,
503 and Sin125 are shown in Fig. 9.

504 For the TCT test (Fig. 9A), a significant decreasing trend was found from Sin250 (mean =
505 16.1 syllables/sec), Sin125 (13.9), GETlargeDR (12.3), GETsmallDR (9.4), to actual CI (6.8)
506 results (Bonferroni adjusted $p < 0.05$), while their standard deviations are comparable within the
507 range from 1.0 to 1.2 syllables/s.

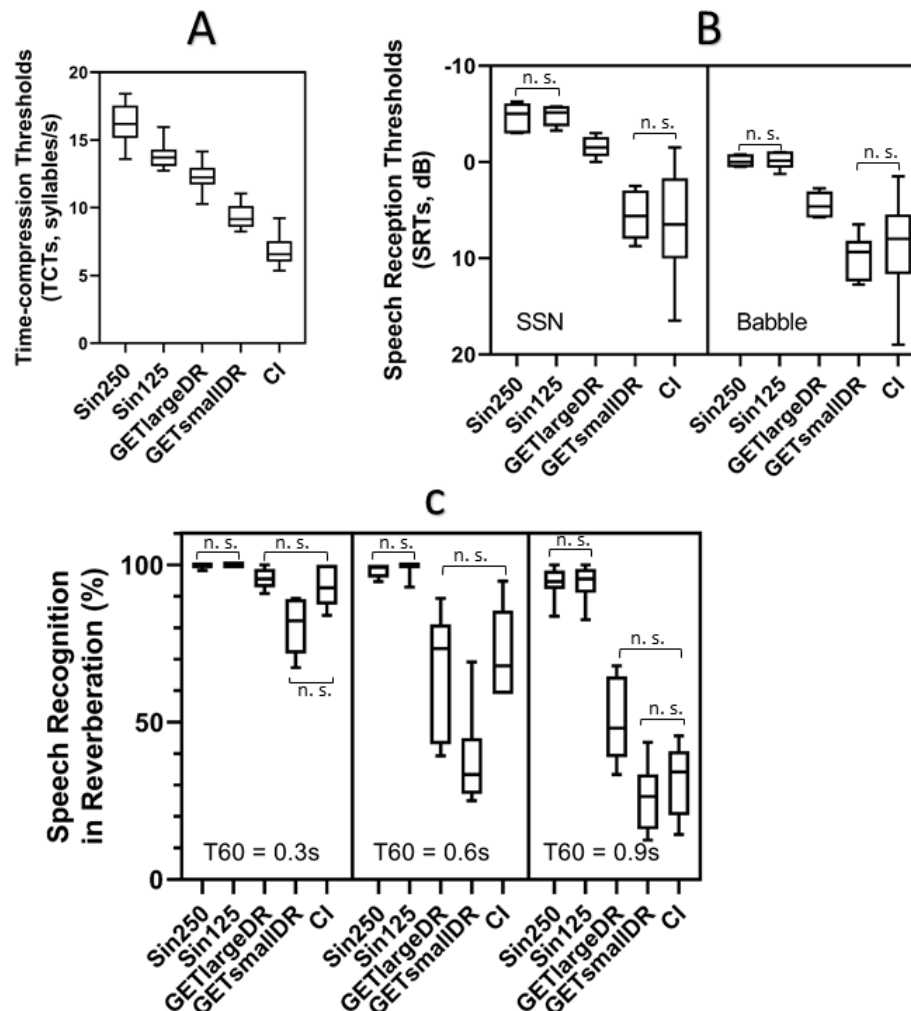


FIG. 9. Results from three speech recognition tasks with two 22-channel sine-wave vocoders (Sin250: 250 Hz cut-off envelope; Sin125: 125 Hz cut-off envelope) and two GET vocoders (GETlargeDR and GETsmallDR; their difference is only in the intensity dynamic range, i.e., 32.7 dB and 5.3 dB for GETlargeDR and GETsmallDR respectively) compared with the results of some CI subjects. There were two groups of normal-hearing participants, each with ten participants. One group used Sin250 and GETlargeDR, and the other group used Sin125 and GETsmallDR. **A.** Time-compression threshold results. **B.** Speech reception threshold results of a speech in noise recognition experiment (SSN and babble noise). **C.** Speech recognition scores in reverberation with T60 = 0.3, 0.6, and 0.9s. Pairwise comparisons with Bonferroni corrections were examined. In each box, “n. s.” denotes the non-significant difference ($p > 0.05$), otherwise, there was a significant difference.

509 For the SRT test (Fig. 9B), there was no significant difference (adjusted $p > 0.05$) between
510 Sin250 (means: -4.7 dB in SSN and -0.1 dB in babble noise) and Sin125 (means: -4.8 dB in
511 SSN and -0.1 dB in babble noise) and between GETsmallDR (means: 5.6 dB in SSN and 10 dB
512 in babble noise) and actual CIs (means: 6.5 dB in SSN and 8.8 dB in Babble noise). The mean
513 results with GETlargeDR (means: -1.5 dB in SSN and 4.5 dB in babble noise) were significantly
514 lower (adjusted $p < 0.05$) than those with Sin250 and Sin125, and significantly higher (adjusted
515 $p < 0.05$) than those with GETsmallDR and CIs. The mean SRTs in babble noise were always
516 significantly lower than those in SSN for all four vocoder conditions (adjusted $p < 0.05$). For CI
517 users, mean SRTs in the two noise types did not show a significant difference (adjusted $p >$
518 0.05).

519 For the reverberant speech recognition test (Fig. 9C), all vocoders and the actual CI
520 condition showed a significant trend of decreased recognition scores when the reverberation
521 time increased. However, the sine vocoder simulations were much less sensitive to reverberation
522 than the CI users. It is shown that even with $T60 = 0.9$ s, the sine vocoders still derived $>94\%$
523 means, which were much higher than CI participants' 32% . The GETlargeDR and GETsmallDR
524 derived significantly lower scores than the sine vocoders did (adjusted $p < 0.05$). Under the $T60$
525 $= 0.3$ s and 0.9 s conditions, there was no significant mean score difference between either GET
526 vocoder and CI (adjusted $p > 0.05$), while GETsmallDR derived significantly lower mean scores
527 than GETlargeDR did (adjusted $p > 0.05$). However, the mean results with CI were closer to
528 GETlargeDR at $T60 = 0.3$ s and to GETsmallDR at $T60 = 0.9$ s. Under the $T60 = 0.6$ s condition,
529 there was no significant mean score difference between GETlargeDR and CI, while
530 GETsmallDR derived significantly higher mean scores than GETlargeDR and CI.

531 In all three tasks, GET vocoders were able to simulate actual CI performance more closely
532 than sine vocoders. In fact, the sine vocoders overestimated CI performance in all tasks. Sin250
533 performed slightly better than Sin125 in mean results but did not show a significant difference.
534 In the time-compression task, all vocoders produced better than CI performance, with
535 GETsmallDR being the closest (Fig. 9A). In the SRT-in-noise test, GETsmallDR and CI
536 produced comparable performance (Fig. 9B). In the reverberation task, GETlargeDR had
537 similar-to-CI performance in all T60 conditions and GETsmallDR in the T60 = 0.3 and 0.9s
538 conditions (Fig. 9C).

539 V. DISCUSSION

540 Sounds are transmitted through air as continuous compression waves, but they are encoded
541 by discrete spikes in the neural system and by pulsatile electric stimuli in CIs. Vocoders have
542 been developed to simulate the signal processing and sound perception of CIs. However, the
543 pulsatile feature, which is acknowledged as critical to the success of modern CIs, has not been
544 simulated until now by the most widely used noise and sine-wave excited vocoder ([Shannon et al., 1995](#);
545 [Dorman et al., 1997](#)). Some studies have proposed pulsatile vocoders using filtered
546 carriers with strong periodicities including noise burst ([Blamey et al., 1984a](#); [Blamey et al.,](#)
547 [1984b](#)) and complex tones ([Deeks and Carlyon, 2004](#); [Hilkuysen and Macherey, 2014](#);
548 [Mesnildrey et al., 2016](#)). Instead of using filtered carriers, some CI manufacturers have provided
549 software to directly map electrograms to vocoded sounds ([Ausili et al., 2019](#); [Stam et al.,](#)
550 [2019](#)). In this study, a GET-based vocoder was proposed, theoretically analyzed, and evaluated
551 for its performance on CI speech perception simulation.

552 **A. GETs and electric pulses**

553 The GET can be used to simulate a “perceivable” atom of sound, which can be traced back
554 to [Gabor \(1947\)](#). More recently, it has been used in many psychoacoustic studies. The GET
555 vocoder model can be a phenomenological one, in which each GET corresponds to an electrical
556 pulse. The amplitude of the GET is scaled proportionally to the pulse current level. Moreover,
557 the GET vocoders can simulate main features in CIs, including the place of stimulation, pulse
558 time, temporal envelope, spectral envelope and spectral interaction, and intensity quantization
559 and maxima-selection, by corresponding features of the acoustic pulses.

560 An inherent limitation with the GETs is the tradeoff between temporal duration and spectral
561 bandwidth. Shortening the GET duration increases the spectral bandwidth, which introduces
562 temporal or spectral overlaps between different GETs, especially at low frequencies (see Fig. 7
563 and related text). Real CIs have no such limitation, in which both pulse duration and pulse rate
564 are the same whether it is a basal or apical electrode.

565 **B. Speech perception with GET vocoders**

566 In this study, two types of GET vocoders (Fig. 3B&C) were proposed to simulate different
567 aspects of CI processing ([Lu et al., 2007](#); [Meng et al., 2018](#)). The first GET vocoder simply
568 replaced the continuous noise or sine-wave carriers in conventional vocoders by a new type of
569 carrier, or GET train. In the first implementation (Fig. 3B), a non-interleaved sampling 100-pps
570 GET carrier was generated to study the effects of spread of excitation by controlling the GET
571 duration according to the time-frequency uncertainty principle. Spread of excitation is an
572 important factor underlying the poor- and large-variance performance for CI participants ([Fu
573 and Nogaki, 2005](#); [Bingabr et al., 2008](#); [Strydom and Hanekom, 2011](#); [Grange et al., 2017](#);
574 [O'Neill et al., 2019](#); [Mehta et al., 2020](#)). Different from the noise-or sine-vocoders that produced

575 performance better than actual CI performance even in the case of the severe channel interaction
576 (i.e., using the low-pass filtered noise carriers), the GET vocoder produced a wide range of
577 vowel and consonant recognition performance encompassing actual CI performance (Fig. 6).
578 One limitation in this experiment was that the spectral spread simulated by GET vocoders at
579 low frequency channels might be influenced by the sparsity of the electric pulses. For example
580 (see Fig.7), at the lowest frequency channel, temporal overlap happens between two GETs and
581 the bandwidth of the two overlapped GETs is narrower than an isolated GET. Fortunately, due
582 to the sparse nature of speech signal and narrower GET durations at higher channels, the effects
583 of this limitation should be limited. Another limitation of Experiment 1 was that all vocoders
584 used a 50-Hz envelope cutoff frequency, which was lower than real CIs.

585 The second vocoder directly mapped individual electric pulses in a CI electrodiagram to
586 individual GETs to simulate the ACE strategy (Fig. 3C). This direct mapping allows simulation
587 of all processing steps including the *n-of-m* maxima selection to amplitude compression and
588 quantization. Compared with the conventional sine-wave vocoder, not only did the GET vocoder
589 better resemble the ACE electrodiagram, but more importantly the GET vocoder produced a
590 mean and range of speech in noise recognition performance similar to that of actual CI users. In
591 particular, the wider dynamic range simulated better CI performance (Fig. 9). Future studies are
592 needed to establish and evaluate individualized CI simulation, in which both the mean and error
593 patterns of phonemic recognition are used to judge the validity and quality of the simulation
594 model (DiNino et al., 2016; Winn, 2020; Bance et al., 2022).

595 The GET vocoder is perhaps a more general vocoder model as it can closely approximate
596 conventional noise (using noise carriers instead of sine waves) and sine-wave vocoders by
597 summing many GETs occurring at high rates or long GET duration and using high-fidelity

598 intensity (or envelope) information. This means that the conventional vocoders can be treated
599 as special cases of GET vocoders.

600 The MATLAB source code of the GET vocoder for the ACE strategy is provided for
601 academic research purposes². Based on this code, more variants could be generated by
602 manipulating the vocoder parameters, e.g., spectral spread, stimulation place or frequency
603 shifting, and carrier types.

604 V. CONCLUSION

605 This study indicates that pulsatile simulation of speech, which is a key to the success of
606 modern CI and has been omitted in previous vocoders, could be realized by using the proposed
607 GET vocoders. The main conclusions include:

- 608 (1) The time-frequency uncertainty principle empowers and imposes constraints on using
609 GETs for CI simulation;
- 610 (2) Many features of modern CIs including pulsatile timing, current spread, *n-of-m*
611 maxima selection, dynamic compression could be implemented in GET vocoders and
612 then used to derive similar sentence recognition performance to actual CI users;
- 613 (3) A GET vocoder framework for arbitrary CI strategy and a package of source code
614 (using ACE as an example) are provided to serve as a general-purpose research tool to
615 generate vocoded sounds (including speech) based on direct pulse-to-pulse mapping.

616 Further experiment studies (e.g., in phoneme confusion patterns) are warranted to systematically
617 examine the performance of GET simulation.

² Currently as an attachment of the submission and will be open at a permanent website before the final version if the manuscript could be accepted for publication in JASA.

618 **ACKNOWLEDGMENTS**

619 We thank all the participants in these experiments. J. Carroll and S. Tiaden helped collect the data
620 in Experiment 1. Fanhui Kong and Yulong Xiao helped collect the data in Experiment 2. This research
621 was supported by NIH R01 DC15587 (F.G.Z.), National Natural Science Foundation of China (11704129
622 and 61771320), Guangdong Basic and Applied Basic Research Foundation Grant (2020A1515010386),
623 and Science and Technology Program of Guangzhou (202102020944) (Q.M.). Thanks to Drew Cappotto
624 for proof-reading this article.

625 **REFERENCES**

- 626 Ausili, S. A., Backus, B., Agterberg, M. J. H., van Opstal, A. J., and van Wanrooij, M. M. (2019).
627 "Sound localization in real-time vocoded cochlear-implant simulations with Normal-
628 Hearing Listeners," *Trends. Hear.* **23**, 2331216519847332.
- 629 Baer, T., Moore, B. C. J., and Glasberg, B. R. (1999). "Detection and intensity discrimination of
630 Gaussian-shaped tone pulses as a function of duration," *J. Acoust. Soc. Am.* **106**, 1907-
631 1916.
- 632 Baer, T., Moore, B. C. J., and Marriage, J. (2001). "Detection and intensity discrimination of brief
633 tones as a function of duration by hearing-impaired listeners," *Hear. Res.* **159**, 74-84.
- 634 Bance, M., Brochier, T., Vickers, D., and Goehring, T. (2022). "From microphone to phoneme: an
635 end-to-end computational neural model for predicting speech perception with cochlear
636 implants," *IEEE Trans Biomed Eng.* Early Access.
- 637 Bingabr, M., Espinoza-Varas, B., and Loizou, P. C. (2008). "Simulating the effect of spread of
638 excitation in cochlear implants," *Hear. Res.* **241**, 73-79.
- 639 Blamey, P. J., Dowell, R. C., Tong, Y. C., Brown, A. M., Luscombe, S. M., and Clark, G. M.
640 (1984a). "Speech processing studies using an acoustic model of a multiple-channel cochlear
641 implant," *J. Acoust. Soc. Am.* **76**, 104-110.
- 642 Blamey, P. J., Dowell, R. C., Tong, Y. C., and Clark, G. M. (1984b). "An acoustic model of a
643 multiple-channel cochlear implant," *J. Acoust. Soc. Am.* **76**, 97-103.
- 644 Brown, A. D., and Stecker, G. C. (2010). "Temporal weighting of interaural time and level
645 differences in high-rate click trains," *J. Acoust. Soc. Am.* **128**, 332-341.
- 646 Buell, T. N., and Hafter, E. R. (1988). "Discrimination of interaural differences of time in the
647 envelopes of high-frequency signals: integration times," *J. Acoust. Soc. Am.* **84**, 2063-2066.
- 648 Croghan, N. B. H., and Smith, Z. M. (2018). "Speech understanding with various maskers in
649 cochlear-implant and simulated cochlear-implant hearing: effects of spectral resolution and
650 implications for masking release," *Trend. Hear.* **22**. 2331216518787276.
- 651 Deeks, J. M., and Carlyon, R. P. (2004). "Simulations of cochlear implant hearing using filtered
652 harmonic complexes: Implications for concurrent sound segregation," *J. Acoust. Soc. Am.*
653 **115**, 1736-1746.
- 654 Dieudonne, B., Van Wilderode, M., and Francart, T. (2020). "Temporal quantization deteriorates
655 the discrimination of interaural time differences," *J. Acoust. Soc. Am.* **148**, 815-828.

- 656 DiNino, M., Wright, R. A., Winn, M. B., and Bierer, J. A. (2016). "Vowel and consonant confusions
657 from spectrally manipulated stimuli designed to simulate poor cochlear implant electrode-
658 neuron interfaces," J. Acoust. Soc. Am. **140**, 4404-4418.
- 659 Dorman, M. F., Loizou, P. C., and Rainey, D. (1997). "Speech intelligibility as a function of the
660 number of channels of stimulation for signal processors using sine-wave and noise-band
661 outputs," J. Acoust. Soc. Am. **102**, 2403-2411.
- 662 Dudley, H. (1939). "Remaking Speech," J. Acoust. Soc. Am. **11**, 169-177.
- 663 Ehlers, E., Kan, A., Winn, M. B., Stoelb, C., and Litovsky, R. Y. (2016). "Binaural hearing in
664 children using Gaussian enveloped and transposed tones," J. Acoust. Soc. Am. **139**, 1724-
665 1733.
- 666 Feichtinger, H. G., and Strohmer, T. (1998). *Gabor analysis and algorithms: theory and*
667 *applications* (Birkhäuser, Boston).
- 668 Friesen, L. M., Shannon, R. V., Baskent, D., and Wang, X. (2001). "Speech recognition in noise as
669 a function of the number of spectral channels: Comparison of acoustic hearing and cochlear
670 implants," J. Acoust. Soc. Am. **110**, 1150-1163.
- 671 Fu, Q. J., and Nogaki, G. (2005). "Noise susceptibility of cochlear implant users: the role of spectral
672 resolution and smearing," J. Assoc. Res. Otolaryngol. **6**, 19-27.
- 673 Fu, Q. J., Zhu, M., and Wang, X. (2011). "Development and validation of the Mandarin speech
674 perception test," J. Acoust. Soc. Am. **129**, EL267-273.
- 675 Gabor, D. (1947). "Acoustical quanta and the theory of hearing," Nature **159**, 591-594.
- 676 Gardner, T. J., and Magnasco, M. O. (2006). "Sparse time-frequency representations," Proc. Natl.
677 Acad. Sci. U.S.A. **103**, 6094-6099.
- 678 Ghosh, R., Ali, H., and Hansen, J. H. L. (2022). "CCi-MOBILE: a portable real time speech
679 processing platform for cochlear implant and hearing research," IEEE Trans. Biomed. Eng.
680 **69**, 1251-1263.
- 681 Goupell, M. J., Majdak, P., and Laback, B. (2010). "Median-plane sound localization as a function
682 of the number of spectral channels using a channel vocoder," J. Acoust. Soc. Am. **127**, 990-
683 1001.
- 684 Goupell, M. J., Stoelb, C., Kan, A., and Litovsky, R. Y. (2013). "Effect of mismatched place-of-
685 stimulation on the salience of binaural cues in conditions that simulate bilateral cochlear-
686 implant listening," J. Acoust. Soc. Am. **133**, 2272-2287.

- 687 Grange, J. A., Culling, J. F., Harris, N. S. L., and Bergfeld, S. (2017). "Cochlear implant simulator
688 with independent representation of the full spiral ganglion," *J. Acoust. Soc. Am.* **142**,
689 EL484-489.
- 690 Greenwood, D. D. (1990). "A cochlear frequency-position function for several species—29 years
691 later," *J. Acoust. Soc. Am.* **87**, 2592-2605.
- 692 Hilkhuisen, G., and Macherey, O. (2014). "Optimizing pulse-spreading harmonic complexes to
693 minimize intrinsic modulations after auditory filtering," *J. Acoust. Soc. Am.* **136**, 1281.
- 694 Johnson, L. A., Della Santina, C. C., and Wang, X. (2017). "Representations of time-varying
695 cochlear implant stimulation in auditory cortex of awake marmosets (*Callithrix jacchus*),"
696 *J. Neurosci.* **37**, 7008-7022.
- 697 Jones, H., Kan, A., and Litovsky, R. Y. (2014). "Comparing sound localization deficits in bilateral
698 cochlear-implant users and vocoder simulations with normal-hearing listeners," *Trend.*
699 *Hear.* **18**. 2331216514554574.
- 700 Kan, A., Stoelb, C., Litovsky, R. Y., and Goupell, M. J. (2013). "Effect of mismatched place-of-
701 stimulation on binaural fusion and lateralization in bilateral cochlear-implant users," *J.*
702 *Acoust. Soc. Am.* **134**, 2923-2936.
- 703 Kong, F., Wang, X., Teng, X., Zheng, N., Yu, G., and Meng, Q. (2019). "Reverberant speech
704 recognition with actual cochlear implants: verifying a pulsatile vocoder simulation method,"
705 in *23rd International Congress on Acoustics* (Aachen, Germany), pp. 3109-3112.
- 706 Laback, B., Balazs, P., Necciari, T., Savel, S., Ystad, S., Meunier, S., and Kronland-Martinet, R.
707 (2011). "Additivity of nonsimultaneous masking for short Gaussian-shaped sinusoids," *J.*
708 *Acoust. Soc. Am.* **129**, 888-897.
- 709 Laback, B., Necciari, T., Balazs, P., Savel, S., and Ystad, S. (2013). "Simultaneous masking
710 additivity for short Gaussian-shaped tones: Spectral effects," *J. Acoust. Soc. Am.* **134**, 1160-
711 1171.
- 712 Lu, T., Carroll, J., and Zeng, F. G. (2007). "On acoustic simulations of cochlear implants," in
713 *Conference on Implantable Auditory Prostheses* (Lake Tahoe, CA).
- 714 Lu, T., Liang, L., and Wang, X. (2001). "Temporal and rate representations of time-varying signals
715 in the auditory cortex of awake primates," *Nat. Neurosci.* **4**, 1131-1138.
- 716 Lu, T., Litovsky, R., and Zeng, F. G. (2010). "Binaural masking level differences in actual and
717 simulated bilateral cochlear implant listeners," *J. Acoust. Soc. Am.* **127**, 1479-1490.

- 718 Lu, T., and Wang, X. (2000). "Temporal discharge patterns evoked by rapid sequences of wide-
719 and narrowband clicks in the primary auditory cortex of cat," *J. Neurophysiol.* **84**, 236-246.
- 720 Mehta, A. H., Lu, H., and Oxenham, A. J. (2020). "The perception of multiple simultaneous pitches
721 as a function of number of spectral channels and spectral spread in a noise-excited envelope
722 vocoder," *J. Assoc. Res. Otolaryngol.* **21**, 61-72.
- 723 Meng, Q., Yu, G., Wan, Y., Kong, F., Wang, X., and Zheng, N. (2018). "Effects of vocoder
724 processing on speech perception in reverberant classrooms," in *2018 Asia-Pacific Signal
725 and Information Processing Association Annual Summit and Conference (APSIPA ASC)*
726 (IEEE), pp. 761-765.
- 727 Meng, Q. L., Wang, X. R., Cai, Y. X., Kong, F. H., Buck, A. N., Yu, G. Z., Zheng, N. H., and
728 Schnupp, J. W. H. (2019). "Time-compression thresholds for Mandarin sentences in
729 normal-hearing and cochlear implant listeners," *Hear. Res.* **374**, 58-68.
- 730 Mesnildrey, Q., Hilkuysen, G., and Macherey, O. (2016). "Pulse-spreading harmonic complex as
731 an alternative carrier for vocoder simulations of cochlear implants," *J. Acoust. Soc. Am.*
732 **139**, 986-991.
- 733 Nizami, L., Reimer, J. F., and Jesteadt, W. (2001). "The intensity-difference limen for Gaussian-
734 enveloped stimuli as a function of level: Tones and broadband noise," *J. Acoust. Soc. Am.*
735 **110**, 2505-2515.
- 736 O'Neill, E. R., Kreft, H. A., and Oxenham, A. J. (2019). "Speech perception with spectrally non-
737 overlapping maskers as measure of spectral resolution in cochlear implant users," *J. Assoc.
738 Res. Otolaryngol.* **20**, 151-167.
- 739 Schneider, B. A., Pichora-Fuller, M. K., Kowalchuk, D., and Lamb, M. (1994). "Gap detection and
740 the precedence effect in young and old adults," *J. Acoust. Soc. Am.* **95**, 980-991.
- 741 Shannon, R. V., Zeng, F. G., Kamath, V., Wygonski, J., and Ekelid, M. (1995). "Speech recognition
742 with primarily temporal cues," *Science* **270**, 303-304.
- 743 Shannon, R. V., Zeng, F. G., and Wygonski, J. (1998). "Speech recognition with altered spectral
744 distribution of envelope cues," *J. Acoust. Soc. Am.* **104**, 2467-2476.
- 745 Singh, S., Kong, Y. Y., and Zeng, F. G. (2009). "Cochlear implant melody recognition as a function
746 of melody frequency range, harmonicity, and number of electrodes," *Ear Hear.* **30**, 160-168.
- 747 Skinner, M. W., Holden, L. K., Holden, T. A., Dowell, R. C., Seligman, P. M., Brimacombe, J. A.,
748 Beiter, A. L. J. E., and Hearing (1991). "Performance of postlinguistically deaf adults with

749 the wearable speech processor (WSP III) and mini speech processor (MSP) of the nucleus
750 multi-electrode cochlear implant," *Ear Hear.* **12**, 3-22.

751 Skinner, M. W., Holden, L. K., Whitford, L. A., Plant, K. L., Psarros, C., and Holden, T. A. (2002).
752 "Speech recognition with the nucleus 24 SPEAK, ACE, and CIS speech coding strategies
753 in newly implanted adults," *Ear Hear.* **23**, 207-223.

754 Souza, P., and Rosen, S. (2009). "Effects of envelope bandwidth on the intelligibility of sine- and
755 noise-vocoded speech," *J. Acoust. Soc. Am.* **126**, 792-805.

756 Stam, L., Goverts, S. T., and Smits, C. (2019). "Effect of cochlear implant n-of-m strategy on signal-
757 to-noise ratio below which noise hinders speech recognition," *J. Acoust. Soc. Am.* **145**,
758 EL417-422.

759 Strydom, T., and Hanekom, J. J. (2011). "An analysis of the effects of electrical field interaction
760 with an acoustic model of cochlear implants," *J. Acoust. Soc. Am.* **129**, 2213-2226.

761 Svirsky, M. A., Capach, N. H., Neukam, J. D., Azadpour, M., Sagi, E., Hight, A. E., Glassman, E.
762 K., Lavender, A., Seward, K. P., Miller, M. K., Ding, N., Tan, C. T., and Fitzgerald, M. B.
763 (2021). "Valid Acoustic Models of Cochlear Implants: One Size Does Not Fit All," *Otol.*
764 *Neurotol.* **42**, S2-S10.

765 Tong, Y. C., Clark, G. M., Seligman, P. M., and Patrick, J. F. (1980). "Speech processing for a
766 multiple-electrode cochlear implant hearing prosthesis," *J. Acoust. Soc. Am.* **68**, 1897-1898.

767 Trehub, S. E., Schneider, B. A., and Henderson, J. L. (1995). "Gap detection in infants, children,
768 and adults," *J. Acoust. Soc. Am.* **98**, 2532-2541.

769 van Schijndel, N. H., Houtgast, T., and Festen, J. M. (1999). "Intensity discrimination of Gaussian-
770 windowed tones: Indications for the shape of the auditory frequency-time window," *J.*
771 *Acoust. Soc. Am.* **105**, 3425-3435.

772 Vandali, A. E., Whitford, L. A., Plant, K. L., and Clarke, G. M. (2000). "Speech perception as a
773 function of electrical stimulation rate: Using the nucleus 24 cochlear implant system," *Ear*
774 *Hear.* **21**, 608-624.

775 Wilson, B. S., Finley, C. C., Lawson, D. T., Wolford, R. D., Eddington, D. K., and Rabinowitz, W.
776 M. (1991). "Better speech recognition with cochlear implants," *Nature* **352**, 236-238.

777 Winn, M. B. (2020). "Accommodation of gender-related phonetic differences by listeners with
778 cochlear implants and in a variety of vocoder simulations," *J. Acoust. Soc. Am.* **147**, 174-
779 190.

- 780 Winn, M. B., and Nelson, P. B. (2021). "Cochlear Implants," (Oxford University Press).
- 781 Wong, L. L., Soli, S. D., Liu, S., Han, N., and Huang, M. W. (2007). "Development of the Mandarin
782 Hearing in Noise Test (MHINT)," *Ear Hear.* **28**, 70S-74S.
- 783 Xi, X., Ching, T. Y. C., Ji, F., Zhao, Y., Li, J. N., Seymour, J., Hong, M. D., Chen, A. T., and
784 Dillon, H. (2012). "Development of a corpus of Mandarin sentences in babble with
785 homogeneity optimized via psychometric evaluation," *Int. J. Audiol.* **51**, 399-404.
- 786 Xu, L., Thompson, C. S., and Pfingst, B. E. (2005). "Relative contributions of spectral and temporal
787 cues for phoneme recognition," *J. Acoust. Soc. Am.* **117**, 3255-3267.
- 788 Zeng, F. G., Rebscher, S., Harrison, W. V., Sun, X., and Feng, H. (2008). "Cochlear Implants:
789 System Design, Integration and Evaluation," *IEEE Rev. Biomed. Eng.* **1**, 115-142.
- 790

791 **TABLEs**

792 **TABLE I.** Cutoff frequencies of the band-pass filters in Exp. 1 according to a Greenwood map

Band number	Cutoff frequencies (Hz)
2	80, 1250, and 7999
4	80, 424, 1250, 3234, and 7999
8	80, 215, 424, 748, 1250, 2028, 3234, 5103, and 7999
16	80, 140, 215, 308, 424, 568, 748, 972, 1250, 1597, 2028, 2565, 3234, 4067, 5103, 6393, and 7999
32	80, 108, 140, 176, 215, 259, 308, 363, 424, 492, 568, 653, 748, 854, 972, 1103, 1250, 1414, 1597, 1801, 2028, 2282, 2565, 2881, 3234, 3628, 4067, 4556, 5103, 5713, 6393, 7152, and 7999

793

794

795 **TABLE II.** Results of multiple comparisons between vowel recognition scores with five band
796 numbers for each of the five vocoders.

vocoder→ number of bands pair↓	sine-wave		noise-separate		noise-spread		GET-separate		GET-spread	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
2 vs. 4	*37.0	0.014	24.1	0.140	13.6	1.000	11.0	0.827	*13.3	0.019
2 vs. 8	**57.0	0.001	**56.6	0.002	31.0	0.085	31.6	0.082	*20.1	0.048
2 vs. 16	***66.1	<0.001	***58.9	<0.001	**44.9	0.001	***62.3	<0.001	*24.3	0.005
2 vs. 32	***69.7	<0.001	***68.7	<0.001	**50.4	0.001	***71.9	<0.001	*25.4	0.012
4 vs. 8	20	0.076	*32.4	0.019	17.4	0.057	20.6	0.209	6.9	0.564
4 vs. 16	*29.1	0.049	*34.7	0.023	***31.3	<0.001	**51.3	0.002	*11.0	0.034
4 vs. 32	*32.7	0.037	*44.6	0.011	***36.9	<0.001	***60.9	<0.001	12.1	0.296
8 vs. 16	9.1	0.199	2.29	1.000	13.9	0.147	30.7	0.058	4.1	1.000
8 vs. 32	12.7	0.288	12.1	0.696	19.4	0.059	*40.3	0.011	5.3	1.000
16 vs. 32	3.6	1.000	9.86	0.131	5.6	0.107	9.6	0.455	1.1	1.000

797

798

799

800

801 **TABLE III.** Results of multiple comparisons between vowel recognition scores with five vocoders
 802 for each of the three band numbers (8, 16, and 32).

number of bands→ vocoder pair↓	8		16		32	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
sine-wave vs. noise-separate	-2.1	1.000	4.7	1.000	-1.6	1.000
sine-wave vs. noise-spread	23.4	0.054	**18.7	0.009	16.7	0.120
sine-wave vs. GET-separate	25.6	0.065	4.0	1.000	-2.0	1.000
sine-wave vs. GET-spread	**42.4	0.002	***47.4	<0.001	***49.9	<0.001
noise-separate vs. noise-spread	*25.6	0.030	*14.0	0.025	*18.3	0.048
noise-separate vs. GET-separate	27.7	0.080	-0.7	1.000	-0.4	1.000
noise-separate vs. GET-spread	**44.6	0.004	***42.7	<0.001	***51.4	<0.001
noise-spread vs. GET-separate	2.1	1.000	-14.7	0.130	-18.7	0.102
noise-spread vs. GET-spread	19.0	0.459	**28.7	0.002	**33.1	0.004
GET-separate vs. GET-spread	16.9	0.156	**43.4	0.001	***51.9	<0.001

803

804 **TABLE IV.** Results of multiple comparisons between consonant recognition scores with five
 805 number of bands for each of the five vocoders

vocoder→ number of bands pair↓	sine-wave		noise-separate		noise-spread		GET-separate		GET-spread	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
2 vs. 4	*37.0	0.014	24.1	0.140	13.6	1.000	11.0	0.827	*13.3	0.019
2 vs. 8	**57.0	0.001	**56.6	0.002	31.0	0.085	31.6	0.082	*20.1	0.048
2 vs. 16	***66.1	<0.001	***58.9	<0.001	**44.9	0.001	***62.3	<0.001	**24.3	0.005
2 vs. 32	***69.7	<0.001	***68.7	<0.001	**50.4	0.001	***71.9	<0.001	*25.4	0.012
4 vs. 8	20.0	0.076	*32.4	0.019	17.4	0.057	20.6	0.209	6.9	0.564
4 vs. 16	*29.1	0.049	*34.7	0.023	***31.3	<0.001	**51.3	0.002	*11.0	0.034
4 vs. 32	*32.7	0.037	*44.6	0.011	***36.9	<0.001	***60.9	<0.001	12.1	0.296
8 vs. 16	9.1	0.199	2.3	1.000	13.9	0.147	30.7	0.058	4.1	1.000
8 vs. 32	12.7	0.288	12.1	0.696	19.4	0.059	40.3	0.011	5.3	1.000
16 vs. 32	3.6	1.000	9.9	0.131	5.6	0.107	9.6	0.455	1.1	1.000

806

807

808

809

810

811

812 **TABLE V.** Results of multiple comparisons between consonant recognition scores with five
 813 vocoders for each of the five band numbers

number of bands→ vocoder pair↓	2		4		8		16		32	
	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>	mean diff	<i>p</i>
sine vs. noi-sep	-4.9	1.000	-3.7	1.000	-3.0	1.000	-0.1	1.000	-0.1	1.000
sine vs. noi-spr	-1.6	1.000	13.7	0.267	21.7	0.292	*24.7	0.016	*24.3	0.022
sine vs. GETsep	30.0	0.065	**35.1	0.002	18.1	0.096	-1.0	1.000	2.4	1.000
sine vs. GETspr	27.6	0.108	**33.1	0.010	**49.0	0.006	**54.4	0.002	**59.6	0.002
noi-sep vs. noi-spr	3.3	1.000	17.4	0.157	24.7	0.111	*24.9	0.021	*24.4	0.038
noi-sep vs. GETsep	**34.9	0.007	**38.9	0.007	21.1	0.115	-0.9	1.000	2.6	1.000
noi-sep vs. GETspr	*32.4	0.018	**36.9	0.002	**52.0	0.007	**54.6	0.003	**59.7	0.004
noi-spr vs. GETsep	*31.6	0.043	**21.4	0.003	-3.6	1.000	**25.7	0.010	*21.9	0.040
noi-spr vs. GETspr	29.1	0.078	**19.4	0.003	27.3	0.065	**29.7	0.006	***35.3	<0.001
GETsep vs. GETspr	-2.4	1.000	-2.0	1.000	*30.9	0.034	**55.4	0.002	**57.1	0.003

814

815

816

817 **FIGURE CAPTIONS**

818 **FIG. 1.** (Color online) A unit-amplitude single pulse with Gaussian-shaped envelope (black line)
819 in both the time (a) and frequency (b) domains. The carrier frequency is 5 kHz (the blue waveform
820 in the left panel and the frequency with maximum amplitude in the right panel). The σ equals to
821 $3/f_c = 0.6$ ms in Eq. (1), producing an effective duration of 0.85 ms and an effective bandwidth of
822 1.2 kHz.

823 **FIG 2.** (Color online) The relationship between the tone carrier frequency and the effective duration $D =$
824 $\sqrt{2}\sigma$ (see Panel **A**) or effective bandwidth $B = 1/D$ (see Panel **B**) of Gaussian-enveloped tones (GETs).
825 All axes are logarithmically scaled. The σ was assumed to be $2/f_c$, $3/f_c$, or $4/f_c$ to demonstrate the effects of
826 different duration of GETs. For certain combinations of f_c and σ , the maximum GET rate that can be
827 transmitted with no temporal interaction between neighboring GETs is $1/D$, which equals in value to the
828 effective bandwidth in Panel B.

829 **FIG. 3.** Block diagrams of conventional channel vocoder (A), the first (B) and second (C) types
830 of GET vocoders. The pulsatile vocoders are using GETs as carriers (the first type; used in Exp. 1)
831 or using a single GET as an impulse response (the second type; used in Exp. 2). The front-end pre-
832 emphasis, bandpass filter, and envelope extraction can be implemented either in the temporal or
833 spectral domain.

834 **FIG. 4.** (Color online) A 100-Hz pulse train, repeating a single pulse every 10 ms, in both the time
835 (left panel) and frequency (right panel) domains. The parameters of the individual pulses are the
836 same as those in Fig. 1.

837 **FIG. 5.** (Color online) Spectrograms of three vowel stimuli encoded by the sine-wave, noise-
838 separate, noise-spread, GET-separate, and GET-spread vocoders with 16 channels.

839 **FIG. 6.** (Color online) Vowel (**A**) and consonant (**B**) recognition as a function of number of bands
840 (channels). Cochlear implant data is adapted from [Friesen et al. \(2001\)](#). Simulation data are
841 averaged from seven normal hearing subjects listening to vocoded speech. For the simulation data,
842 standard errors are indicated by the vertical bars. For the CI data, the bars show the entire ranges
843 of performance across all their 19 participants.

844 **FIG. 7.** Mapping a CI electrodiagram to a sound using the second type GET vocoder. **A.** An artificial
845 10-channel CI electrodiagram, including two pulse sweeps with a 10-ms difference between **a** and

846 **b**, as well as two additional sweeps with a 1-ms difference between **c** and **d**, corresponding to
847 stimulation rates of 100 pps and 1000 pps, respectively. **B**. GETs mimicking the electric pulse trains.
848 **C**. The final GET waveform resulting from the sum of ten band-specific GET trains in **B**.

849 **FIG. 8.** (Color online) Speech stimulus demonstrations for the ACE-GET simulation experiment.
850 Left: Spectrogram; middle: band-specific signal; right: zoom in of the boxed signals. **A**.
851 Spectrogram and ACE electrodiagram of a clear sentence of speech. **B-E**. Spectrogram and band-
852 specific waveforms of vocoded speech using two GET vocoders (GETlargeDR, and GETsmallDR)
853 and two conventional sine-wave vocoders (Sin250 and Sin125), respectively.

854 **FIG. 9.** Results from three speech recognition tasks with two 22-channel sine-wave vocoders
855 (Sin250: 250 Hz cut-off envelope; Sin125: 125 Hz cut-off envelope) and two GET vocoders
856 (GETlargeDR and GETsmallDR; their difference is only in the intensity dynamic range, i.e., 32.7
857 dB and 5.3 dB for GETlargeDR and GETsmallDR respectively) compared with the results of some
858 CI subjects. There were two groups of normal-hearing participants, each with ten participants. One
859 group used Sin250 and GETlargeDR, and the other group used Sin125 and GETsmallDR. **A**. Time-
860 compression threshold results. **B**. Speech reception threshold results of a speech in noise recognition
861 experiment (SSN and babble noise). **C**. Speech recognition scores in reverberation with T60 = 0.3,
862 0.6, and 0.9s. Pairwise comparisons with Bonferroni corrections were examined. In each box, “n.
863 s.” denotes the non-significant difference ($p > 0.05$), otherwise, there was a significant difference.

864