

1 **Revised submission -3**

2
3 **Effects of noise suppression and envelope dynamic range compression on the intelligibility of**
4 **vocoded sentences for a tonal language**

5
6 **Fei Chen** ^{a)}

7 *Department of Electrical and Electronic Engineering, Southern University of Science and*
8 *Technology, Xueyuan Road 1088#, Xili, Nanshan District, Shenzhen 518055, China*

9
10 **Dingchang Zheng**

11 *Health and Wellbeing Academy, Faculty of Medical Science, Anglia Ruskin University, Chelmsford,*
12 *UK*

13
14 **Yu Tsao**

15 *Research Center for Information Technology Innovation, Academia Sinica, Taipei, Taiwan*

16
17 Wednesday, August 2, 2017

18
19 ^{a)} Author to whom correspondence should be addressed:

20 Fei Chen, Ph. D.

21 Department of Electrical and Electronic Engineering

22 Southern University of Science and Technology

23 Xueyuan Road 1088#, Xili, Nanshan District, Shenzhen 518055, China

24 Phone: 0086-755-88018554 Email: fchen@sustc.edu.cn

25
26 PACS number: 43.71.Gv, 43.71.Es

27 Keywords: Noise suppression, envelope dynamic range compression.

28
29 Running title: Intelligibility of vocoded sentences

31 **Abstract:** Vocoder simulation studies have suggested that the carrier signal type employed affects
32 the intelligibility of vocoded speech. The present work further assessed how carrier signal type
33 interacts with additional signal processing, namely, single-channel noise suppression and envelope
34 dynamic range compression, in determining the intelligibility of vocoder simulations. In
35 Experiment 1, Mandarin sentences that had been corrupted by speech spectrum-shaped noise (SSN)
36 or two-talker babble (2TB) were processed by one of four single-channel noise-suppression
37 algorithms before undergoing tone- (TV) or noise-vocoded (NV) processing. In Experiment 2,
38 dynamic ranges of multiband envelope waveforms were compressed by scaling of the
39 mean-removed envelope waveforms with a compression factor before undergoing TV or NV
40 processing. TV Mandarin sentences yielded higher intelligibility scores with normal-hearing (NH)
41 listeners than did noise-vocoded sentences. The intelligibility advantage of noise-suppressed
42 vocoded speech depended on the masker type (SSN vs. 2TB). NV speech was more negatively
43 influenced by envelope dynamic range compression than was TV speech. These findings suggest
44 that an interactional effect exists between the carrier signal type employed in the vocoding process
45 and envelope distortion caused by signal processing.

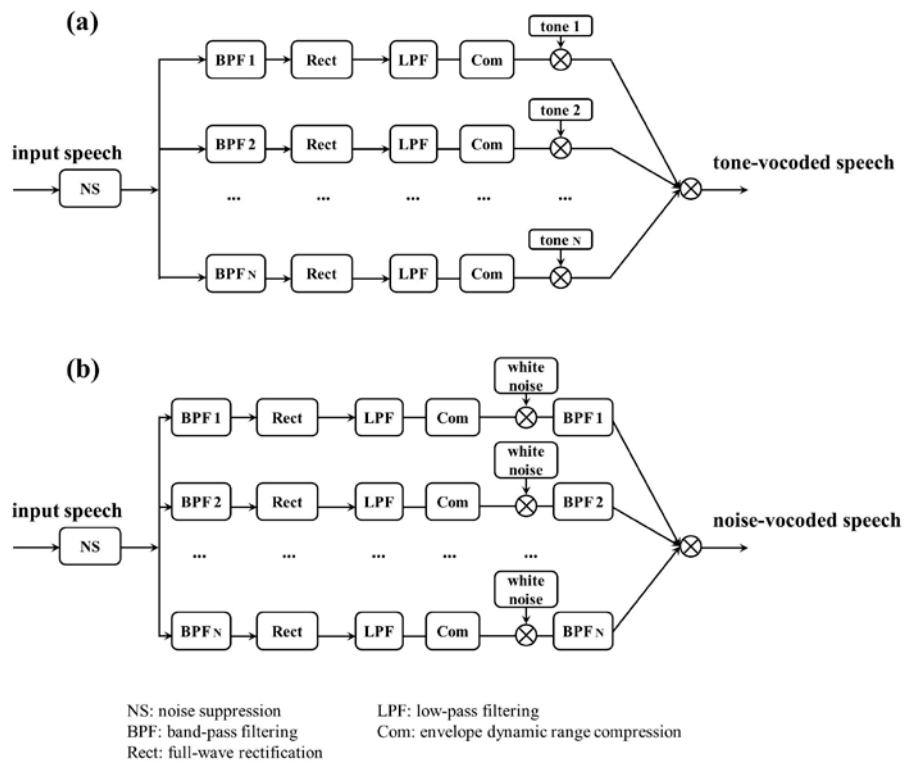
46 I. INTRODUCTION

47 The perceptual contribution of the temporal envelope has attracted enduring research interest.
48 Many studies have assessed the importance of the temporal envelope for speech intelligibility
49 under various conditions (e.g., Shannon *et al.*, 1995; Dorman *et al.* 1997; Chen and Loizou,
50 2011a). Vocoder simulations have long been used to extract the multiband temporal envelope
51 waveforms while removing the underlying fine-structure information to synthesize
52 envelope-based vocoded speech (e.g., Shannon *et al.*, 1995; Dorman *et al.* 1997; Chen and Loizou,
53 2011a). With envelope information from up to four bands, normal-hearing (NH) listeners can have
54 near-perfect speech understanding in quiet condition (Shannon *et al.*, 1995).

55 In a cochlear implant (CI) device, incoming sound signals are received via a microphone and
56 fed into a speech processor. Most of the existing CI speech processors capture multi-channel
57 temporal envelopes of sound signal inputs, and then generate electric stimulations that excite
58 patients' residual auditory nerves directly. Vocoder simulations aim to transfer only those acoustic cues that
59 are present for CI users, so they simulate the signal processing of a CI. Vocoder simulations have
60 been applied to examine numerous factors that influence the intelligibility of envelope-based
61 vocoded speech, including the number of channels (Shannon *et al.*, 1995; Dorman *et al.* 1997),
62 carrier signal type (Dorman *et al.* 1997; Fu *et al.*, 2004; Gonzalez and Oliver, 2005; Whitmal, *et*
63 *al.*, 2007; Chen and Lau, 2014), envelope cutoff frequency (Shannon *et al.* 1995; Xu *et al.* 2005;
64 Souza and Rosen 2009), and frequency spacing (Kasturi and Loizou, 2007), among other factors.
65 For this reason, vocoder simulations have been used widely to assess the potential of new
66 speech-processing and coding strategies for CIs before large-scale clinical evaluations with users
67 are conducted. Vocoder simulation remains a valuable tool in the field of CI research because it
68 can be used to assess the effects of acoustic factors in the absence of patient-specific confounds.

69 When performing vocoder simulations, the envelope waveform is extracted by steps of
70 bandpass filtering (BPF), waveform rectification and low-pass filtering (LPF) (see Fig. 1). The
71 envelope waveform is used to modulate a carrier signal. There are two common types of carrier
72 signals used in synthesizing vocoded speech; pure-tone and white-noise signals yield tone- (TV)
73 and noise-vocoded (NV) speech stimuli, respectively. A limited number of studies have compared
74 the relative performance of these two vocoder types on speech intelligibility in English (e.g.,
75 Dorman *et al.* 1997; Whitmal, *et al.*, 2007; Souza and Rosen, 2009; Rosen *et al.*, 2015) and on the

76 listener's ability to distinguish gender and speaker identity (in English, Fu *et al.*, 2004; in Spanish,
 77 Gonzalez and Oliver, 2005). Dorman *et al.* (1997) compared English speech intelligibility using a
 78 tone or noise vocoder with varying numbers of channels and found only small differences that did
 79 not reach statistical significance under most test conditions with vowels, consonants, and
 80 sentences. Their findings suggested that neither of the two vocoder types was superior to the other.
 81 However, in a more recent study, Whitmal *et al.* (2007) examined the intelligibility of English
 82 sentences and vowel-consonant-vowel syllables using a six-band vocoder and found that a tone
 83 vocoder produced more intelligible speech than a noise vocoder under both quiet and noisy
 84 conditions across different signal-to-noise ratios (SNRs) and two different masker types, namely,
 85 speech spectrum-shaped noise (SSN) and two-talker babble (2TB). In their study on the
 86 interaction between carrier type and cutoff frequency in the vocoding process, Souza and Rosen
 87 (2009) found that TV speech was less intelligible than NV speech for a low envelope cutoff
 88 frequency of 30 Hz, but more intelligible for a high envelope cutoff of 300 Hz. And Rosen *et al.*
 89 (2015) reported recently that using tone carriers with a denser spectrum improved the
 90 intelligibility of TV speech considerably over typical tone vocoders, equating and even surpassing
 91 the performance observed with noise vocoders.



92

93

FIG 1. Block diagrams of (a) tone-vocoder and (b) noise-vocoder processes.

94 Studies on gender and speaker identification, in which good performance depends heavily on
95 cues such as fundamental frequency (F0) and formant structure, have shown better performance
96 for TV speech than for NV speech. Using one- and four-band noise vocoders, Fu *et al.* (2004)
97 observed poor voice gender discrimination (approximately chance level). However, with a tone
98 vocoder, they obtained better results that were more consistent with those of real CI users. In
99 another study of gender and speaker identification in Spanish, Gonzalez and Oliver (2005) found
100 that the tone vocoder performed substantially better than did the noise vocoder across conditions
101 with different numbers of channels. Recently, Chen and Lau (2014) evaluated the effect of
102 vocoder carrier signal type on the intelligibility of Mandarin Chinese, a tonal language, and found
103 an advantage of tone over noise carriers on the intelligibility of vocoded Chinese speech.

104 Noise is prevalent in our daily lives and poses a great challenge to human speech perception.
105 Alleviation of background noise interference is the goal of many single-channel noise-suppression
106 algorithms, such as the spectral-subtraction (Kamath and Loizou 2002), statistical-model-based
107 (Ephraim and Malah, 1985), and subspace (Hu and Loizou, 2003) algorithms. However, noise
108 suppression may cause undesirable distortion (e.g., “musical noise”) of speech, which is
109 detrimental to speech perception (Loizou, 2007). Certain noise-suppression algorithms (e.g.,
110 statistically based Wiener filtering) have been shown to improve speech quality *per se* without
111 improving speech intelligibility for NH listeners (Hu and Loizou, 2007; Li *et al.*, 2011).

112 Vocoded speech is synthesized with multiband envelope information from the original speech
113 signal. When the original speech signal contains distortions due to noise-reduction processing, it
114 is unclear how this distortion affects the intelligibility of envelope-based vocoded speech. In
115 Experiment 1, our aim was to investigate whether the intelligibility advantage of tone over noise
116 vocoders persists when noise-suppression processing is used. Comparisons between tone and
117 noise vocoders have shown no inherent fluctuation in a tone carrier compared to a noise carrier.
118 The combination of envelope distortion (caused by noise-suppression processing) and inherent
119 fluctuation in a noise carrier may have a negative influence on the intelligibility of NV speech.
120 Hence, we hypothesized that carrier signal type may affect the intelligibility of vocoded speech in
121 the context of noise-suppression processing. In other words, we are supposing that the
122 intelligibility advantage of tone over noise vocoders may occur when the vocoding process
123 involves noise-suppression processing. In addition, the effect of noise masking is commonly

124 represented by two mechanisms: energetic masking by steady-state noise, and informational
125 masking by other speech characteristics, such as fluctuations, in the noise (Carhart *et al.*, 1967;
126 Watson, 2005). Hence, in Experiment 1, we also examined whether the interactional effect of the
127 vocoder type and noise-suppression processing depends on the masker type.

128 Dynamic range plays an important role in speech perception (e.g., Zeng *et al.*, 2002). This fact
129 provides a partial explanation for why CI users have poor speech perception (i.e., reduced hearing
130 dynamic range of 5–10 dB), especially in adverse listening environments. Fitting the wide
131 dynamic range of speech signals into the narrow range of the residual hearing of CI users requires
132 dynamic range compression. Several vocoder simulation studies have assessed the effect of
133 envelope dynamic range on speech intelligibility (Fu and Shannon, 1999; Loizou *et al.*, 2000;
134 Chen *et al.*, 2013; Lai *et al.*, 2015). Similarly, we are interested in clarifying whether reducing the
135 dynamic range has a negative effect on envelope-based vocoded speech, and to what extent
136 carrier signal type affects the intelligibility of vocoded speech with a dynamic range-compressed
137 envelope. In Experiment 2, our aim was to investigate whether the intelligibility advantage of tone
138 over noise vocoders persists in the context of envelope dynamic range compression. Earlier work
139 has shown that the spectral sidebands contained in TV speech (due to the multiplication of
140 pure-tone carrier and envelope waveform) carries additional cue which is beneficial for speech
141 intelligibility (e.g., Whitmal *et al.*, 2007; Stone *et al.*, 2008). In addition, white-noise carrier has
142 intrinsic envelope fluctuations that are absent in pure-tone carrier. Multiplying the white-noise
143 carrier by the envelope waveform may have an additional temporal influence on the envelope
144 waveform, which is detrimental to speech understanding (Stone *et al.*, 2011). Given the potential
145 negative effect of dynamic range compression and the intelligibility disadvantage of NV relative
146 to TV speech, this work hypothesized that when envelope dynamic range compression is included
147 in the vocoding process, the intelligibility of NV speech would drop at a higher rate than that of
148 TV speech. NV speech would be far less intelligible than TV speech.

149

150 **II. EXPERIMENT 1: EFFECT OF NOISE SUPPRESSION ON THE INTELLIGIBILITY** 151 **OF VOCODED SENTENCES**

152 The purpose of Experiment 1 was to examine the effect of noise suppression on the
153 intelligibility of TV and NV Mandarin sentences.

154

155 **A. Methods**

156 **1. Subjects**

157 Eight (five males and three females) native-Mandarin-Chinese listeners (18–23 years old)
158 participated in the experiment. All participants were undergraduate students at Southern
159 University of Science and Technology, and were paid for their participation. All subjects had NH,
160 as determined by having measured pure-tone thresholds (250–8000 Hz) better or equal to a 20-dB
161 hearing level. The study protocol was approved by the Human Research Ethics Committee for
162 Non-Clinical Faculties of Southern University of Science and Technology.

163

164 **2. Materials**

165 The speech material consisted of sentences taken from the Mandarin Hearing in Noise Test
166 (MHINT) database (Wong *et al.*, 2007), which includes 24 lists of 10 sentences, with each
167 sentence containing 10 key words. All of the sentences were produced by a male speaker with an
168 F0 range of 75–180 Hz.

169 Two types of masking were used to corrupt the sentences: steady-state SSN and 2TB. For SSN
170 masking, a finite impulse response filter was designed based on the average spectrum of the
171 MHINT sentences, and a white noise was filtered and scaled to the same long-term average
172 spectrum and level as the sentences. The 2TB masker contained two equal-level interfering male
173 talkers. A random noise segment of the same length as the clean speech signal was cut out of the
174 noise recordings, appropriately scaled to reach the desired input SNR level, and finally added to
175 the speech signals at -2-dB and 6-dB input SNR levels for the SSN and 2TB maskers, respectively.
176 The input SNR levels were chosen based on known performance from a pilot study.

177

178 **3. Signal processing**

179 The noise-suppressed vocoded speech generation processes are summarized in block diagrams
180 in Figure 1. Input noise-corrupted speech signals were first processed by existing single-channel
181 noise-suppression algorithms, followed by the tone- or noise-vocoding process. To process
182 noise-corrupted sentences, we used four representative noise-suppression algorithms: the
183 generalized Karhunen-Loeve transform (KLT) approach (Hu and Loizou, 2003), the Log

184 Minimum Mean Square Error (logMMSE) algorithm (Ephraim and Malah, 1985), the multiband
185 spectral subtraction (MB) algorithm (Kamath and Loizou, 2002), and the Wiener algorithm based
186 on *a priori* SNR estimation (Scalart and Filho, 1996). These four algorithms encompass the four
187 most commonly used types of single-channel noise-suppression methods, namely the subspace,
188 statistical-modeling, spectral-subtraction, and Wiener-filtering approaches, respectively (see
189 review in Loizou, 2007).

190 For the KLT method, the noise-corrupted speech signal is projected into orthogonal subspaces;
191 KLT parts representing the signal subspace are modified by a gain function, determined by the
192 estimator; remaining KLT parts representing the noise subspace are nulled; and the enhanced
193 signal is obtained from the inverse KLT of the modified parts (Hu and Loizou, 2003). The
194 statistical-modeling approach employs statistical models with optimization criteria (e.g.,
195 minimum mean square error) to estimate the magnitude spectrum of the speech signal (Ephraim
196 and Malah, 1985). The spectral-subtractive algorithm is implemented with an estimate of the
197 clean signal spectrum, generated by subtracting an estimate of the noise spectrum from a
198 noise-corrupted speech spectrum (Kamath and Loizou 2002). The Wiener filter uses *a priori* SNR
199 statistics to design a gain function that suppresses low-SNR segments, while preserving high-SNR
200 ones. Detailed descriptions of the algorithms including the exact parameters used in the current
201 study can be found in Hu and Loizou (2007) and Loizou (2007). The Matlab code used to
202 implement the four noise-suppression algorithms was obtained from Loizou (2007).

203 All noise-suppressed materials were further processed by a tone or noise vocoder (Fig. 1). To
204 implement the tone vocoder, speech signals were first processed through a pre-emphasis filter
205 (first-order high-pass filter with 1200-Hz cutoff frequency). Then, signals were bandpass-filtered
206 into eight frequency bands between 80 Hz and 6000 Hz with sixth-order Butterworth filters. The
207 cutoff frequencies for the channel allocation of bandpass filters were (in Hz): 80, 221, 426, 724,
208 1158, 1790, 2710, 4050, and 6000. From each band, the envelope was extracted by full-wave
209 rectification and low-pass filtering with a 200-Hz cutoff frequency by way of a fourth-order
210 Butterworth filter. Sine waves at the center frequencies of the bandpass filters were generated
211 with amplitudes modulated by the extracted envelopes. All amplitude-modulated sine waves from
212 the resultant set of bands were summed to generate a TV stimulus, whose amplitude was adjusted
213 to have the same root-mean-square (RMS) energy as the original speech signal. RMS energy

214 scaling was performed with respect to the noisy and noise-suppressed input speech signals under
215 the noisy and noise-suppressed conditions, respectively. Noise-suppression processing causes an
216 RMS energy difference between noisy and noise-suppressed speech signals. The RMS energy
217 scaling was done with respect to the energy of each original speech signal. Experimental results
218 may vary when RMS energy is scaled with respect to the same energy (of either the noisy or
219 noise-suppressed speech signal); this possibility warrants further investigation.

220 Implementation of the noise vocoder was similar to that of the tone vocoder, except that a
221 white noise instead of a sine wave was used as the carrier signal, and amplitude-modulated by the
222 extracted envelope. Output from each band was further band-limited with the same bandpass filter
223 at that band. All amplitude-modulated noises (with band-limiting processing) were summed to
224 generate the NV stimulus, with its amplitude adjusted to have the same RMS power as the
225 original signal. Again, RMS energy scaling was performed with respect to the noisy and
226 noise-suppressed input speech signals under the noisy and noise-suppression conditions,
227 respectively. The envelope dynamic compression block (labeled 'Com' in Fig. 1) was deactivated
228 (compression factor $\alpha = 1$; see Experiment 2) in the vocoding process.

229

230 **4. Procedure**

231 The experiment was performed in a sound booth, and stimuli were played to listeners
232 diotically through an HD 650 circumaural headphone (Sennheiser, Germany) set at a comfortable
233 listening level. Before the actual testing session, each subject participated in a 10-min training
234 session and was given four lists of 10 MHINT sentences. The training session familiarized the
235 subjects with the testing procedure and conditions. During the training session, the subjects were
236 allowed to read transcriptions of the training sentences while they were listening to the sentences.
237 Four testing conditions [= 2 masker types (i.e., SSN at -2 dB SNR and 2TB at 6 dB SNR) \times 2
238 vocoder types (i.e., TV and NV) \times 1 signal processing condition (i.e., noisy)] were used during
239 training. In the testing session, the order of the conditions was randomized across subjects, and
240 the subjects were asked to repeat orally all of the words they heard. In addition, the lists were
241 randomized across listeners. The sentences used during testing were not the same as any of the
242 training sentences. Each subject participated in a total of 20 conditions [= 2 masker types (i.e.,
243 SSN at -2 dB SNR and 2TB at 6 dB SNR) \times 2 vocoder types (i.e., TV and NV) \times 5 signal

244 processing conditions (i.e., KLT, logMMSE, Wiener, MB, and noisy)]. One list of 10 Mandarin
245 sentences was used per tested condition, and none of the sentences was repeated across conditions.
246 Subjects were allowed to listen to each stimulus a maximum of three times, and were asked to
247 repeat as many words as they could recognize. A simple custom software interface was designed
248 for the listening experiment, which each participant used to control the auditory delivery of the
249 processed stimuli. During the testing session, a tester accompanied the participant and scored
250 his/her response in the computer. A 5-minute break was given every 30 minutes to avoid listening
251 fatigue. The intelligibility score for each condition was computed as the ratio between the number
252 of correctly recognized words and the total number of words contained in each MHINT list. The
253 total testing time was one hour and ten minutes (10-minute training and 60-minute testing).

254

255 5. *Data analysis*

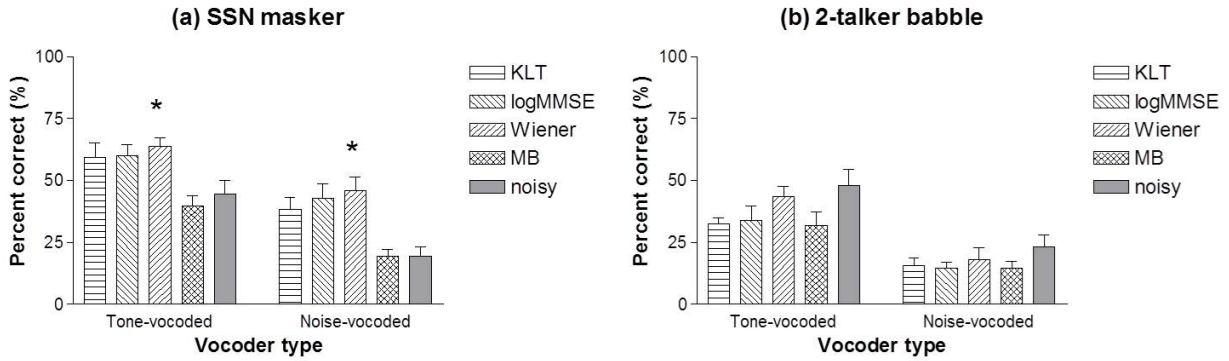
256 The data were subjected to two-way repeated measures analyses of variance (rmANOVAs)
257 with recognition score as the dependent variable and vocoder type and signal processing condition
258 as within-subject factors. Recognition scores were first converted to rational arcsine units using
259 the rationalized arcsine transform (Studebaker, 1985). A one-way rmANOVA was conducted for
260 each type of vocoder to further analyze the effect of signal processing condition; the ANOVA
261 alpha level was Bonferroni corrected, and only those tests with p values lower than 0.0125 (= $0.05/4$)
262 were considered significant. Paired t -tests were conducted in each signal processing
263 condition to further analyze vocoder-type effects.

264

265 B. Results

266 Mean recognition scores for all conditions in Experiment 1 are shown in Figure 2, with data
267 for the SSN and 2TB maskers shown in panels a and b, respectively. For the results of the SSN
268 masker at -2 dB SNR condition (Fig. 2a), a two-way rmANOVA indicated significant effects of
269 vocoder type ($F_{1,7} = 34.13, p < 0.005$) and signal processing condition ($F_{4,28} = 19.83, p < 0.001$),
270 but no significant interaction between these two variables ($F_{4,28} = 0.397, p = 0.81$). One-way
271 rmANOVAs showed significant differences in performance between Wiener-processed and noisy
272 (i.e., no noise suppression) vocoded speech for both vocoder types ($p < 0.01$), and paired t -tests
273 revealed performance differences ($p < 0.05$) between paired TV and NV speech under all signal

274 processing conditions.



275

276 **FIG 2.** Sentence recognition scores for all conditions with and without noise reduction algorithms
277 with (a) a -2-dB SNR SSN masker and (b) a 6-dB SNR 2TB masker. The error bars denote ± 1
278 standard error of the mean. The asterisk denotes that the intelligibility score is significantly
279 ($p < 0.01$) larger than that in the noisy condition.

280

281 For the results of the 2TB masker at 6 dB SNR condition (Fig. 2b), a two-way rmANOVA
282 indicated significant effects of vocoder type ($F_{1,7} = 69.14, p < 0.001$) and signal processing
283 condition ($F_{4,28} = 3.65, p < 0.05$), but not a significant interaction ($F_{4,28} = 0.40, p = 0.81$) between
284 vocoder type and signal processing condition. Again, a one-way rmANOVA revealed no
285 significant performance difference ($p > 0.02$) between noise-suppressed and noisy vocoded
286 speech. Paired t -tests revealed significant performance differences ($p < 0.05$) between paired TV
287 and NV speech under all signal processing conditions.

288

289 **III. EXPERIMENT 2: EFFECT OF ENVELOPE DYNAMIC RANGE COMPRESSION** 290 **ON THE INTELLIGIBILITY OF VOCODED SENTENCES**

291 The purpose of Experiment 2 was to examine the effect of envelope dynamic range
292 compression on the intelligibility of TV and NV Mandarin sentences.

293

294 **A. Methods**

295 **1. Subjects and Materials**

296 Seven (four males and three females, 19–20 years old) new (i.e., did not participate in

297 Experiment 1) NH native-Mandarin listeners participated in this experiment. All participants were
298 undergraduate students at Southern University of Science and Technology, and were paid for their
299 participation.

300 The speech materials were the same as in Experiment 1, and the SSN masker was used to
301 corrupt the MHINT sentences at 3 dB and -3 dB input SNR levels.

302

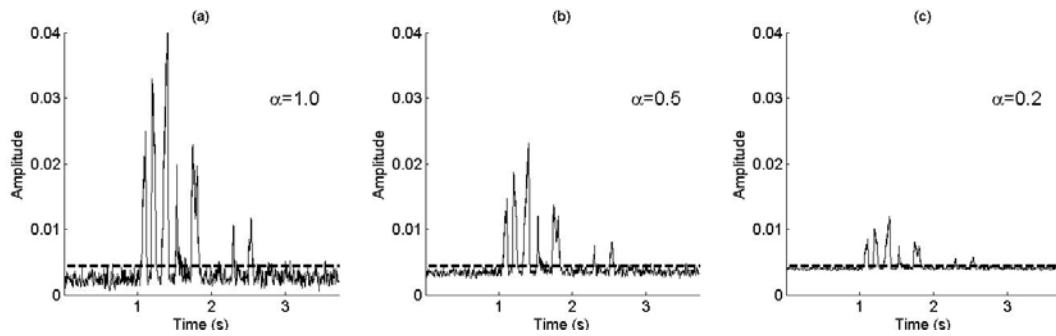
303 2. *Signal processing*

304 We implemented a simple compression method (Chen *et al.*, 2013). Letting x and y denote
305 input and output amplitude envelopes, respectively, the output compressed amplitude envelope y
306 was computed as:

$$307 \quad y = \alpha \times (x - \bar{x}) + \bar{x}, \quad (1)$$

308 where \bar{x} is the mean of the input amplitude envelope x , and α is the compression factor constant
309 chosen for compressing the output amplitude envelope dynamic range. Mean values of the output
310 and input amplitude envelopes were equal (i.e., $\bar{y} = \bar{x}$), regardless of the value of α . A small
311 compression factor α denotes a large compression ratio and *vice versa*. When $\alpha = 0$ in Eq. (1), the
312 compressed amplitude envelope becomes a DC signal with a constant value of \bar{x} (i.e., $\bar{y} = \bar{x}$), and
313 the dynamic range is 0 dB. When $\alpha = 1.0$, the output amplitude envelope maintains the original
314 dynamic range of the input (i.e., no envelope compression). Figure 3 shows the three compressed
315 amplitude envelope waveforms, with compression factor $\alpha = 1.0, 0.5$, and 0.2 , respectively. Note
316 that the three compressed amplitude envelope waveforms have the same mean values (dashed lines
317 in the three panels in Fig. 3). We employed α values of 1, 0.5, and 0.2, which reduced the input
318 envelope dynamic range by 0 dB, 6 dB, and 14 dB, respectively.

319 The compressed envelope was multiplied by the carrier signal (i.e., tone or noise) to generate
320 vocoded stimuli, as in Experiment 1. The noise suppression (NS) block in Fig. 1 was deactivated in
321 the vocoding process. The compression strategy in Eq. (1) was motivated by preserving the
322 loudness of processed speech signals, while reducing the dynamic range of envelope variation
323 selectively (Chen *et al.*, 2013). The compression strategy in Eq. (1) is different from those used in
324 actual CIs, wherein a nonlinear function is used to limit the speech envelope into the range
325 restricted by the threshold and most comfortable levels of a CI listener.



326

327 **FIG 3.** Example waveforms of compressed amplitude envelope with compression factor α of
 328 values (a) 1.0, (b) 0.5, and (c) 0.2. The dashed line in each panel denotes the mean value of the
 329 amplitude envelope waveform.

330

331 3. Procedure

332 The experimental procedure used in Experiment 2 was essentially the same as that used in
 333 Experiment 1. Again, in the training session in which subjects were familiarized with the testing
 334 procedure and conditions, each subject was given four lists of 10 sentences (different from those
 335 used in the testing session) and allowed to read transcriptions while listening to the sentences.
 336 However, in Experiment 2, each subject was exposed to a total of 12 conditions [= 2 input SNR
 337 levels (i.e., 3 dB and -3 dB) \times 2 vocoder types (i.e., TV and NV) \times 3 values of compression factor
 338 (i.e., $\alpha=1.0, 0.5$ and 0.2)], which were randomized across the subjects. As in Experiment 1, one
 339 list of 10 sentences was presented per condition, and none of the sentences was repeated across
 340 the conditions. The total testing time was 50 minutes (10-minute training and 40-minute testing).

341

342 6. Data analysis

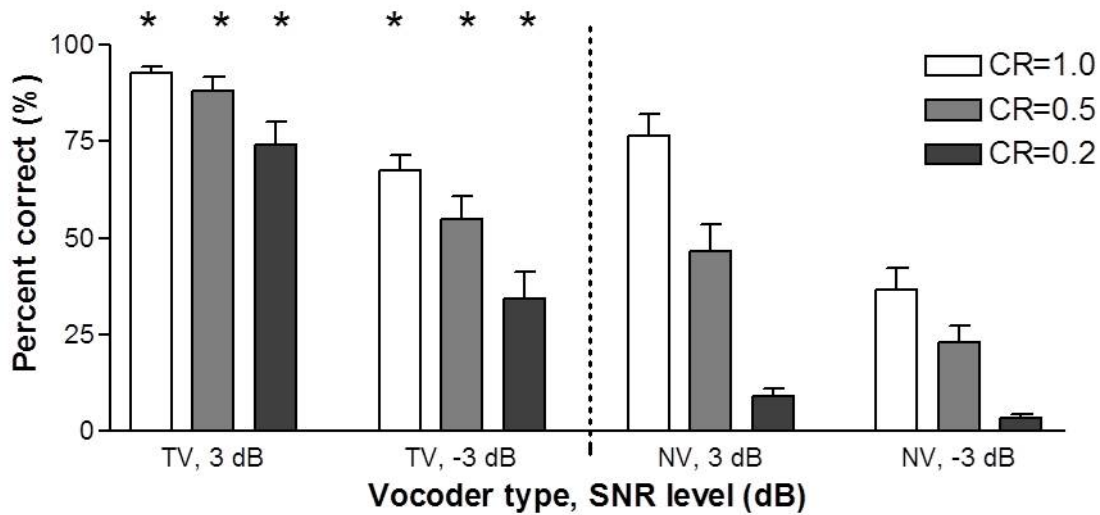
343 The data were analyzed as in Experiment 1. The three within-subject factors for the three-way
 344 rmANOVAs were vocoder type, SNR level and compression factor; and paired t -tests were
 345 conducted in each compression condition to further analyze vocoder-type effects.

346

347 B. Results

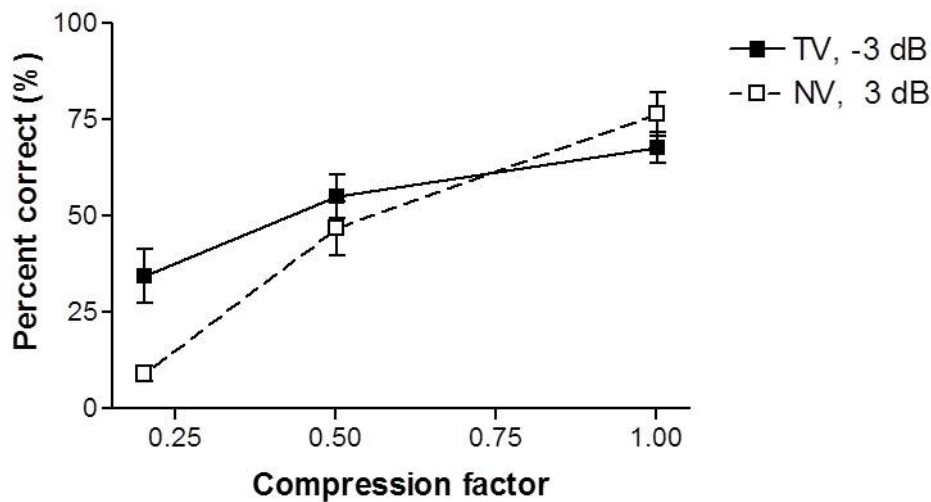
348 The mean recognition scores of Mandarin sentences for all conditions are shown in Figure 4.
 349 A three-way rmANOVA indicated significant effects of vocoder type ($F_{1,6} = 167.47, p < 0.005$),
 350 SNR level ($F_{1,6} = 230.66, p < 0.005$) and compression factor ($F_{2,12} = 107.75, p < 0.005$), as well

351 as a non-significant interaction between vocoder type and SNR level ($F_{1,6} = 5.9, p = 0.06$), a
 352 non-significant interaction between SNR level and compression factor ($F_{2,12} = 153.80, p = 0.3$), a
 353 significant interaction between vocoder type and compression factor ($F_{1,6} = 6.69, p < 0.05$), and a
 354 non-significant interaction among vocoder type, SNR level and compression factor ($F_{2,12} = 3.08,$
 355 $p = 0.09$). Paired t -tests showed that performance differed significantly ($p < 0.001$) between TV
 356 and NV speech under all test conditions with the same SNR level and compression factor value.



357
 358 **FIG 4.** Sentence recognition scores for all test conditions. The error bars denote ± 1 standard error
 359 of the mean. The asterisk denotes that the intelligibility score of tone-vocoded speech is
 360 significantly ($p < 0.005$) larger than that of noise-vocoded speech.

361
 362 The significant interaction between vocoder type and compression factor appears to be due to
 363 the ceiling/flooring effect on the intelligibility scores of TV/NV speech in Fig. 4. To further
 364 analyze the interactional effect between vocoder type and compression factor, Figure 5 displays
 365 the scores of TV and NV speech in near-linear range (as a function of compression factor), and
 366 excludes the effect of ceiling/flooring on data analysis. The SNR levels in Fig. 5 are -3 dB and 3
 367 dB for TV and NV speech, respectively. It is seen in Fig. 5 that at uncompressed condition (i.e.,
 368 CR=1.0), the intelligibility scores of TV and NV speech are similar. Envelope dynamic range
 369 compression causes decreased intelligibility to both TV and NV speech. However, it is noted that
 370 the intelligibility score of NV speech drops at a higher rate than that of TV speech does, indicating
 371 that the effect of compression is different between TV and NV speech.



372

373 **FIG 5.** Sentence recognition scores for selected test conditions in Fig. 4. The solid and dashed
 374 lines show the scores for TV speech at -3 dB and NV speech at 3 dB, respectively. The error bars
 375 denote ± 1 standard error of the mean.

376

377 **IV. DISCUSSION AND CONCLUSIONS**

378 Prior vocoder simulation studies have demonstrated a perceptual contribution of the temporal
 379 envelope to speech intelligibility (e.g., Shannon *et al.* 1995; Dorman *et al.* 1997; Whitmal *et al.*,
 380 2007; Stone *et al.*, 2008; Stone *et al.*, 2011; Chen and Loizou, 2011a). Several factors can be
 381 manipulated to control the amount of information that is included in the multiband envelope. In
 382 the present work, we assessed how noise suppression and envelope dynamic range compression
 383 affect the intelligibility of vocoded speech. We also investigated the effect of carrier signal type
 384 on the intelligibility of noise-suppressed and envelope dynamic range-compressed vocoded
 385 speech.

386

387 **A. Intelligibility advantage of tone- over noise-vocoded speech**

388 In Experiments 1 and 2, TV speech was found to be more intelligible than NV speech under
 389 the same signal-processing conditions, consistent with earlier studies reporting a perceptual
 390 advantage of a tone vocoder over a noise vocoder (e.g., Whitmal *et al.*, 2007; Chen and Lau,
 391 2014). We further showed that this advantage persisted even under conditions of envelope
 392 waveform distortion (i.e., noise suppression and narrowing of the dynamic range). Taken together,

393 these results provide evidence for the notion that there is a perceptual advantage of TV speech
394 not only when processed by the raw vocoder simulation model, but also when there is additional
395 signal processing, such as noise suppression and envelope dynamic range compression.

396 Two mechanisms may account for the perceptual advantage of TV speech. The first potential
397 mechanism concerns the spectral sidebands that are contained in TV speech when a pure tone is
398 multiplied by the envelope waveform (e.g., Whitmal *et al.*, 2007; Stone *et al.*, 2008). The
399 amplitude-modulated tone carrier has two spectral sidebands, and these sidebands impose a
400 periodic temporal structure in voiced speech segments on the tone-vocoder's output, with the
401 talker's pitch being preserved over most voiced segments (Whitmal *et al.*, 2007). Hence, the
402 spectral sidebands contain an additional cue that is beneficial for speech intelligibility, even when
403 the noise-suppressed envelope contains nonlinear distortions due to noise-suppression processing.
404 The second potential mechanism is related to the difference in intrinsic temporal fluctuations
405 between sine-wave and white-noise carriers. White-noise carriers have intrinsic envelope
406 fluctuations that are absent in sine-wave carriers. Hence, the white-noise carrier, when multiplied
407 by the envelope waveform, may have an additional temporal influence on the envelope waveform
408 and cause a detrimental effect on speech understanding (Stone *et al.*, 2011).

409 Another factor that might account for the intelligibility difference observed between TV and
410 NV speech might be the tonal quality of Mandarin Chinese. Mandarin differs from English in that
411 a syllable's tone (or F0 contour) is used to differentiate meaning between otherwise similar lexical
412 items (Howie, 1976; Chen and Loizou, 2011b;). Although the F0 contour is the primary cue for
413 lexical tone identification, the tonal envelope waveform also carries important information for
414 tone identification (Luo and Fu, 2004). In this study, the F0 of the target MHINT sentence ranged
415 from 75 Hz to 180 Hz, and the envelope cutoff frequency was set to 200 Hz. Hence, the envelope
416 waveform and the spectral sidebands of tone-vocoded speech may carry important tonal
417 information. However, noise-vocoded speech, due to its use of noise carriers, may influence or
418 distort the envelope waveform.

419

420 **B. Dependence of masker type on the intelligibility of noise-suppressed vocoded speech**

421 For noise-suppressed wideband speech signals, no improvements in speech intelligibility
422 have been observed for NH listeners (Hu and Loizou, 2007; Li *et al.*, 2011). When the

423 noise-suppressed wideband speech signal was processed by a vocoder, we observed
424 masker-dependent intelligibility performance. With an SSN masker, a single-channel
425 noise-suppression algorithm (i.e., Wiener filtering, see Loizou and Kim, 2011) may improve the
426 intelligibility of vocoded speech, regardless of vocoder type. However, when the masker is
427 competing speech (i.e., 2TB), no intelligibility improvement was observed with processing by
428 existing noise-suppression algorithms.

429 The exact mechanism underlying the presently observed masker-dependent intelligibility
430 performance in vocoded speech is unclear. We hypothesize that noise-suppression processing
431 causes less envelope distortion of speech signals with an SSN masker than with a 2TB masker.
432 When the noise-suppression algorithm was integrated into the vocoding process, this difference in
433 envelope distortion could have accounted for the beneficial effect of noise suppression with
434 steady-state noise corruption and the lack of intelligibility improvement by noise suppression with
435 competing masker corruption. When Chen *et al.* (2015) evaluated the performance of
436 noise-suppression (i.e., the same four single-channel noise-suppression algorithms used in this
437 work) for improving speech recognition by Mandarin-speaking CI users, they tested three types of
438 maskers: SSN, babble, and car noise. They found that although most noise-suppression algorithms
439 could improve Mandarin speech recognition in the presence of noise (e.g., SSN), the algorithms
440 performed differently across different environmental noise conditions. They used an
441 envelope-distortion based objective intelligibility measure (i.e., the normalized covariance
442 measure) to predict CI speech recognition scores and found that an envelope-distortion based
443 intelligibility index could predict the intelligibility of noisy and noise-suppressed speech by CI
444 listeners modestly well (i.e., correlation coefficient 0.81). Similarly, when Baumgärtel *et al.* (2015)
445 evaluated the performance of single-channel noise reduction in the listening scenarios of
446 stationary speech-shaped noise and competing speech, they also found better and worse
447 performance in the stationary noise and competing speech scenarios, respectively. They attributed
448 this masker-differentiated performance to the errors to estimate speech and noise power based on
449 the speech presence probability in single-channel noise reduction processing. Baumgärtel *et al.*
450 noted that in the stationary noise condition, speech and noise power estimates (or the separation of
451 a noisy signal into speech and noise components) worked quite well, whereas in the nonstationary
452 noise (e.g., competing speech) condition, estimation errors occurred and little performance

453 improvement was found. Future studies should investigate the degree of envelope waveform
454 distortion generated by processing with existing single-channel noise-suppression algorithms. In
455 addition, two different SNR levels were used for the SSN and 2TB conditions in Experiment 1,
456 i.e., a negative SNR of -2 dB for SSN and a positive SNR of 6 dB for 2TB. It remains to be
457 resolved how this SNR level difference interacts with masker type in determining the
458 intelligibility of noise-suppressed vocoded speech.

459

460 **C. Influence of envelope dynamic range compression on the intelligibility of vocoded speech**

461 In addition to demonstrating a perceptual advantage of employing a tone carrier over
462 employing a noise carrier in the vocoding process, the present work showed that these two types
463 of vocoded speech were associated with differing responses to envelope dynamic range
464 compression. Dynamic range narrowing has been shown repeatedly to impede speech
465 intelligibility (Fu and Shannon, 1999; Loizou *et al.*, 2000; Chen *et al.*, 2013). Fu and Shannon
466 (1999) measured phoneme recognition in CI users when the dynamic range of the input speech
467 signals was reduced by either peak clipping or center clipping. The compression strategy in Eq. (1)
468 in the present study follows that in Chen *et al.* (2013), and is similar to that developed in Loizou
469 *et al.* (2000). That is, both compression strategies use a linear transformation to convert the range
470 of the input amplitude envelope to a smaller range of the output amplitude envelope; however, the
471 main differences lie in (1) how the minimum envelope amplitude of the input signal is determined
472 and (2) how the linear transformation is designed. In addition, the compression strategy in Eq. (1)
473 preserves the loudness of the processed speech signal.

474 The present work further showed that noise-vocoded speech was more negatively affected by
475 reducing the envelope dynamic range. With the same compressed envelope waveform (e.g.,
476 compression factor of $\alpha = 0.5$ or 0.2 in Fig. 4), noise-vocoded speech showed a much larger drop
477 in intelligibility than did tone-vocoded speech relative to the uncompressed condition (i.e.,
478 compression factor of $\alpha = 1.0$). For instance, at 3 dB SNR level, compared to the uncompressed
479 condition, a 6-dB drop of envelope dynamic range reduced intelligibility by 4.6% and 29.8%, and
480 a 14-dB drop reduced intelligibility by 18.8% and 67.5%, for TV and NV speech, respectively
481 (Fig. 5). This result indicates that narrowing the envelope dynamic range has a more negative
482 influence on noise- than on tone-vocoded speech. This finding may not be fully attributed to the

483 confounding factor of saturation or flooring effect when comparing the intelligibility of TV and
484 NV speech (see Fig. 4). Analysis in Fig. 5 excluded the effect of saturation/flooring in
485 intelligibility scores by choosing two different SNR levels for TV and NT speech (i.e., -3 dB and
486 3 dB, respectively). Again, it is observed in Fig. 5 that NV speech is more susceptible to the
487 influence of reduced envelope dynamic range than TV speech, and its intelligibility score drops at
488 a higher range than TV speech does.

489

490 **D. Implications of vocoder-based acoustic simulation for studies with CIs**

491 Vocoder simulations have been used for inferring systematically effects of noise suppression
492 and dynamic range compression on speech intelligibility for the purpose of implications in CI
493 listeners (e.g., Lai *et al.*, 2015). Researchers have also developed speech-processing strategies for
494 tonal languages (e.g., Mandarin Chinese) and have applied vocoder simulations for assessing their
495 performance (e.g., Luo and Fu, 2006, Lan *et al.*, 2004). Establishing an optimal vocoder for
496 acoustic simulation in CI studies remains an important issue. Although tone-vocoding yields an
497 intelligibility advantage over noise-vocoding, both simulation types may reflect speech
498 intelligibility performance trends with respect to manipulations of acoustic cues. The present
499 study provides evidence of an intelligibility difference between NV and TV speech for NH
500 listeners when an extra signal-processing block is involved in the vocoding process.

501 The better intelligibility of TV sentences relative to NV sentences may be due, at least in part,
502 to the spectral sidebands contained in TV speech and the absence of intrinsic temporal
503 fluctuations in sine wave carriers. Accordingly, when the vocoding process is combined with
504 another signal-processing block, such as noise suppression or envelope dynamic range
505 compression, it is necessary to consider potential interactions between the nature of the carrier
506 signal and distortion produced during signal processing and how such interactions may impede
507 the performance of the signal. The lower intelligibility of the NV speech might be attributable to
508 envelope distortion caused by noise suppression and/or increased envelope distortion when the
509 noise-suppressed envelope is multiplied by a noise carrier containing noise-like amplitude
510 fluctuation. Conversely, the higher intelligibility of TV Mandarin speech may be due in part to a
511 potential contribution of the spectral sidebands in the tone-vocoding process.

512 Although the tone and noise vocoders implemented in this work mimic speech processing in a
513 CI device, many patient- and device-specific confounds were not addressed, including electrode
514 array insertion, spread of the electrical field generated by the implant, etc. Williges *et al.* (2015)
515 used a modified vocoder to sample the envelope waveform in each channel with either sequential
516 or randomized pulse train. Spatial spread of the electrical field was simulated by multiplying each
517 pulse with a two-sided exponential decaying function; additionally, an auralization step was
518 implemented to mimic the transfer of signals in each channel to their respective positions along
519 the cochlea. This vocoder implementation provides a realistic simulation of the technical and
520 physiological steps of signal processing in CI listeners. Future work should investigate the effect
521 of modelling such physiologically-inspired features on the results presented here.

522

523 **E. Limitations of the present work**

524 First, the present work was focused selectively on the effects of noise-suppression and
525 envelope dynamic range compression on the intelligibility of vocoded sentences. Many other
526 factors that may affect the performance of these two vocoder types were not considered, such as
527 envelope cutoff frequency, the number of channels, and filter width. Notably, Rosen *et al.* (2015)
528 showed that a noise vocoder yielded a higher intelligibility than a tone vocoder for a small
529 number of channels (i.e., 2–5). Second, the contribution of the selected cutoff frequency (200 Hz
530 in this study) for extracting the envelope waveform needs to be further investigated. With a
531 200-Hz cutoff frequency, the original signal (envelope, and a portion of full-wave rectified fine
532 structure waveform) is preserved through channel one (with cutoff frequencies of 80 Hz and 221
533 Hz) with the tone vocoder, but not with the noise vocoder which adds noisy fluctuations. Third,
534 the tone vocoder modulates the carrier sinusoids in each frequency channel, i.e., the narrow-band
535 signals. The noise vocoder, however, modulates white noise in each channel and then the
536 amplitude-modulated white noises are bandpass-filtered; or the noise vocoder modulates
537 wideband signals. Hence, it is possible that the aforementioned relative intelligibility deficit of
538 noise-vocoded speech simulations may be due, perhaps in part, to the additional narrowband
539 filtering of the amplitude modulation that occurs at the end of the noise vocoding process. Fourth,
540 the present work used the envelope dynamic range compression strategy developed by Chen *et al.*

541 (2013). It is possible that a different pattern of results would be obtained with the use of
542 alternative compression strategies.

543

544 In conclusion, the present work assessed the effects of noise suppression and envelope
545 dynamic range compression on the intelligibility of vocoded Mandarin sentences, and compared
546 the intelligibility of tone- versus noise-vocoded speech. The following conclusions can be drawn:

547 1) Under all test conditions, tone-vocoded Mandarin sentences showed higher intelligibility
548 scores than did noise-vocoded sentences. This perceptual advantage is consistent with
549 earlier findings. The present study extends this result to vocoded speech that was
550 processed through a noise-suppression algorithm and through envelope dynamic range
551 compression. The perceptual advantage of tone-vocoded Mandarin speech might be
552 attributable to the spectral sidebands contained in tone-vocoded speech and the influence
553 of the amplitude fluctuation of a noise carrier.

554 2) The intelligibility benefit of noise suppression on both tone- and noise-vocoded speech
555 was dependent upon the masker type employed. When corrupted by a steady-state noise,
556 existing single-channel noise-reduction algorithms (e.g., Wiener filtering) might cause
557 intelligibility improvement. However, when corrupted by a competing masker (e.g.,
558 two-talker babble), most existing noise-suppression algorithms did not yield
559 intelligibility improvement.

560 3) While the envelope dynamic range was narrowed, both tone- and noise-vocoded speech
561 showed reduced intelligibility performance. However, noise-vocoded speech was more
562 negatively influenced by envelope dynamic range compression, yielding a substantial
563 intelligibility gap between tone- and noise-vocoded speech.

564 4) When additional signal processing is involved in vocoder simulations, interpreting the
565 functional contribution of this processing should be done cautiously. The nature of the
566 carrier signal in the vocoding process and the envelope distortion caused during signal
567 processing may jointly affect the intelligibility of vocoded speech.

568

569 **ACKNOWLEDGEMENTS**

570 This work was supported by the National Nature Science Foundation of China (Grant No.
571 61571213), and the Basic Research Foundation of Shenzhen (Grant No.
572 JCYJ20160429191402782).

573

574 Baumgärtel, R. M., Krawczyk-Becker, M., Marquardt, D., Völker, C., Hu, H., Herzke, T.,
575 Coleman, G., Adiloğlu, K., Ernst, S. M., Gerkmann, T., Doclo, S., Kollmeier, B., Hohmann, V.,
576 Dietz, M. (2015). “Comparing Binaural Pre-processing Strategies I: Instrumental Evaluation,”
577 Trends Hear. **19**, 1–16.

578 Carhart, R., Tillman, T. W., and Johnson, K. R. (1967). “Release of masking for speech through
579 interaural time delay,” J. Acoust. Soc. Am. **42**, 124–138.

580 Chen, F., Hu, Y., and Yuan, M. (2015). “Evaluation of noise reduction methods for speech
581 recognition by Mandarin-speaking cochlear implant listeners,” Ear Hear. **36**, 61–71.

582 Chen, F., and Lau, A. H. Y. (2014). “Effect of vocoder type to Mandarin speech recognition in
583 cochlear implant simulation,” in *International Symposium on Chinese Spoken Language*
584 *Processing*, pp. 551–554.

585 Chen, F., Wong, L. L., Qiu, J., Liu, Y., Azimi, B., and Hu, Y. (2013). “The contribution of matched
586 envelope dynamic range to the binaural benefits in simulated bilateral electric hearing,” J.
587 Speech Language Hear. Research **56**, 1166–1174.

588 Chen, F., and Loizou, P. C. (2011a). “Predicting the intelligibility of vocoded speech,” Ear Hear. **32**,
589 3281–3290.

590 Chen, F., and Loizou, P. C. (2011b). “Predicting the intelligibility of vocoded and wideband
591 Mandarin Chinese,” J. Acoust. Soc. Am. **129**, 3281–3290.

592 Dorman, M. F., Loizou, P. C., and Rainey, D. (1997). “Speech intelligibility as a function of the
593 number of channels of stimulation for signal processors using sine-wave and noise-band
594 outputs,” J. Acoust. Soc. Am. **102**, 2403–2411.

595 Ephraim, Y., and Malah, D. (1985). “Speech enhancement using a minimum mean-square error
596 log-spectral amplitude estimator,” IEEE Trans. Acoust., Speech, Signal Process. ASSP-**33**,
597 443–445.

598 Fu, Q. J., Chinchilla, S., and Galvin, J. J. (2004). “The role of spectral and temporal cues in voice
599 gender discrimination by normal-hearing listeners and cochlear implant users,” J. Asso.

600 Research Oto. **5**, 253–260.

601 Fu, Q. J., and Shannon, R. V. (1999). “Effect of acoustic dynamic range on phoneme recognition in
602 quiet and noise by cochlear implant users,” J. Acoust. Soc. Am. **106**, EL65–EL70.

603 Gonzalez, J., and Oliver, J. C. (2005). “Gender and speaker identification as a function of the
604 number of channels in spectrally reduced speech,” J. Acoust. Soc. Am. **118**, 461–470.

605 Howie, J.M. (1976). *Acoustical studies of Mandarin Vowels and Tones* (Cambridge University
606 Press, Cambridge, England), pp. 1–308.

607 Hu, Y., and Loizou, P. (2003). “A generalized subspace approach for enhancing speech corrupted by
608 colored noise,” IEEE Trans. Speech Audio Process. **11**, 334–341.

609 Hu, Y., and Loizou, P. C. (2007). “A comparative intelligibility study of single-microphone noise
610 reduction algorithms,” J. Acoust. Soc. Am. **122**, 1777–1786.

611 Kamath, S., and Loizou, P. (2002). “A multi-band spectral subtraction method for enhancing speech
612 corrupted by colored noise,” in *Proceedings of the IEEE International Conference on Acoustics,
613 Speech, and Signal Processing*, pp. IV–4164.

614 Kasturi, K., and Loizou, P. C. (2007). “Effect of filter spacing on melody recognition: acoustic and
615 electric hearing,” J. Acoust. Soc. Am. **122**, 29–34.

616 Lai, Y. H., Tsao, Y., and Chen, F. (2015). “Effects of adaptation rate and noise suppression on the
617 intelligibility of compressed-envelope based speech,” PLOS ONE **10**, e133519.

618 Lan, N., Nie, K., Gao, S., and Zeng, F. G. (2004). “A novel speech-processing strategy
619 incorporating tonal information for cochlear implants,” IEEE Trans. Biomed. Eng. **52**,
620 752–760.

621 Li, J., Yang, L., Zhang, J., Yan, Y., Hu, Y., Akagi, M., and Loizou, P. C. (2011). “Comparative
622 intelligibility investigation of single-channel noise-reduction algorithms for Chinese, Japanese,
623 and English,” J. Acoust. Soc. Am. **129**, 3291–3301.

624 Loizou, P. C. (2007). *Speech Enhancement: Theory and Practice* (CRC, Boca Raton, FL), pp.
625 1–689.

626 Loizou, P. C., and Kim, G. (2011). “Reasons why current speech-enhancement algorithms do not
627 improve speech intelligibility and suggested solutions,” IEEE Trans. Audio Speech Lang
628 Process. **19**, 47–56.

629 Loizou, P. C., Dorman, M., and Fitzke, J. (2000). “The effect of reduced dynamic range on speech

630 understanding: implications for patients with cochlear implants,” *Ear Hear.* **21**, 25–31.

631 Luo, X., and Fu, Q. J (2006). “Contribution of low-frequency acoustic information to Chinese
632 speech recognition in cochlear implant simulations,” *J. Acoust. Soc. Am.* **120**, 2260–2266.

633 Luo, X., and Fu, Q. J (2004). “Enhancing Chinese tone recognition by manipulating amplitude
634 envelope: Implications for cochlear implants,” *J. Acoust. Soc. Am.* **116**, 3659–3667.

635 Rosen, S., Zhang, Y., and Speers, K. (2015). “Spectral density affects the intelligibility of
636 tone-vocoded speech: Implications for cochlear implant simulations,” *J. Acoust. Soc. Am.* **138**,
637 EL318–EL323.

638 Scalart, P., and Filho, J. (1996). “Speech enhancement based on a priori signal to noise estimation,”
639 in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal*
640 *Processing*, pp. 629–632.

641 Shannon, R. V., Zeng, F. G., Kamath, V., Wygonski, J., and Ekelid, M. (1995). “Speech recognition
642 with primarily temporal cues,” *Science* **270**, 303–304.

643 Souza, P., and Rosen, S. (2009). “Effects of envelope bandwidth on the intelligibility of sine- and
644 noise-vocoded speech,” *J. Acoust. Soc. Am.* **126**, 792–805.

645 Stone, M. A., Füllgrabe, C., Mackinnon, R. C., and Moore, B. C. (2011). “The importance for
646 speech intelligibility of random fluctuations in ‘steady’ background noise,” *J. Acoust. Soc. Am.*
647 **130**, 2874–2881.

648 Stone, M. A., Füllgrabe, C., Moore, B. C. (2008). “Benefit of high-rate envelope cues in vocoder
649 processing: Effect of number of channels and spectral region,” *J. Acoust. Soc. Am.* **124**,
650 2272–2282.

651 Studebaker, G. A. (1985). “A ‘rationalized’ arcsine transform,” *J. Speech Hear Research* **28**,
652 455–462.

653 Watson, C. S. (2005). “Some comments on informational masking,” *Acta Acust* **91**, 502–512.

654 Whitmal, N. A., Poissant, S. F., Freyman, R. L., and Helfer, K. S. (2007). “Speech intelligibility in
655 cochlear implant simulations: Effects of carrier type, interfering noise, and subject experience,”
656 *J. Acoust. Soc. Am.* **122**, 2376–2388.

657 Williges, B., Dietz, M., Hohmann, V., and Jürgens, T. (2015). “Spatial release from masking in
658 simulated cochlear implant users with and without access to low-frequency acoustic hearing,”
659 *Trends Hear.* **19**, 1–14.

- 660 Wong, L. L., Soli, S. D., Liu, S., Han, N., and Huang, M. W. (2007). "Development of the Mandarin
661 Hearing in Noise Test (MHINT)," *Ear Hear.* **28**, 70S–74S.
- 662 Xu, L., Thompson, C. S., and Pfingst, B. E. (2005). "Relative contributions of spectral and temporal
663 cues for phoneme recognition," *J. Acoust. Soc. Am.* **117**, 3255–3267.
- 664 Zeng, F. G., Grant, G., Niparko, J., Galvin, J., Shannon, R., Opie, J., and Segel, P. (2002). "Speech
665 dynamic range and its effect on cochlear implant performance," *J. Acoust. Soc. Am.* **111**,
666 377–386.