

1 **Publication details:** *Proceedings of the Royal Society of London, Series B: Biological*
2 *Sciences, 2011, 278, 1595-1600.*

3
4
5
6
7
8
9

10 **The intelligibility of noise-vocoded speech: Spectral information**
11 **available from across-channel comparison of amplitude envelopes**

12
13

14 Brian Roberts^{1*}, Robert J. Summers¹, and Peter J. Bailey²

15

16 1. Psychology, School of Life and Health Sciences, Aston University, Birmingham B4
17 7ET, UK.

18

19 2. Department of Psychology, University of York, Heslington, York YO10 5DD, UK.

20

21

22

23

24

25 Running title: Spectral cues in noise-vocoded speech

26

27

28

29

30

31 * **Correspondence to:** Professor B. Roberts at Psychology, School of Life and Health
32 Sciences, Aston University, Birmingham, B4 7ET, UK

33 Email: b.roberts@aston.ac.uk

34 **Summary**

35

36 Noise-vocoded (NV) speech is often regarded as conveying phonetic information
37 primarily through temporal-envelope cues rather than spectral cues. However, listeners
38 may infer the formant frequencies in the vocal-tract output – a key source of phonetic
39 detail – from across-band differences in amplitude when speech is processed through a
40 small number of channels. The potential utility of this spectral information was assessed
41 for NV speech created by filtering sentences into six frequency bands, and using the
42 amplitude envelope of each band (≤ 30 Hz) to modulate a matched noise-band carrier
43 (N). Bands were paired, corresponding to F1 ($\approx N1+N2$), F2 ($\approx N3+N4$), and the higher
44 formants (F3' $\approx N5+N6$), such that the frequency contour of each formant was implied
45 by variations in relative amplitude between bands within the corresponding pair.
46 Three-formant analogues (F0=150 Hz) of the NV stimuli were synthesised using
47 frame-by-frame reconstruction of the frequency and amplitude of each formant. These
48 analogues were less intelligible than the NV stimuli or analogues created using contours
49 extracted from spectrograms of the original sentences, but more intelligible than when
50 the frequency contours were replaced with constant (mean) values. Across-band
51 comparisons of amplitude envelopes in NV speech can provide phonetically important
52 information about the frequency contours of the underlying formants.

53

54

55

56 **Keywords:** noise-vocoded speech, spectral cues, formant frequencies, intelligibility.

57

58

59 1. INTRODUCTION

60 Speech is highly redundant and so it can remain intelligible even after substantial
61 distortion or simplification of the signal. A commonly used simplification is vocoding,
62 which involves filtering speech into one or more frequency bands, using the amplitude
63 envelope of each band to modulate a carrier shaped by the corresponding filter, and
64 reconstructing the simplified signal by summing the modulated carrier bands. The
65 technique was originally devised by Dudley (1939) for speech transmission through
66 telecommunications systems, particularly for encrypted communications, and has since
67 been used widely for voice processing in popular music. Shannon et al. (1995) first
68 introduced noise vocoding, in which the carrier for each channel is filtered Gaussian
69 noise. Their study demonstrated that the intelligibility of noise-vocoded speech can be
70 high when only three or four channels are used, at least when all stimuli are derived
71 from the speech of a single talker. Dorman et al. (1997) obtained comparable results
72 with sine-vocoded speech, a closely related stimulus consisting of a set of amplitude
73 modulated sinusoids instead of noise bands. Vocoding has since become a standard
74 research tool for simulating listening to speech through a cochlear implant; many
75 contemporary studies use noise-vocoded speech (e.g., Li and Loizou 2009; Loebach et
76 al. 2009; Chatterjee et al. 2010; Eisner et al. 2010) or sine-vocoded speech (e.g., Chen
77 and Loizou 2010; Hopkins and Moore 2010) for this purpose.

78 Interpreting the results of perceptual experiments using vocoded speech requires
79 an understanding of the nature of, and weight attached to, sources of phonetic
80 information in the signal. Processing speech through a noise vocoder with a small
81 number of channels implies a considerable loss of spectral information. Hence, this type
82 of stimulus is often regarded as conveying phonetic information primarily through
83 temporal-envelope cues rather than spectral cues (Shannon et al. 1995; Nitttrouer et al.
84 2009). In this conception, the intelligibility of noise-vocoded speech depends mainly on
85 the within-channel analysis of low-rate changes in amplitude over time; an account of
86 the types of linguistic information potentially available from temporal-envelope cues
87 has been provided by Rosen (1992). Other studies of noise- and sine-vocoded speech
88 have tended to characterise the relative contributions of spectral and temporal cues to
89 intelligibility in terms of the effects of varying the number of channels and the low-pass
90 envelope cut-off, respectively, and their trade-off (e.g., Xu et al. 2005; Xu and Zheng
91 2007; Xu and Pfingst 2008). What is often overlooked in this characterisation is the

92 spectral information that can potentially be retrieved through comparing the levels of
93 adjacent channels (Dorman et al. 1997; Loizou et al. 1998). In particular, changes in
94 relative amplitude across channels over time can potentially carry information about the
95 underlying frequency contours of the spectral prominences in the signal, and this
96 derived spectral information may contribute more (and temporal-envelope information
97 perhaps less) to the intelligibility of noise- and sine-vocoded speech than is commonly
98 supposed.

99 Spectral prominences in speech – called *formants* – are perceptually important,
100 because they arise as a result of resonances in the air-filled cavities of the talker’s vocal
101 tract. Variation in the centre frequency of a formant is an inevitable consequence of
102 change in the size of its associated cavity as the vocal-tract articulators – particularly the
103 tongue, lips, and jaw – are moved by the talker. Thus, knowledge of formant
104 frequencies and their change over time is likely to be of considerable benefit to listeners,
105 as it provides salient information about the configuration and kinematics of the talker’s
106 vocal tract. The experiment reported here demonstrates that the formant-frequency
107 contours implied by variations in relative amplitude between adjacent spectral bands
108 can be extracted from noise-vocoded signals and can support intelligibility in
109 synthetic-formant analogues of speech.

110

111 **2. METHODS**

112 *(a) Participants*

113 Twenty listeners (10 males) took part; their mean age was 23.2 years (range = 19.2 –
114 54.7). All listeners were native speakers of British English, naïve to the purpose of the
115 experiment, and had audiometric thresholds better than 20 dB hearing level at 0.5, 1, 2,
116 and 4 kHz. Each listener gave written consent to participate in the experiment, which
117 was approved by the Aston University Ethics Committee.

118

119 *(b) Stimuli and conditions*

120 All stimuli were derived from 24 BKB sentences (Bench et al. 1979), spoken by a
121 British male talker of Received Pronunciation English and low-pass filtered at 5 kHz.
122 There were four conditions in the experiment, corresponding to the four speech
123 analogues described below. Figure 1 shows the spectrogram of an example sentence and
124 of the four analogues derived from it. For each listener, the sentences were divided

125 equally across conditions (i.e., six per condition) using an allocation that was
126 counterbalanced by rotation across each set of four listeners tested. Each sentence group
127 was balanced so as to contain 95 or 96 phonemes in total. Examples of the stimuli are
128 available in the electronic supplementary material.

129 -----
130 Figure 1 near here
131 -----

132 *Noise-vocoded (NV) stimuli* were created from the original sentences using Praat
133 software (Boersma & Weenink 2008). The speech was first filtered, using a 16th-order
134 Butterworth filter (96 dB/octave roll-off), into six logarithmically spaced bands with
135 cut-off frequencies of 200, 362, 655, 1186, 2147, 3885, and 7032 Hz. Pairs of bands
136 were tailored to correspond quite closely with the formant ranges of the talker (B1+B2
137 \approx F1; B3+B4 \approx F2; B5+B6 \approx F3 and above, denoted F3'). The amplitude envelope (\leq 30
138 Hz) of each band was then extracted by half-wave rectification and used to modulate a
139 Gaussian noise source with the same lower and upper cut-off frequencies; increasing the
140 low-pass corner frequency above 30 Hz does not further improve the intelligibility of
141 NV speech (Souza and Rosen 2009). Each band (N1-N6) was scaled to have the same
142 RMS level as that of the corresponding band in the original speech and the bands were
143 summed to create the modulated noise-band speech analogues.

144 *Extracted-formant (EF) stimuli* were created from the original sentences using
145 Praat to estimate automatically from the waveform the frequency contours of the first
146 three formants every 1 ms; a 25-ms-long Gaussian window was used. During phonetic
147 segments with frication the third-formant contour often corresponded to the fricative
148 formant rather than to F3. Gross errors in formant-frequency estimates were
149 hand-corrected using a graphics tablet; amplitude contours corresponding to the
150 corrected frequencies were extracted from spectrograms for each sentence. The
151 frequency and amplitude contours were used to generate three-formant analogues of the
152 sentences by means of simple parallel-formant synthesis, using second-order resonators
153 and an excitation pulse modelled on the glottal waveform (Rosenberg 1971). The pitch
154 was monotonous (F0 = 150 Hz), and the 3-dB bandwidths of F1, F2, and F3 were 50, 70,
155 and 90 Hz, respectively.

156
157

158

159

Figure 2 near here

160

161

162

163

164

165

166

167

168

169

170

171

Reconstructed-formant (RF) stimuli were created from the NV sentences using a simple procedure designed to retrieve the information about formant frequency and amplitude carried by each pair of bands (i.e., N1+N2 for F1, N3+N4 for F2, N5+N6 for F3'). For each pair, the amplitude contour of the reconstructed formant was computed frame-by-frame as the mean amplitude across both bands. The frequency contour was derived from frame-by-frame changes in the relative amplitudes of the two bands within each pair. Figure 2 depicts the reconstruction of the F2 frequency contour from the band pair N3+N4 for an example sentence. The reconstructed contours were used to generate three-formant analogues of the sentences by parallel synthesis, as described above for the EF stimuli. At a particular time, the implied frequency, F , is given by:

172

$$\log F = \log(g) + kw \log\left(\frac{f_{hi}}{g}\right), \quad (1)$$

173

$$w = \frac{a_{hi} - a_{lo}}{a_{hi} + a_{lo}} \quad (-1 \leq w \leq +1), \quad (2)$$

174

175

176

177

178

179

180

181

182

183

184

where a_{lo} and a_{hi} are the amplitudes of the lower and upper bands, f_{hi} is the upper cut-off frequency of the upper band, k ($0 < k \leq 1$) is a scale factor determining the maximum possible frequency range, and g is the geometric mean frequency of the lower and upper bands. The value of k used here was 0.9; this was to ensure that the frequency range available for formant excursions in the reconstructions was substantial, but not so great as to have allowed unnaturally close approaches between neighbouring formants. Note that low-pass filtering the original sentences at 5 kHz lowers the amplitude of band N6 in the NV stimuli, which tends to lower the frequency, as well as the amplitude, of the reconstructed F3', particularly during fricative segments. This improves the overall quality of the RF stimuli by reducing the "buzziness" of these segments.

185

186

187

188

Constant-formant (CF) stimuli differed from their RF counterparts only in that the frequency of each formant was set to be constant at the geometric mean frequency of the whole reconstructed track. For all conditions, the speech analogues were played at a sample rate of 22.05 kHz and 16-bit resolution over Sennheiser HD 480-13II

189 earphones, via a sound card, programmable attenuators (Tucker-Davis Technologies
190 PA5), and a headphone buffer (TDT HB7). Output levels were calibrated using a
191 sound-level meter (Brüel and Kjaer, type 2209) coupled to the earphones by an artificial
192 ear (type 4153). All stimuli were shaped using 10-ms raised-cosine onset and offset
193 ramps and presented diotically at 75 dB SPL.

194

195 *(c) Procedure*

196 Listeners were tested whilst seated in front of a computer screen and a keyboard
197 in a sound-attenuating booth. There were two phases to the study, training and the main
198 experiment, which together took less than an hour to complete. Stimuli were presented
199 in quasi-random order in both phases of the study. Listeners first completed a training
200 session to familiarise them with synthetic-formant and noise-vocoded speech analogues,
201 in that order. The former were examples of EF stimuli, but differed from those used in
202 the main experiment in that the natural pitch contour was used in the resynthesis;
203 listeners were not exposed to RF or CF stimuli during training. The stimuli for each part
204 of the training were derived from 40 sentences taken from commercially available
205 recordings of the IEEE sentence lists (IEEE, 1969). On each of the 40 trials in each part,
206 participants were able to listen to the stimulus up to a maximum of six times before
207 typing in their transcription of the sentence. After each transcription was entered,
208 feedback to the listener was provided by playing the original recording followed by a
209 repeat of the speech analogue. Davis et al. (2005) found this “degraded-clear-degraded”
210 presentation strategy to be an efficient way of enhancing the perceptual learning of
211 speech analogues. All listeners who obtained scores of $\geq 60\%$ keywords correct in the
212 second half of each set of training trials were included in the main experiment. As in the
213 training, participants in the main experiment were able to listen to each stimulus up to
214 six times before typing in their transcription, and the time available to respond was not
215 limited. However, this time the listeners did not receive feedback of any kind on their
216 responses.

217

218 *(d) Data analysis*

219 For each listener, the intelligibility of each sentence was quantified in terms of the
220 overall percentage of phonetic segments identified correctly. Phonetic scores are usually
221 more effective at distinguishing performance between conditions for which there is

222 limited intelligibility, owing to floor effects in keyword scores. Listeners' typed
223 responses were converted automatically into phonetic representations using eSpeak
224 (Duddington 2008) for comparison with stored phonetic representations of the original
225 sentences. Phonetic scores were computed using HResults, part of the HTK software
226 (Young et al. 2006). HResults uses a string alignment algorithm to find an optimal
227 match between two strings.

228

229 **3. RESULTS AND DISCUSSION**

230 Figure 3 shows the mean percentage of phonetic segments identified correctly across
231 conditions, with inter-subject standard errors. A within-subjects analysis of variance
232 (ANOVA) showed a highly significant effect of condition on intelligibility [$F(3,57)=$
233 278.9 , $p<0.001$, $\eta^2=0.936$]. Paired-samples comparisons (two-tailed) were computed
234 using the restricted least-significant-difference test (Snedecor and Cochran 1967); the
235 scores for each condition differed significantly from those for every other condition
236 ($p<0.001$, in all cases). Scores were very high for the NV speech, given that there were
237 only six spectral bands (cf. Shannon et al., 1995). This may reflect the tailored
238 alignment of each pair of bands in relation to the talker's ranges of formant frequencies.
239 More generally, the intelligibility of NV speech tends to be lower when the inventory of
240 stimuli is derived from multiple talkers (Loizou et al. 1999); this reflects the need for
241 listeners to accommodate acoustic-phonetic variability across talkers (see, e.g.,
242 Mullennix et al., 1989). Scores were somewhat lower for the EF speech, probably as a
243 result of two sources of error in recreating phonetically relevant acoustic detail. First,
244 estimation of formant frequencies from fluent speech is a technical challenge and prone
245 to inaccuracy, even when the output of an algorithm for the automatic extraction of
246 formant frequencies is subject to hand correction. Such errors in the formant-frequency
247 parameters fed to the synthesiser would be expected to impair intelligibility (Assmann
248 and Katz, 2005). Second, the use of a minimal model for the formant synthesiser,
249 incorporating only three fixed-bandwidth formants, will have introduced synthesis
250 errors, notably in the reproduction of phonetic segments having significant amounts of
251 high-frequency energy, such as voiceless fricatives. These sources of error do not
252 contribute to the process of creating NV speech.

253

254

255

256

257

Figure 3 near here

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

Scores approached 30% when the three-formant analogues were created using frequency and amplitude contours that were reconstructed from the amplitude-envelope information carried by the three band pairs comprising the NV analogues of the original speech. Hence, RF speech was still nearly half as intelligible as EF speech, even using such a simplistic approach to reconstructing the formant-frequency contours from the NV speech. The frequency resolution of normal-hearing listeners far exceeds that required to retrieve this information from the representation of noise-vocoded speech in the peripheral auditory system (see, e.g., Moore 2003). Performance was halved again for CF speech, for which the reconstructed frequency contours were replaced with constant values set to the geometric mean of each formant track. At least in part, the non-zero performance for the CF speech might be because the reconstructed amplitude contours still convey useful information about vocal tract dynamics. Note, however, that simulations comparing randomly generated text strings with those specifying the stimuli used here suggest that baseline phonetic scores can be in the region of 15% for entirely unintelligible speech; the mean score for the CF stimuli was 12%. Remez and Rubin (1990) explored the relative contributions of variations in the frequency and amplitude contours of formants to the intelligibility of sine-wave speech, created by adding together pure tones that follow the frequency and amplitude contours of the lower formants (Bailey et al 1977; Remez et al. 1981). They concluded that frequency variation is far more important than amplitude variation for maintaining intelligibility; this is also true for across-formant grouping in sine-wave speech (Roberts et al. 2010).

The higher recognition scores observed for the RF relative to the CF condition support the notion that changes in the relative amplitudes of different bands in NV speech convey useful phonetic information about formant frequency variation. Consistent with this view, the effect of quantising the amplitude envelope into a small number of steps (<8) has a much greater impact on the intelligibility of sine-vocoded speech processed through a small (6) rather than a large (16) number of channels, presumably because the reduced information available from across-channel amplitude comparisons makes it more difficult to infer the underlying formant frequencies (Loizou et al 1999). The importance of combining information across a small number of

288 channels to reconstruct signal properties important for intelligibility is also evident in
289 Apoux and Healy's (2009) demonstration that phonemes can be identified from
290 relatively few randomly selected channels, even when noise is present in other channels.
291 Dorman et al. (1997) suggested that the mechanism mediating the high degree of
292 intelligibility achievable with a small number of channels may be the same as that
293 mediating the recognition of speech produced by talkers with a high fundamental
294 frequency.

295 Recently, more direct evidence that the frequency contours of formants can be
296 inferred from across-channel amplitude comparisons, at least for single formant
297 transitions, has been provided by Fox et al. (2008). Their study explored the role of F3
298 transitions in distinguishing the place of articulation of initial stops in the syllable pairs
299 [da]-[ga] and [ta]-[ka]. They compared actual F3 transitions with virtual ones, where the
300 percept of a frequency transition was cued by a dynamic change in spectral centre of
301 gravity over 50 ms arising from a smooth but rapid change in the relative amplitude of
302 two noise-excited formants with constant frequency (1907 Hz and 2861 Hz). These
303 frequencies are easily resolvable by the peripheral auditory system, but fall within the
304 much larger bandwidth of about 3.5 critical bands (roughly 5 equivalent rectangular
305 bandwidths; Glasberg and Moore 1990) over which the central auditory system appears
306 to integrate phonetic information (e.g., Delattre et al. 1952; Carlson et al. 1975;
307 Chistovich 1985). Virtual F3 transitions were broadly comparable with actual F3
308 transitions in supporting the correct identification of initial stops; listeners could also
309 distinguish the direction of the F3 transitions when heard in isolation as rising or falling,
310 whether actual or virtual. Fox et al. (2008) concluded that amplitude and frequency
311 information can be combined in the perception of formant transitions.

312 To conclude, across-band comparisons of amplitude envelopes in NV speech can
313 provide phonetically important information about the implied frequency contours of the
314 underlying formants for sentence-length utterances. In principle, this dynamic spectral
315 information is easily accessible to most listeners even when the number of channels
316 available is relatively limited.

317

318 **Acknowledgements:** This work was supported by EPSRC (UK). Grant Reference
319 EP/F016484/1 (Roberts & Bailey).

320 **REFERENCES**

- 321 Apoux, F., Healy, E.F. 2009 On the number of auditory filter outputs needed to
322 understand speech: Further evidence for auditory channel independence. *Hear.*
323 *Res.* **255**, 99-108.
- 324 Assmann, P.F., Katz, W.F. 2005 Synthesis fidelity and time-varying spectral change in
325 vowels. *J Acoust Soc Am.* **117**, 886-895.
- 326 Bailey, P.J., Summerfield, Q., Dorman, M. 1977 On the identification of sine-wave
327 analogues of certain speech sounds. *Haskins Lab. Status Rep.* **SR-51/52**, 1-25.
- 328 Bench, J., Kowal, A., Bamford, J. 1979 The BKB (Bamford-Kowal-Bench) sentence
329 lists for partially-hearing children. *Brit. J. Audiol.* **13**, 108-112.
- 330 Boersma, P., Weenink, D. 2008 Praat: Doing phonetics by computer [software package],
331 version 5.0.18, retrieved 1 April 2008 from <http://www.praat.org/>
- 332 Carlson, R., Fant, G., Granstrom, B. 1975 Two-formant models, pitch and vowel
333 perception. In *Auditory Analysis and Perception of Speech* (eds. G. Fant, M.A.A.
334 Tatham), pp. 55-82. London: Academic Press.
- 335 Chatterjee, M., Peredo, F., Nelson, D., Başkent, D. 2010 Recognition of interrupted
336 sentences under conditions of spectral degradation. *J. Acoust. Soc. Am.* **127**,
337 EL37-EL41.
- 338 Chen, F., Loizou, P.C. 2010 Contribution of consonant landmarks to speech recognition
339 in simulated acoustic-electric hearing. *Ear Hear.* **31**, 259-267.
- 340 Chistovich, L.A. 1985 Central auditory processing of peripheral vowel spectra. *J.*
341 *Acoust. Soc. Am.* **77**, 789-805.
- 342 Davis, M.H., Johnsrude, I.S., Hervais-Adelman, A., Taylor, K., McGettigan, C. 2005
343 Lexical information drives perceptual learning of distorted speech: evidence from
344 the comprehension of noise-vocoded sentences. *J. Exp. Psychol. Gen.* **134**,
345 222-241.
- 346 Delattre, P.C., Liberman, A.M., Cooper, F.S., Gerstman, L.J. 1952 An experimental
347 study of the acoustic determinants of vowel color; observations on one- and
348 two-formant vowels synthesized from spectrographic patterns. *Word* **8**, 195-210.
- 349 Dorman, M., Loizou, P., Rainey, D. 1997 Speech intelligibility as a function of the
350 number of channels of stimulation for signal processors using sine-wave and
351 noise-band outputs. *J. Acoust. Soc. Am.* **102**, 2403-2411.
- 352 Duddington, J. 2008 eSpeak 1.36, <http://espeak.sourceforge.net/>

- 353 Dudley, H. 1939 Remaking speech. *J. Acoust. Soc. Am.* **11**, 169-177.
- 354 Eisner, F., McGettigan, C., Faulkner, A., Rosen, S., Scott, S.K. 2010 Inferior frontal
355 gyrus activation predicts individual differences in perceptual learning of
356 cochlear-implant simulations. *J. Neurosci.* **30**, 7179-7186.
- 357 Fox, R.A., Jacewicz, E., Feth, L.L. 2008 Spectral integration of dynamic cues in the
358 perception of syllable-initial stops. *Phonetica* **65**, 19-44.
- 359 Glasberg, B.R., Moore, B.C.J. 1990 Derivation of auditory filter shapes from
360 notched-noise data. *Hear. Res.* **47**, 103-138.
- 361 Hopkins, K., Moore, B.C.J. 2010 The importance of temporal fine structure information
362 in speech at different spectral regions for normal-hearing and hearing-impaired
363 subjects. *J. Acoust. Soc. Am.* **127**, 1595-1608.
- 364 Institute of Electrical and Electronics Engineers (IEEE) 1969 IEEE recommended
365 practice for speech quality measurements. *IEEE Trans. Audio Electroacoust.*
366 **AU-17**, 225-246.
- 367 Li, N., & Loizou, P. C. 2009 Factors affecting masking release in cochlear-implant
368 vocoded speech. *J. Acoust. Soc. Am.* **126**, 338-346.
- 369 Loebach, J.L., Pisoni, D.B., Svirsky, M.A. 2009 Transfer of auditory perceptual learning
370 with spectrally reduced speech to speech and nonspeech tasks: Implications for
371 cochlear implants. *Ear Hear.* **30**, 662-674.
- 372 Loizou, P., Dorman, M., Powell, V. 1998 The recognition of vowels produced by men,
373 women, boys and girls by cochlear implant patients using a six-channel CIS
374 processor. *J. Acoust. Soc. Am.* **103**, 1141-1149.
- 375 Loizou, P.C., Dorman, M., Tu, Z. 1999 On the number of channels needed to understand
376 speech. *J. Acoust. Soc. Am.* **106**, 2097-2103.
- 377 Mullennix, J., Pisoni, D., Martin, C. 1989 Some effects of talker variability on spoken
378 word recognition. *J. Acoust. Soc. Am.* **85**, 365-378.
- 379 Moore, B.C.J. 2003 *An Introduction to the Psychology of Hearing*, 5th ed. London:
380 Academic Press.
- 381 Nittrouer, S., Lowenstein, J.H., Packer, R. 2009 Children discover the spectral skeletons
382 in their native language before the amplitude envelopes. *J. Exp. Psychol. Hum.*
383 *Percept. Perform.* **35**, 1245-1253.
- 384 Remez, R.E., Rubin, P.E. 1990 On the perception of speech from time-varying acoustic
385 information: Contributions of amplitude variation. *Percept. Psychophys.* **48**,

- 386 313-325.
- 387 Remez, R.E., Rubin, P.E., Pisoni, D.B., Carrell, T.D. 1981 Speech perception without
388 traditional speech cues. *Science* **212**, 947-950.
- 389 Roberts, B., Summers, R.J., Bailey, P.J. 2010 The perceptual organization of sine-wave
390 speech under competitive conditions. *J. Acoust. Soc. Am.* **128**, 804-817.
- 391 Rosen, S. 1992 Temporal information in speech: acoustic, auditory and linguistic
392 aspects. *Phil. Trans. R. Soc. Lond. B* **336**, 367-373.
- 393 Rosenberg, A.E. 1971 Effect of glottal pulse shape on the quality of natural vowels. *J.*
394 *Acoust. Soc. Am.* **49**, 583-590.
- 395 Shannon, R.V., Zeng, F.-G., Kamath, V., Wygonski, J., Ekelid, M. 1995 Speech
396 recognition with primarily temporal cues. *Science* **270**, 303-304.
- 397 Snedecor, G.W., Cochran, W.G. 1967 *Statistical Methods*, 6th ed. Ames, Iowa: Iowa
398 University Press.
- 399 Souza, P., Rosen, S. 2009 Effects of envelope bandwidth on the intelligibility of sine-
400 and noise-vocoded speech. *J. Acoust. Soc. Am.* **126**, 792-805.
- 401 Xu, L., Pfingst, B.E. 2008 Spectral and temporal cues for speech recognition:
402 Implications for auditory prostheses. *Hear. Res.* **242**, 132-140.
- 403 Xu, L., Thompson, C.S., Pfingst, B.E. 2005 Relative contributions of spectral and
404 temporal cues for phoneme recognition. *J. Acoust. Soc. Am.* **117**, 3255-3267.
- 405 Xu, L., Zheng, Y. 2007 Spectral and temporal cues for phoneme recognition in noise. *J.*
406 *Acoust. Soc. Am.* **122**, 1758-1764.
- 407 Young, S.J., Evermann, G., Gales, M.J.F., Kershaw, D., Moore, G., Odell, J.J., Ollason,
408 D.G., Povey, D., Valtchev, V., Woodland, P.C. 2006 *The HTK book, version 3.4*
409 *Manual*. Cambridge, UK: Department of Engineering, University of Cambridge.

410 **Figure Captions**

411

412 **Figure 1** Spectrograms of an exemplar original sentence, “The oven door was open”,
413 and of the four experimental versions derived from it. The horizontal dashed lines in the
414 panel depicting the noise-vocoded (NV) stimulus indicate the band cut-off frequencies.
415 Note that the most striking discrepancy between the extracted-formants (EF) and the
416 reconstructed-formants (RF) stimuli corresponds to the voiced fricative [z] in “was”. In
417 the EF case, the formant contour extracted by Praat corresponds to F3, but in the RF
418 case the reconstructed formant contour is dominated by the energy in the fricative
419 formant.

420

421

422 **Figure 2** Reconstruction of formant-frequency contours. This schematic illustrates the
423 reconstruction of the frequency contour of F2 from the noise-vocoded (NV) version of
424 the exemplar sentence “The oven door was open”. The reconstructed contour (dashed
425 line) is governed by changes over time in the relative amplitudes of noise bands 3 and 4;
426 the amplitude modulation of each band is depicted by a filled contour centred on the
427 geometric mean frequency. The frequency contour was computed frame by frame using
428 equations 1 and 2 (see main text).

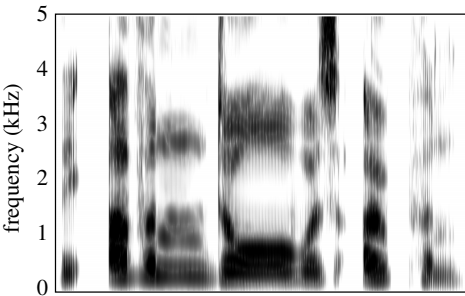
429

430

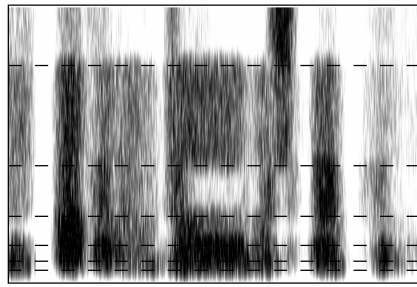
431 **Figure 3** Intelligibility of the four analogues derived from the original sentences. These
432 correspond to the noise-vocoded (NV, 81.6%), extracted-formants (EF, 64.4%),
433 reconstructed-formants (RF, 28.1%), and constant-formants (CF, 12.0%) conditions.
434 Each histogram bar shows the mean phonetic score and corresponding inter-subject
435 standard error (n=20).

436

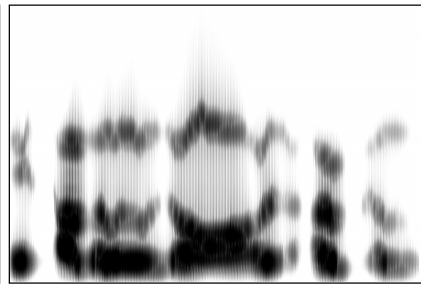
natural



noise vocoded (NV)

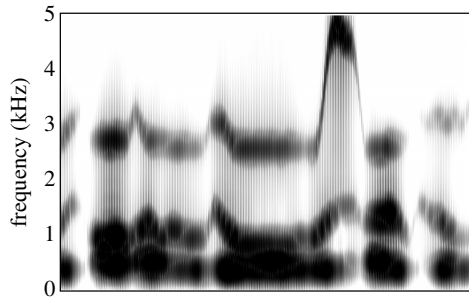


extracted formants (EF)

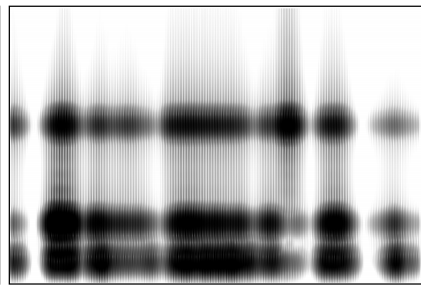


time

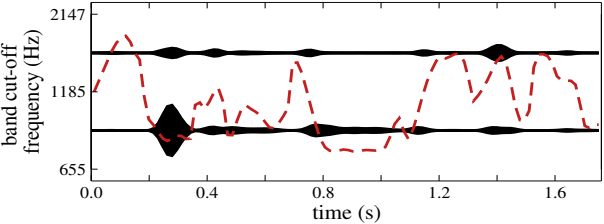
reconstructed formants (RF)



constant formants (CF)



time



noise-vocoded and three-formant speech analogues
diotic presentation (n=20)

