1

2

3        Perception of degraded speech by chinchillas (*Chinchilla laniger*):

4                    Word-level stimulus generalization

5

6    *Authors*: William P. Shofner, Nicole Yacko and Kristina Bowdrie

7

8    Department of Speech and Hearing Sciences

9    Indiana University

10   200 S. Jordan Ave.

11   Bloomington, IN 47405

12

13

15

16   Running Head: Speech perception by chinchillas

17

22

23

24  Correspondence regarding the article should be addressed to:

25   William P. Shofner

26  Department of Speech and Hearing Sciences

27  Indiana University

28  200 S. Jordan Ave.

29  Bloomington, IN 47405

30  Email: wshofner@indiana.edu

31

**Abstract**

One characteristic of human speech perception is a remarkable ability to recognize speech when

the speech signal is highly degraded.  It has been argued that this ability to perceive highly

degraded speech reflects speech-specific mechanisms.  The present study tested this hypothesis

by measuring the ability of chinchillas to recognize noise-vocoded versions of naturally-spoken

monosyllabic words using operant conditioning in a stimulus generalization paradigm.

Chinchillas do not generalize the vocoded words to be perceptually equivalent to the naturally-

spoken words.  The responses from chinchillas to the vocoded words fall well below their

responses to the naturally-spoken words.  In this case, pitch cues rather than speech cues may be

controlling the behavioral responses.  In order to reduce pitch cues, chinchillas were re-trained

using 64-channel noise-vocoded words. The responses from chinchillas to the vocoded test

words were now similar to those of the 64-channel versions and were similar to those obtained

from human listeners.  However, responses obtained from chinchillas to time-reversed versions

were high and similar to responses obtained to time-normal versions suggesting that the cue

controlling behavioral responses was the phonetic structure of the words.  These results show

that chinchillas used different acoustic cues than human listeners.  The ability of chinchillas to

recognize noise-vocoded words as being perceptually equivalent to the naturally-spoken versions

is inferior compared to that of human listeners.  The findings suggest that the ability of human

listeners to recognize highly degraded speech is unlikely to be based solely on the general

auditory and perceptual mechanisms that are common among mammals.

**Introduction**

Humans possess a remarkable ability to recognize and understand speech when the speech signal is highly degraded, and the existence of specialized speech mechanisms could account for this ability (e.g. Remez, Rubin, Pisoni & Carrell, 1981; Remez, Rubin, Berns, Pardo & Lang, 1994).  However, the ability to recognize degraded speech alone provides insufficient evidence for the existence of specialized speech mechanisms.  In order to argue for the existence of speech-specific mechanisms, it is necessary to demonstrate an inability to perceive degraded speech based solely on general auditory processing mechanisms (Fitch, 2011).  A comparison of speech perception in humans to speech perception based solely on general auditory processing mechanisms in non-human mammals provides constructive insights into understanding the adaptations that may be enhanced or specialized for speech processing in humans.  Only a few studies have examined the perception of degraded speech by non-human mammals (Heimbauer, Beran & Owren, 2011; Ranasinghe, Vrana, Matney & Kilgard, 2012; Shofner, 2014).  Of particular interest here, is the study of Heimbauer et al. (2011), in which they report on noise-vocoded (NV) word recognition in Panzee, a linguistically trained chimpanzee.  Noise-vocoding is a common method for degrading speech sounds used in human perceptual studies, in part, because of the relationship between vocoding and the processing in cochlear implants.

Heimbauer et al. (2011) showed a parallel in NV word recognition performance for the same set of words between Panzee and a group of human listeners, suggesting that the mechanisms underlying degraded speech perception may have been present in the common ancestor of humans and chimpanzees.  One conclusion of their paper is that linguistic experience plays a critical role in speech perception in humans and Panzee.  Although the role of linguistic experience may seem important intuitively, its importance for degraded speech perception is

4

75  challenging to test directly.  Specifically, it is critical to measure degraded word recognition

76  performance in the absence of linguistic experience, but the absence of linguistic experience is

77  difficult to achieve.  For example, one approach might be to study vocoded word recognition in

78  human listeners using nonsense words.  However, Remez et al. (1981) demonstrated that when

79  listeners were presented with sentences based on sine-wave speech, they did not recognize the

80  sounds as speech, but when listeners were simply told the sound they were hearing was

81  computer-generated speech, listeners could then recognize and understand a substantial number

82  of words in the highly degraded, sine wave sentences.  Thus, if listeners are presented with

83  natural and vocoded nonsense words, there is no guarantee that they would not be tapping into

84  existing linguistic mechanisms.  Another approach might be to compare speech perception by

85  adult cochlear implant users with pre-lingual deafness to those with post-lingual deafness (Teoh

86  et al., 2004a).  However, this approach may be confounded by deafness-induced degeneration of

87  central auditory structures along the auditory pathway (Teoh et al., 2004b).  An obvious

88  approach would be to repeat the experiments of Heimbauer et al. (2011) with chimpanzees that

89  are not linguistically trained as Panzee was.  However, this is challenging given the present

90  limitations of using chimpanzees in biomedical and behavioral research (Institute of Medicine,

91  2011).  We argue that the chinchilla offers a good alternative animal model to this latter

92  approach.

93        Chinchillas are rodents, but unlike many other rodents, they have a range of hearing

94  similar to human listeners (Heffner and Heffner, 1991).  Although psychophysically measured

95  thresholds are generally higher in chinchillas than humans, functional relationships are often

96  similar between chinchillas and humans as with frequency discrimination (Nelson & Kiester,

97  1978) and noise intensity discrimination (Shofner, Yost & Sheft, 1993, Shofner & Sheft, 1994).

98    Chinchilla auditory filters derived from simultaneous masking using notched-noise are similar to

99    those of humans (Niemiec, Yost & Shofner, 1992), and chinchillas appear to possess a spectral

100   dominance region and missing fundamental percept for pitch that are similar to those of humans

101   (Shofner & Yost, 1997; Shofner, 2011).  Chinchillas show phonetic boundaries consistent with

102   categorical perception of voice onset time that are similar to the boundaries of humans (Kuhl &

103   Miller, 1975; 1978; Ohlemiller, Jones, Heidbreder, Clark & Miller, 1999).  However, unlike

104   humans and Panzee, chinchillas lack linguistic experience.  Thus, we argue that behavioral

105   responses obtained from chinchillas to degraded words will reflect speech perception based

106   solely on the general auditory and perceptual mechanisms that are common between chinchillas

107   and humans, and are presumably common among mammals.  The present study reports on the

108   recognition of NV monosyllabic words by chinchillas as measured in a stimulus generalization

109   paradigm.

110        In stimulus generalization paradigms, an animal is trained to respond to a specific

111   stimulus and then responses are measured to test stimuli that vary systematically along one or

112   more stimulus dimensions (Malott & Malott, 1970).  The signal is presented frequently and is

113   thus expected whereas test stimuli are presented infrequently and are unexpected.  A systematic

114   change in behavioral response along the physical dimension of the stimulus is known as a

115   generalization gradient and is consistent with the hypothesis that the animal possesses a

116   perceptual or psychological dimension related to the physical dimension (Guttman, 1963).  Thus,

117   data from stimulus generalization paradigms can indicate what acoustic features control the

118   behavioral response of the animal.  Generalization data are often interpreted to indicate

119   similarities in an animal's perception between test and signal stimuli (Guttman, 1963).  Test

120   stimuli that evoke similar behavioral responses as the signal stimulus indicate a perceptual

121 equivalence (Hulse, 1995) among these stimuli. Rock, Lasker & Simon (1969) argue that

122 generalization by animals occurs through recognition processes. Thus, we interpret

123 generalization to NV words to reflect the recognition of NV words.

124 How then does an animal like the chinchilla, which has the same basic auditory system as

125 a human but lacks any speech-specific mechanisms or linguistic experience, perceive a word like

126 "cut", for example? To be more specific, what are the acoustic cues available in the word that

127 control the behavioral response of the animal in the generalization paradigm? Words will have

128 no meaning to the chinchillas and will simply be a type of complex sound. If there is a

129 perceptual equivalence among vocoded test words and the naturally-spoken words *and* the

130 acoustic cues controlling the behavioral responses in the generalization paradigm are the

131 phonetic structures of the words, then it would suggest that the mechanisms underlying word

132 recognition are similar between humans and chinchillas. The results of the present study indicate

133 that the acoustic cues controlling the behavioral responses of the chinchillas differ from those

134 used by human listeners suggesting that general auditory and perceptual mechanisms alone are

135 inadequate to account for word recognition.

136

137 **General Methods**

138 The procedures used were approved by the Institutional Animal Care and Use Committee

139 and the Institutional Review Board for Indiana University. All human participants provided

140 informed consent.

141

142 *Subjects*

143        Five adult chinchillas (*Chinchilla laniger*) and 16 human listeners served as subjects in

144    these experiments.  All 5 chinchillas had experience with the stimulus generalization task

145    (Shofner, 2014).  Chinchillas received food pellet rewards during behavioral testing, and their

146    body weights were maintained between 85-90% of their normal weight.  American English was

147    the first language for all 16 human participants, and audiometric thresholds for 125-8000 Hz

148    were $\leq$ 25 dB HL.  Listeners were paid an hourly stipend for participation.

149

150    ***Acoustic Stimuli***

151        Naturally-spoken monosyllabic words spoken by a female voice were obtained from the

152    Lexical Neighborhood Test (Kirk, Pisoni & Osberger, 1995).  Table 1 summarizes the formant

153    and fundamental frequencies (F0s) of the naturally-spoken words used in the present study; the

154    spectral analysis of the words was carried out using Praat (http://www.fon.hum.uva.nl/praat).

155    NV versions of these words were generated using Tiger CIS version 1.05.02 developed by Qian-

156    Jie Fu (http://tigerspeech.com).  A naturally-spoken word was first passed through a series of

157    bandpass filters from 200-7000 Hz.  Default filter slopes of 24 dB/octave and center frequencies

158    based on the Greenwood function were used.  The number of channels was fixed at 1, 2, 4, 8, 16,

159    32, 64 or 128.  The envelope was extracted for each channel with the lowpass cut-off frequency

160    fixed at 160 Hz.  The carrier type of the vocoder was set to white noise; in this mode, the

161    extracted envelope from a given channel is used to modulate a wideband noise.  This modulated

162    wideband noise was then bandpass filtered with a center frequency equal to that of the analysis

163    channel.  The number of contiguous bandpass noises used for resynthesis equaled the number of

164    channels used for analysis.  The bandpass-filtered, modulated noises were then summed to yield

165    the vocoded version of the word.  NV words were stored as wav files that were later converted to

8

166    16-bit integer files at a sampling rate of 50 kHz using Adobe Audition in order to be played

167    through Tucker-Davis System II modules.  The natural and vocoded words had durations of

168    approximately 500 ms.  Example waveforms and spectrograms for the word "sit" are illustrated

169    in Figure 1.  Example envelopes extracted from the waveforms by half-wave rectification and

170    low-pass filtering with a cutoff frequency of 100 Hz are illustrated in Figure 2.

171

172    *Testing Procedures*

173         Chinchillas were placed into a testing cage located in a single-walled sound attenuating

174    chamber.  A pellet dispenser was located at one end of the cage with a reward chute attached to a

175    response lever.  A loudspeaker was located next to the pellet dispenser approximately 30º to the

176    right of center at a distance of 6" in front of the chinchilla.  The sound pressure level was fixed at

177    73 dB SPL (A-weighted) for all sounds.  A standard sound was presented continually in 500-ms

178    bursts at a rate of one per second, regardless of whether or not a trial was initiated.  The 2-

179    channel version was used as the standard instead of the 1-channel version in order to avoid

180    introducing overall spectral shape as a confounding variable (see Shofner, 2014).  Chinchillas

181    initiated a trial by pressing down on the response lever; the holdtime varied randomly for each

182    trial ranging from 1.15-8.15 seconds for 4 chinchillas and from 1.15-6.15 seconds for a 5[th]

183    chinchilla.  After the lever was depressed for the required holdtime, two 500-ms bursts of a

184    selected sound were presented for that trial.  The response window was coincident with the

185    duration of the two 500-ms bursts, but began 150 ms after the onset of the first burst.  Thus, the

186    duration of the response window was 1850 ms.

187         The sounds presented during the response window could be signals, test sounds, or

188    standards.  A signal trial consisted of two bursts of the word the animal was trained to

9

189  discriminate (e.g. naturally-spoken word).  If the animal released the lever during the response

190  window of a signal trial, then this positive response was treated as a hit and was rewarded with a

191  food pellet.  A standard trial consisted of two additional bursts of the 2-channel standard.  If the

192  animal continued to depress the lever throughout the response window of a standard trial, then

193  this negative response was treated as a correct rejection.  Food pellet rewards for correct

194  rejections were generally not necessary to reinforce continued lever depression for 4/5

195  chinchillas; one animal did receive a food pellet for a correct rejection.  A test trial consisted of

196  two bursts of a test sound which was generally a NV version of the naturally-spoken word based

197  on 4- to 128-channels.  Chinchillas did not receive food pellet rewards for positive responses to

198  test stimuli.

199       Chinchillas were tested in blocks of either 40 or 10 trials.  Trials were presented

200  randomly with a given block of trials.  Within a block of trials, the signal (e.g. naturally-spoken

201  "cut") was presented on 60% of the trials, the standard (e.g. 2-channel NV "cut") was presented

202  on 20% of the trials, test stimulus #1 (e.g. 32-channel NV "cut") was presented on 10% of the

203  trails, and test stimulus #2 (e.g. 4-channel NV "cut") was presented on 10% of the trials.

204  Responses were collected for a minimum of 2000 total trials, which results in a minimum of 200

205  trials for each test sound.  The two test stimuli were changed after the minimum number of trials

206  was completed.  Animals completed all 6 test versions for one word before moving on to the next

207  word.  Responses obtained from the stimulus generalization task are in terms of percent.

208  Because percent is not normally distributed, percent responses were converted into percent

209  rationalized arcsine units (RAUs) using the formula described by Sherbecoe and Studebaker

210  (2004) in order to carry out inferential statistical analyses on the responses.

211  Human listeners sat in front of a computer keyboard and monitor in a double-walled

212  sound attenuating chamber and were tested in a single-interval, forced-choice procedure.  At the

213  beginning of a block of 40 test trials, listeners first heard the naturally-spoken target word and

214  were instructed to reply "yes" when they recognized that specific target word.  Listeners then

215  heard a Gaussian noise and were instructed to reply "no" if they did not recognize that target

216  word.  NV versions of 1, 2, 4, 8, 16, 32, 64 and 128 channels of the target word were presented

217  randomly 5 times each at 50 dB SL through Koss ESP/950 electrostatic earphones monaurally to

218  the left ear.  Listeners did not receive feedback for their responses.  Recognition functions were

219  collected for 200-520 total trials for each target word, which results in 25-65 trials for each test

220  word.

221

222  **Experiment 1: Responses to NV-test words**

223  *Introduction*

224  The purpose of this experiment was to determine whether chinchillas show a perceptual

225  equivalence between naturally-spoken words and NV words.  That is, are the behavioral

226  responses to vocoded test words similar to their responses to naturally-spoken words?

227

228  *Methods*

229  Chinchillas were trained to discriminate naturally-spoken words from the 2-channel NV

230  versions.  In this experiment, a naturally-spoken word was the signal, the 2-channel NV version

231  was the standard, and test sounds were 4-, 8-, 16-, 32-, 64-, and 128-channel NV versions of the

232  word.  The following six monosyllabic words were tested: "ball", "cut", "hot", "meat", "sit" and

233  "wet".

11

234

*Results*

236     The mean responses for individual chinchillas to the naturally-spoken words are high,

237 ranging from 93.7-98.9%, whereas the mean responses to the 2-channel versions are low,

238 ranging from 0.7-7.3% (Fig.3). However, the responses to the NV test words are highly variable

239 across animals. For example, the responses obtained from C12 and C47 to the vocoded test

240 words are generally below 20% even when the number of vocoder channels is as high as 32 or

241 above. In these cases, there are no systematic increases in behavioral responses as the number of

242 channels increases. The generalization gradients are relatively "flat" for these 2 chinchillas for

243 the 6 words tested. In contrast, for C24, C36 and C15, the generalization gradients obtained for

244 some of the words show a systematic increase as the number of channels increases (e.g. "ball"),

245 but the gradients obtained for other words are relatively "flat" (e.g. "meat"). That is, the

246 responses are more variable within these animals. Note that for these animals, the responses to

247 vocoded test words based on 4-16 channels are generally below 40% for C24 and C36, and

248 below 60% for C15. The largest responses obtained for C15, C24 and C36 are around 50-80%

249 for some, but not all, vocoded test words based on 32, 64 and 128 channels (Fig. 3). Thus, the

250 generalization gradients obtained from these 5 chinchillas indicate that some chinchillas can

251 show relatively large behavioral responses to some NV monosyllabic words, but only for a high

252 number of vocoder channels (i.e. 32-128 channels). However, even the largest responses

253 obtained to vocoded test words are generally well below the responses to the naturally-spoken

254 words.

255     The responses to NV monosyllabic test words averaged across animals and across words

256 appears to show a shallow gradient with a logarithmic increase in response as the number of

257 channels increases from 2 to 128 (Fig. 4A). The responses averaged across animals and words

258 ranges from 8% for 4-channels to 31.7% for 128-channels and fall well below the averaged

259 response of 96.9% obtained to the naturally-spoken words. A repeated-measures analysis of

260 variance (ANOVA) based on the chinchilla RAUs (Fig. 4B) shows a significant effect of

261 stimulus, $F(7, 21) = 26.9$, $p < 0.001$, $\eta^2 = 0.72$. A regression analysis was applied to the

262 averaged data covering the range of 2 to 128 channels. The regression line had the form of $Y =$

263 $m*\log_{10}(X) + b$, where Y is the percent RAU, m is the slope, X is the number of channels and b

264 is the Y-intercept. The slope of the mean generalization gradient is 26.08, and a two-tailed t-test

265 showed this is significantly different from a slope of 0 ($t = 5.815$; $p = 0.0021$).

266 There is a substantial difference between the average vocoded word recognition function

267 obtained from human listeners and the generalization gradient obtained from chinchillas (Fig.

268 4A). The percent responses averaged across individuals and across words for these specific

269 words show that recognition is around 90% with as few as 8 channels for the human listeners

270 which is well above the percent responses averaged across chinchillas and across words. The

271 goals of the remaining experiments are to determine specifically the reasons for the differences

272 in performance between chinchillas and humans.

273

274 **Experiment 2: Testing the role of listening experience**

275 *Introduction*

276 The purpose of this experiment is to test the role of listening experience in chinchillas,

277 but how can chinchillas acquire additional listening experience with words? Clearly, it is not

278 feasible to provide chinchillas with experience that is comparable to the lifetime of experience

279 that human listeners have with words. One would have to breed and raise chinchillas and then

280    attempt to provide the young animals with linguistic experience in a manner similar to that

281    provided to Panzee (Heimbauer et al., 2011).  Moreover, if chinchillas use the waveform period

282    of the vowel (i.e. vowel 1/F0) as the acoustic cue controlling their behavioral responses when the

283    signal in the generalization task is the naturally-spoken version (see Shofner, 2014), it is unlikely

284    that additional experience with the naturally-spoken words will change the use of periodicity

285    cues to the use of phonetic structure cues.  What is necessary is to provide chinchillas with

286    listening experience in terms of the phonetic structure of words by reducing the pitch cues.

287

288    *Methods*

289        Chinchillas were re-trained to discriminate 64-channel NV words from the 2-channel

290    versions.  There was no fixed number of training trials.  Once an animal reached a performance

291    of approximately 85%, training with a specific 64-channel word lasted over 560-4200 trials in

292    order to be assured that the animal had a sufficient amount of listening experience with the

293    degraded 64-channel word prior to testing in the generalization task.   Following re-training, 64-

294    channel NV words were used as the signals in this generalization task.  The same six

295    monosyllabic words used in Experiment 1 were used in Experiment 2. Figure 1 illustrates

296    example spectrograms for the naturally-spoken and 64-channel NV versions of "sit."  It can be

297    observed that the harmonic structure is reduced in the 64-channel NV version compared to that

298    of the naturally-spoken version. Although the first few harmonics are still represented in the 64-

299    channel version, chinchillas are unable to use this spectral information (Shofner and Chaney,

300    2013). Consequently, this means that the periodicity cues for pitch are reduced in the 64-channel

301    NV words.  However, the wide, dark bands indicating the formant frequencies appear to be

302    clearly present in the 64-channel NV words.  Because 64-channel NV words lack the strong

303 harmonic structures and waveform periodicities of the naturally-spoken words, this re-training

304 should provide animals with additional experience listening to the speech cues of the words

305 when pitch cues are reduced.

306

307 **Results**

308 Responses from three chinchillas averaged across words and animals are shown in Figure

309 5A. When 64-channel NV words are used as the signal, responses from chinchillas to other

310 vocoded versions as well as responses to the naturally-spoken versions are high and less variable

311 (red open circles and red lines in Fig. 5A) than when the naturally-spoken versions were used as

312 the signal (Fig. 4A). Note that not only are the responses of chinchillas to 8 or more vocoder

313 channels higher than before ($\geq 87\%$), but they are now similar to those of human listeners (Fig.

314 5A). This similarity in the responses of chinchillas and humans argues that the differences

315 observed previously between chinchillas and humans (Fig. 4A) are not simply a reflection of

316 different behavioral tasks. That is, chinchillas do give large behavioral responses to NV words

317 in the stimulus generalization task under the appropriate stimulus conditions. The differences in

318 behavioral responses by individual chinchillas when naturally-spoken words and 64-channel NV

319 words are used as signals are shown in Figure 5B-D in terms of RAUs. A 2-factor ANOVA was

320 carried out for RAUs obtained for 4-, 8, 16-, 32, 64- and 128- NV-words and showed significant

321 effects of the signal used (i.e. naturally-spoken vs. 64-channel), significant effects of the

322 stimulus, and significant interactions (see Table 2). The significant effect of stimulus appears to

323 largely be due to the lower RAUs generally observed for 4- and 8-channel NV words. The

324 significant effect observed for the signal suggests differences in the acoustic cues controlling the

325 behavioral responses when the signal is the naturally-spoken word or the 64-channel NV word.

15

326     That is, the acoustic cues controlling chinchilla behavioral responses are different in the context

327     of the naturally-spoken words and in the context of the 64-channel NV words.  Thus, listening

328     experience with reduced pitch cues does appear to lead to improvement in degraded word

329     recognition in chinchillas.

330

331     **Experiment 3: Testing the role of phonetic structure**

332     *Introduction*

333       The purpose of this experiment was to determine specifically if the acoustic cue

334     controlling the behavioral response when the 64-channel NV words were used as the signals was

335     phonetic structure.  The data described above in Experiment 2 suggest listening experience with

336     reduced F0 cues improves degraded word recognition by chinchillas to an extent that it is similar

337     to that of humans.  However, it has been argued that similarity in performance does not equate to

338     similarity of mechanisms (e.g. Trout, 2001).  For example, although the responses between

339     humans and chinchillas are similar after chinchillas were re-trained using 64-channel NV words

340     as signals (see Fig. 5A), it is possible chinchillas could be processing different cues than humans

341     use and thus could be using different mechanisms (Trout, 2001).  Consequently, the results of

342     Experiment 2 do not demonstrate that behavioral responses in the chinchilla are controlled by

343     phonetic structure.  Saberi and Perrott (1999) showed that speech is unintelligible when segments

344     of speech 200 ms or longer are time reversed.  Time reversing a word changes its phonetic

345     structure and consequently changes the perception of the word.

346

347     *Methods*

348       In the generalization task used in this experiment, 64-channel NV words were used as

16

349  signals and the 2-channel versions were used as the standards.  Monosyllabic words used were

350  "ball", "hot" and "cut".  Test words included 16-, 32- and 128-channel NV versions of the words

351  and time-reversed versions of 16-, 32-, 64- and 128-channel NV versions.  Examples of

352  waveforms as well as the temporal envelopes for time-normal and time-reversed 64-channel

353  "cut" are illustrated in Figure 6, and example spectrograms are illustrated in Figure 7.

354

355  *Results*

356  We verified in a group of 8 human listeners from the original 16 using a single-interval

357  task that when the vocoded words are time reversed (approximately 500 ms), they are no longer

358  recognized as the time-normal vocoded words.  The responses to 16-, 32-, 64- and 128-channel

359  time-normal NV words are high (see red filled triangles and red dashed lines in left-hand column

360  in Fig. 8), because listeners recognize the vocoded versions as being perceptually equivalent to

361  the naturally-spoken target word.  However, responses to time-reversed versions of the NV

362  words are virtually 0 (see blue filled circles and blue dashed solid line in left-hand column in Fig.

363  8).  Listeners do not recognize time-reversed vocoded versions of words as being perceptually

364  equivalent to the naturally-spoken words, because time reversing alters the detailed phonetic

365  structure of the words.  If chinchillas show similar behavioral responses in the stimulus

366  generalization paradigm, it would argue that the acoustic cues controlling their behavioral

367  responses to the vocoded words are the phonetic structures of the words.

368  Three chinchillas were tested in the stimulus generalization task using the time-normal

369  64-channel version as the signal with both time-normal and time-reversed NV words as test

370  stimuli (see right-hand column of Fig. 8).  If chinchillas are responding to the phonetic structure

371  of the words, in the context of the 64-channel signal, then responses to the 16-, 32-, and 128-

372 channel versions should be high, whereas responses to the time-reversed 16-, 32-, 64- and 128-

373 channel versions should be low. That is, the time-reversed versions should not be perceptually

374 equivalent to the time-normal versions. Clearly, behavioral responses of chinchillas to the time-

375 reversed versions are high and appear to be well above the responses to the 2-channel versions

376 (right-hand column in Fig. 8), but they do not appear to be equal to those of the time-normal

377 words. The mean responses to the time-reversed NV "ball", "cut" and "hot" averaged across

378 channels and animals are 88%, 70.2%, and 72.4%, respectively, whereas the mean responses to

379 the time-normal NV "ball", "cut" and "hot" are 96.2%, 95.6%, and 97.1%, respectively. A two-

380 factor ANOVA showed that the differences between RAUs for each of time-normal and time-

381 reversed words were significant (see Table 3). Because there was not a significant effect of the

382 number of channels (Table 3), the responses of time-reversed 16-, 32-, 64- and 128-channel were

383 combined, and a one-tailed t-test showed that the mean RAUs for the time-reversed NV test

384 words was significantly larger than the mean RAUs for the 2-channel standards (see Table 4).

385 The effect sizes in terms of Cohen's d (Table 4) were converted into $\eta^2$ using formulae described

386 by Cohen (1988) in order to compare to the effect sizes given in Table 3. It can be observed that

387 the effect sizes for the differences between RAUs for time-reversed words and the 2-channel

388 standards (Table 4) are larger than the effect sizes for the differences in RAUs between time-

389 reversed and time-normal words (Table 3). Thus, unlike the human listeners, the large responses

390 to time-reversed words by chinchillas indicates that there is some degree of perceptual

391 equivalence between the time-reversed and time-normal NV words in the chinchillas suggesting

392 that the detailed phonetic structure of the words is not the acoustic cue controlling their

393 behavioral responses.

394

18

**Experiment 4: Testing with naturally-spoken words**

*Introduction*

Experiment 1 and our previous work (Shofner, 2014) suggest that fundamental

frequency (F0) of naturally-spoken vowels is a dominant acoustic cue.  Experiments 2 and 3

above suggest that following additional listening experience with the 64-channel NV versions,

chinchillas may be responding to some speech cues, but not phonetic structure per se.  The

purpose of this experiment was to test the saliency of pitch cues over speech cues in chinchillas.


*Methods*

Chinchillas were re-trained to discriminate naturally-spoken "cap" from the 2-channel

NV version.  Chinchillas had not been previously exposed to the naturally-spoken "cap."

Chinchillas were then tested in the generalization task using the naturally-spoken "cap" as the

signal and the 2-channel NV "cap" as the standard;  test stimuli were naturally-spoken "cut",

"hot", "ball" and "wet" as well as the musical note $G^b3$ as played on a piano and cello.  Musical

notes were obtained from the University of Iowa Electronic Music Studios

(http://theremin.music.uiowa.edu/).  Although there was some variation in stimuli in terms of

phonetic structures and formant structures, all test sounds had approximately the same F0 and

fundamental waveform period (Table 1).  Although the specific F0 difference limen (F0DL) is

unknown for 193 Hz F0, which is the F0 of the naturally-spoken "cap", the F0DL in chinchillas

for a closely comparable 250-Hz F0 harmonic tone complex is 30 Hz (Shofner, 2000).  Note,

however, that the differences between the F0 for "cap" and the F0s of the other stimuli are less

than 30 Hz.  Thus, because the F0 differences between the signal and test stimuli are below the

apparent F0DL, it is unlikely that F0 can be used by chinchillas as a cue to discriminate among

418     these sounds.

419

420     *Results*

421         Chinchillas gave high responses to all of the test sounds, including the two musical notes,

422     although there is some individual variability (Fig. 9A).  For example, the responses of C12 are

423     consistently > 90% for all test stimuli.  In contrast, although the responses of C24 are generally

424     high to the test words, they are lower for the two musical notes (Fig. 9A).  A repeated-measures

425     ANOVA on the RAUs averaged across the 3 chinchillas (Fig. 9B) showed a significant effect of

426     stimulus, $F(7, 14) = 33.39$, $p < 0.001$, $\eta^2 = 0.872$.  The 95% confidence intervals indicate that

427     responses to the natural-spoken "cap" signal and all test sounds are well above those of the 2-

428     channel "cap" standard.  Although the confidence intervals for the cello $G^b3$ fall short of

429     overlapping with those of the "cap" signal, the confidence intervals for all other test sounds,

430     including piano $G^b3$, do overlap with those of the signal (Fig. 9B).  Thus, there is a perceptual

431     equivalence among the signal and test sounds, in general.

432

433     **Discussion**

434     *Motivation in the context of each stimulus generalization task*

435         In the generalization task, animals are trained to discriminate a specific signal from a

436     standard.  There will be some acoustic cue(s) which will control the behavioral responses.

437     Animals are presented with the signal for most trials and receive a food reward for correct

438     responses.  Thus, animals will learn to expect the signal and receive reinforcement.  Test stimuli

439     are presented infrequently.  Given the expectation of a signal and reinforcement, if an animal

440     perceives that the test stimulus sufficiently contains the acoustic cue, it will be motivated to

441  respond, expecting a food reward. If the animal does not perceive the acoustic cue to be

442  sufficient in the test stimulus, there is no motivation for it to respond, because it would have

443  learned that these sounds are not reinforced with food. Thus, in the context of a specific signal

444  and standard, an animal will presumably respond to the acoustic cue that will maximize its food

445  reinforcement while minimizing its effort. Consequently, the stimulus generalization paradigm

446  allows us to deduce what acoustic cue is controlling the behavioral response. For each of the

447  generalization experiments, we want to deduce the cue being used by the chinchillas in the

448  context of discriminating the signal from the standard that will presumably maximize their food

449  reward and minimize their effort.

450      In the context of discriminating a naturally-spoken word from the 2-channel version

451  (Experiment 1), the responses obtained to NV test versions were generally low and were not

452  close to those obtained for the naturally-spoken version (Fig 4B). That is, the chinchillas showed

453  virtually no perceptual equivalence of the vocoded words to the naturally-spoken words,

454  suggesting that the phonetic structure is not the acoustic cue controlling the behavioral responses.

455  It is also unlikely that the temporal envelope is the acoustic cue, because the shapes of the

456  temporal envelopes for a naturally-spoken word and the 2-channel version are extremely similar

457  (e.g. Fig. 2). The formant structure for vocoded words having a high number of channels (e.g.

458  Fig. 1) would suggest that the responses should be higher to these test words if the formant

459  structure was the acoustic cue. However, our previous work with speech tokens (Shofner, 2014)

460  and NV harmonic tone complexes (Shofner and Chaney, 2013) argues that in this context, the

461  periodicity of the F0 is the most likely acoustic cue being used. The variability observed for

462  some chinchillas (e.g. C 24 and C36 in Fig. 3) which show 'flat' responses for some words, but

463  gradients in responses to other words, suggests that the acoustic cue(s) controlling behavioral

21

464 responses in these animals may be F0 for some words but a combination of F0 and formant

465 structure for other words. Thus, in the context of the naturally-spoken word, chinchillas

466 recognize the strong harmonic structures or more specifically the fundamental periods of the

467 vowels as well as voiced consonants (e.g. /m/, /l/, and /j/). Since this stimulus is presented most

468 of the time, chinchillas learn to expect this sound and receive food reinforcement. When the

469 fundamental period is weak or absent, as in the case of the vocoded test sounds, there is little or

470 no motivation for the chinchillas to respond, because they have learned that these sounds are not

471 reinforced.

472 When the chinchillas are re-trained to discriminate 64-channel noise-vocoded words from

473 the 2-channel version, then the responses to the other NV test words and the naturally-spoken

474 versions are high (Experiment 2). That is, there is now a perceptual equivalence among vocoded

475 words and the naturally-spoken words. In this case, the acoustic cue controlling the behavioral

476 response must be strongly similar among the vocoded test words and 64-channel signal in order

477 for the chinchillas to be motivated to respond. Given the similarity between the envelopes of the

478 2-channel standards and the 64-channel signals (e.g. Fig. 6), it is unlikely that the temporal

479 envelope is the acoustic cue. Since the harmonic structures of the vowels in the 64-channel

480 versions are weak, we argue that the acoustic cue is not likely to be F0, but rather is either the

481 phonetic structures of the words or the formant structures of the vowels. This was tested by

482 time-reversing the vocoded test words (Experiment 3).

483 Time-reversing changes the phonetic structure and the time-domain envelope (Fig. 6), but

484 it can also change the formant structure. For naturally-spoken words, the vocal tract resonance

485 can change over time resulting in dynamic changes in the trajectories of formants which are then

486 turned around when the word is time-reversed. Thus, if phonetic structure is the acoustic cue,

487    then chinchillas should recognize the difference between the time-reversed test words and the

488    time-normal 64-channel signal. Consequently, they should not be motivated to release the lever,

489    because they will learn that they will receive no food reward for this stimulus. Likewise, if

490    chinchillas are responding to dynamic changes in formant structure, they should recognize the

491    difference between time-normal and time-reversed words and not be motivated to respond to the

492    time-reversed words. However, if the dynamic changes in formant frequencies for the specific

493    words tested are too small to be detected by the chinchillas, such that the chinchillas only

494    recognize the *average* formant structure, then chinchillas will not recognize the difference

495    between the time-reversed test words and the time-normal 64-channel signal. Consequently,

496    they will be motivated to release the lever, because they will be expecting a food reward. Hence,

497    we conclude that average formant structure is the dominant acoustic cue controlling behavioral

498    responses for most trials in the context of the 64-channel signals, but acknowledge that

499    chinchillas may recognize dynamic changes in formant structure for some trials.

500    Finally, when chinchillas are re-trained using a naturally-spoken "cap" as the signal and

501    tested with other naturally-spoken words as well as a musical note played by two different

502    musical instruments (Experiment 4), the responses to all test sounds were high and not

503    substantially different from the signal as indicated by the 95% confidence intervals. In this

504    context, the phonetic structure, formant structures (Table 1) and temporal envelopes of the test

505    words are all clearly different from those of the signal. If any of these features are the acoustic

506    cue(s) being used to control behavioral responses, then the chinchillas should recognize the

507    difference between the test word and the signal. Consequently, they will not be motivated to

508    release the lever, because they will have learned that no food reward will be given for "non-cap"

509    sounds. However, all of these words have similar fundamental period of voice pitch. If

23

510   fundamental period of the vowels (i.e. 1/F0) is the acoustic cue, then the chinchillas should

511   recognize the test words to be equivalent to the signal word and thus, will be motivated to release

512   the lever expecting a food reward.  We conclude based on the high responses to test words that

513   the fundamental period is a highly salient cue that is controlling the behavioral responses in the

514   context of naturally-spoken words.  This conclusion is reinforced by the high responses to the

515   piano and cello notes.  However, we acknowledge that other phonemes may also play a

516   secondary role in some chinchillas.

517

518   *Comparison with Panzee and humans: Implications for the evolution of speech perception*

519   The current results are in contrast to results obtained from a chimpanzee, Panzee, in

520   which recognition of 7-channel NV words paralleled that of a group of human listeners

521   (Heimbauer et al., 2011).  The recognition for the chinchillas averaged across animals and all six

522   8-channel NV words tested was around 15%, whereas the recognition averaged across all 7-

523   channel NV words tested in Panzee was around 55% (Heimbauer et al., 2011).  A major

524   difference between these two studies is that the words presented to the chinchillas clearly lacked

525   any linguistic meaning, whereas Panzee was a linguistically trained chimpanzee (see Heimbauer

526   et al., 2011 for details) and the words presented to Panzee did have linguistic meaning.  Note,

527   however, that increasing the *listening* experience of chinchillas by re-training on a 64-channel

528   NV word does not equate to *linguistic* experience.  Linguistic experience implies that the word

529   has meaning to the listener and is part of a listener's lexicon.  For example, "ball" is simply a

530   complex sound and possesses no meaning for the chinchillas.

531   Chinchillas can discriminate vowels (Burdick & Miller, 1975), and appear to respond to

532   the average formant structure of NV words when they are re-trained using 64-channel NV words

24

533    in the generalization paradigm.  However, in the context of naturally-spoken words, the formant

534    cues appear to be less salient than the F0 cues in chinchillas.  Even when F0 cues are reduced, as

535    when chinchillas were re-trained using 64-channel NV words, subsequent testing with time-

536    reversed NV words suggests that chinchillas were responding to the average formant structure,

537    not to the detailed phonetic structure.  We conclude that the acoustic cues available in these

538    words that are being used by chinchillas are different from those used by humans and Panzee.

539    Whether or not chinchillas can learn to use phonetic structure as a cue is certainly an interesting

540    question and would be the next logical step in this line of research.

541    Chinchillas are phylogenetically further removed from humans (Huchon et al., 2002) than

542    are chimpanzees (Uddin et al., 2004).  Because the common mammalian ancestor of chinchillas

543    and humans is more distant than that of humans and chimpanzees, the perceptual and cognitive

544    mechanisms underlying speech perception by chinchillas will be based on mechanisms far more

545    ancestrally shared than those common to chimpanzees and humans.  Thus, we contend that

546    chinchilla behavioral responses illustrate degraded speech recognition based solely on the

547    general auditory and perceptual mechanisms that are common among mammals.  It has been

548    argued that the ability of humans to recognize highly degraded speech cannot be taken as

549    evidence for the existence of speech-specific mechanisms without corresponding data

550    establishing that this ability is absent in animals (Fitch, 2011).  The lack of equivalent

551    generalization among vocoded words to their naturally-spoken versions reported in the present

552    study for chinchillas suggests an inability as compared to humans.  More importantly, the results

553    of the present study show that chinchillas appear to use different cues than human listeners.

554    Chinchillas appear to learn to use formant structure as a cue when F0 cues are reduced

555    (Experiments 2 and 3), but in the present study chinchillas never learned to use phonetic

556  structure as a cue.  The use of different cues by chinchillas argues that general auditory and

557  perceptual mechanisms are insufficient to fully account for the ability to recognize degraded

558  speech that is observed for human listeners (and Panzee).

559      The present findings provide evidence that linguistic experience is critical for degraded

560  word recognition, consistent with the conclusions of Heimbauer et al. (2011).  Whereas Panzee

561  had linguistic experience and could recognize degraded words, the chinchillas lacked linguistic

562  experience and could not recognize degraded words.  Shannon (2005) concludes that "…speech

563  recognition is primarily a top-down process of pattern recognition that is highly overlearned

564  from a lifetime of experience."  One might argue that this top-down processing reflects

565  mechanisms that exists beyond the general auditory mechanisms and is an adaptation for

566  learning to recognize the phonetic structure of words.  The results of the present study suggest

567  that these top-down processing mechanisms important for word learning are inadequate or non-

568  existent in chinchillas.  Although this conclusion would seem to be consistent with the 'speech-

569  is-special' hypothesis, it does not imply that the mechanisms important for word recognition are

570  unique to humans.  Indeed, this adaptation also appears to have evolved in chimpanzees

571  (Heimbauer et al., 2011) and domestic dogs (e.g. Pilley & Reid, 2011).  Given that chimpanzees

572  are the closest relative to humans genetically (Uddin et al., 2004), one could argue that this

573  adaptation presumably appeared in the common ancestor of humans and chimpanzees

574  (Heimbauer et al., 2011); that is, these top-down mechanisms are homologous in humans and

575  chimpanzees.  The evolution of domestic dogs parallels the evolution of humans (Wang et al.,

576  2013), but it is likely that these top-down adaptations are the result of convergent evolution.

577  That is, these mechanisms presumably evolved independently in dogs and humans.  Thus, this

578  adaptation for top-down processing may not be unique to humans, but its existence appears to be

579  only in mammals that have demonstrated abilities to associate words with objects.

580

581  *Limitations of the present study*

582  The design of the present study is one that is typically used in psychoacoustics studies,

583  namely a repeated measures design in which a limited number of subjects are each tested with a

584  fixed set of stimulus conditions.  This general approach has been criticized as lacking

585  independent sampling and has been referred to as pseudoreplication (Hurlbert, 1984).  Although

586  this approach does have limitations (Kroodsma, 1990), inferential statistical analysis can still

587  provide useful information (Oksanen, 2001).  Indeed, Oksanen (2001) questions whether

588  pseudoreplication is a spurious issue.  In order to avoid pseudoreplication for both stimuli and

589  subjects in the current study, we would need a total number of subjects of 6 words x N word-

590  samples (i.e. groups) per word x n subjects per group.  If, for example, we use 4 different

591  samples of each specific word (e.g. recorded from two female speakers and two male speakers),

592  and have 3 subjects in each group, then this design would avoid pseudoreplication by providing

593  independent sampling in terms of both stimuli and subjects, and thus, would be a stronger test

594  theoretically (see Kroodsma, 1990; see also Macgregor, 2000).  However, this example would

595  require a total of 72 trained chinchillas and would be untenable for the current study given the 3-

596  4 months typically required to train a naïve chinchilla, the per diem animal costs and housing

597  requirements for 72 chinchillas, the cost of buying 72 chinchillas and the fact that not every

598  naïve chinchilla is capable of learning the behavioral task.  Thus, a compromise must often be

599  made in animal psychoacoustics experiments.  The choice is either to recognize and accept some

600  level of pseudoreplication in an effort to obtain some insight into the question (as in the present

601  study) or not to attempt to answer the question at all (given that the alternative approach is

602    untenable).  As stated by Kroodsma (1990) "There is, of course, nothing technically 'wrong' with

603    doing experiments that have no replication of treatments, but one must be aware that the

604    hypotheses actually being tested are about specific exemplars…"  In fact, the hypotheses tested

605    in the present study are limited to the specific words used.

606         Finally, one other limitation is in regards to the behavioral procedures used.  Caution

607    should be exercised in comparing directly the responses of humans and chinchillas due to

608    differences in the behavioral tasks.  Humans responded in a single-interval, forced choice task

609    whereas chinchillas responded in the stimulus generalization task.  We have previously described

610    that when humans are tested in the generalization task, their responses are similar to those

611    obtained in the single interval task (Shofner, 2014).  The low responses obtained from chinchillas

612    when the naturally-spoken words were used as signals does not mean that the generalization task

613    will inherently produce low responses to test stimuli, because high responses were obtained to

614    the same test stimuli when 64-channel NV words were used as the signals.  The main difference

615    between tasks, however, is that humans were verbally instructed to respond when they recognize

616    the vocoded test word as the naturally-spoken target word.  As such, we know that the responses

617    obtained from humans reflect recognition of these specific vocoded words based on phonetic

618    structure.  Of course, we cannot instruct the chinchillas to respond to phonetic structure, but

619    rather must deduce what acoustic cue is controlling the behavioral responses of the chinchillas in

620    the generalization task in the context of the signal and standard.

621         Thus, in the present study, a limited number of animals and a limited set of monosyllabic

622    words were tested.  The specific word-stimuli used were not meant to represent ideal exemplars

623    of those words, but rather were used as samples in which specific acoustic features could be

624    analyzed.  What we are deducing in the present study are the acoustic cues that chinchillas are

625    using under the specific conditions in each generalization task.  The hypotheses tested and the

626    conclusions reached are limited to understanding the acoustic cues controlling behavioral

627    responses of chinchillas to these specific word stimuli and comparing whether these cues are the

628    same as the cues known to be used by human listeners for the identical stimuli.

629

630 **References**

631

632 Burdick, C.K., & Miller, J.D. (1975). Speech perception by the chinchilla: discrimination of

633 sustained /a/ and /i/. *Journal of the Acoustical Society of America,* 58: 415-427.

634 https://doi.org/10.1121/1.380686

635

636 Cohen, J. (1988). *Statistical Power Analysis for the Behavioral Sciences*, 2nd edition. Lawrence

637 Erlbaum Associates.

638

639 Fitch, W.T. (2011) Speech perception: A language-trained chimpanzee weighs in. *Current*

640 *Biology,* R544-R546. http://dx.doi.org/10.1016/j.cub.2011.06.035

641

642 Guttman, N. (1963). Laws of behavior and facts of perception. In S. Koch (Ed.), *Psychology: A*

643 *study of a science: vol. 5. The process areas, the person, and some applied fields: Their place in*

644 *psychology and in science* (pp. 114-178). New York, NY: McGraw-Hill.

645

646 Heffner, R.S., & Heffner, H.E. (1991). Behavioral hearing range of the chinchilla. *Hearing*

647 *Research,* 52,13-16. http://dx.doi.org/10.1016/0378-5955(91)90183-A

648

649 Heimbauer, L.A., Beran, M.J. & Owren, M.J., (2011) A chimpanzee recognizes synthetic speech

650 with significantly reduced acoustic cues. *Current Biology,* 21, 1210-1214.

651 http://dx.doi.org/10.1016/j.cub.2011.06.007

652

653   Huchon, D., Madsen, O., Sibbald, M.J.J.B., Ament, K., Stanhope, M.J., Catzeflis, F., de Jong,

654   W.W., & Douzery, E.J.P. (2002).  Rodent phylogeny and a timescale for the evolution of glires:

655   Evidence from an extensive taxon sampling using three nuclear genes.  *Molecular Biology and*

656   *Evolution,* 19, 1053-1065. http://dx.doi.org/10.1093/oxfordjournals.molbev.a004164

657

658   Hulse, S.H. (1995).  The discrimination-transfer procedure for studying auditory perception and

659   perceptual invariance in animals.  In G.M. Klump, R.J. Dooling, R.R. Fay, & W.C. Stebbins

660   (Eds.), *Methods in Comparative Psychoacoustics* (pp. 319-330). Basel, Switzerland: Birkhauser-

661   Verlag. http://dx.doi.org/10.1007/978-3-0348-7463-2_27

662

663   Hurlbert, S.H. (1984). Pseudoreplication and the Design of Ecological Field Experiments.

664   *Ecological Monographs,* 54, 187-211. http://www.jstor.org/stable/1942661

665

666   Institute of Medicine (2011). *Chimpanzees in biomedical and behavioral research: Assessing the*

667   *necessity.* Washington, DC: The National Academies Press.

668

669   Kirk, K.I., Pisoni, D.B., & Osberger, M.J. (1995). Lexical effects on spoken word recognition by

670   pediatric cochlear implant users.  *Ear & Hearing,16,* 470-481.

671   http://dx.doi.org/10.1097/00003446-199510000-00004

672

673   Kroodsma, D.E. (1990). Using appropriate experimental designs for intended hypotheses in

674   'song' playbacks, with examples for testing effects of song repertoire sizes.  *Animal Behavior*,

675   40, 1138-1150. https://doi.org/10.1016/S0003-3472(05)80180-0

676

677     Kuhl, P.K., & Miller, J.D. (1975). Speech perception by the chinchilla: Voiced-voiceless

678     distinction in alveolar plosive consonants. *Science,* 190, 69-72.

679     http://dx.doi.org/10.1126/science.1166301

680

681     Kuhl, P.K., & Miller, J.D. (1978). Speech perception by the chinchilla: Identification functions

682     for synthetic VOT stimuli. *Journal of the Acoustical Society of America,* 63, 905-917.

683     http://dx.doi.org/10.1121/1.381770

684

685     Macgregor, P.K. (2000). Playback experiments: Design and analysis. *Acta Ethologica.* 3, 3-8.

686     DOI: 10.1007/s102110000023.

687

688     Malott, R. W., & Malott, M. K. (1970). Perception and stimulus generalization.

689     In W. C. Stebbins (Ed.), *Animal psychophysics: The design and*

690     *conduct of sensory experiments* (pp. 363–400). New York, NY: Appleton-

691     Century-Crofts.

692

693     Nelson, D.A. & Kiester, T.E. (1978). Frequency discrimination in the chinchilla. *Journal of the*

694     *Acoustical Society of America,* 64, 114-126. https://doi.org/10.1121/1.381977

695

696     Niemiec, A.J., Yost, W.A., & Shofner, W.P. (1992). Behavioral measures of frequency

697     selectivity in the chinchilla. *Journal of the Acoustical Society of America*, 92: 2636-2649.

698     https://doi.org/10.1121/1.404380

699

Ohlemiller, K.K., Jones, L.B., Heidbreder, A.F., Clark, W.W., & Miller, J.D. (1999). Voicing

judgements by chinchillas trained with a reward paradigm. *Behavioural Brain Research,* 100,

185-195. https://doi.org/10.1016/S0166-4328(98)00130-2

Oksanen, L. (2001). Logic of experiments in ecology: is pseudoreplication a pseudoissue?

*Oikos*, 94, 27–38.  https://doi.org/10.1034/j.1600-0706.2001.11311.x

Pilley, J.W., & Reid, A.K. (2011). Border collie comprehends object names as verbal referents.

*Behavioural Processes,* 86: 184-195. https://doi.org/10.1016/j.beproc.2010.11.007

Ranasinghe, K.G., Vrana, W.A., Matney, C.J., & Kilgard, M.P. (2012). Neural mechanisms

supporting robust discrimination of spectrally and temporally degraded speech. *Journal of the*

*Association for Research in Otolaryngology,* 13: 527-542. https://doi.org/10.1007/s10162-012-

0328-1

Remez, R.E., Rubin, P.E., Pisoni, D.B., & Carrell, T.D. (1981). Speech perception without

traditional speech cues.  *Science,* 212, 847-950. http://dx.doi.org/10.1126/science.7233191

Remez, R.E., Rubin, P.E., Berns, S.M., Pardo, J.S. & Lang, J.M. (1994).  On the perceptual

organization of speech. *Psychological Review*, 101, 129-156. http://dx.doi.org/10.1037/0033-

295X.101.1.129

721

722 Rock, I., Lasker, A., & Simon, J. (1969). Stimulus 'generalization' as a process of recognition.

723 *American Journal of Psychology*, 82, 1-22. http://www.jstor.org/stable/1420604

724

725 Saberi, K., & Perrott, D.R. (1999). Cognitive restoration of reversed speech. *Nature,* 398: 760.

726 doi:10.1038/19652

727

728 Shannon, R.V. (2005) Speech and music have different requirements for spectral resolution.

729 *International Review of Neurobiology,* 70, 121-134. http://dx.doi.org/10.1016/S0074-

730 7742(05)70004-0

731

732 Sherbecoe, R.L. & Studebaker , G.A. (2004). Supplementary formulas and tables for calculating

733 and interconverting speech recognition scores in transformed arcsine units. *International*

734 *Journal of Audiology,* 43, 442-448. DOI: 10.1080/14992020400050056

735

736 Shofner, W.P. (2000). Comparison of frequency discrimination thresholds for complex and

737 single tones in chinchillas. *Hearing Research*, 149, 106-114. https://doi.org/10.1016/S0378-

738 5955(00)00171-4

739

740 Shofner, W.P. (2011). Perception of the missing fundamental by chinchillas in the presence of

741 low-pass masking noise. *Journal of the Association for Research in Otolaryngology*, 12, 101-

742 112. DOI: 10.1007/s10162-010-0237-0

743

744     Shofner, W.P. (2014). Perception of degraded speech sounds differs in chinchilla and human

745     listeners. *Journal of the Acoustical Society of America, 1*35, 2065-2077.

746     http://dx.doi.org/10.1121/1.4867362

747

748     Shofner, W.P. and Chaney, M. (2013). Processing Pitch in a Nonhuman Mammal (Chinchilla

749     laniger). *Journal of Comparative Psychology*, 127, 142–153.

750     http://dx.doi.org/10.1037/a0029734

751

752     Shofner, W.P. & Sheft, S. (1994). Detection of bandlimited noise masked by wideband noise in

753     the chinchilla. *Hearing Research*, 77, 231-235, https://doi.org/10.1016/0378-5955(94)90271-2

754

755     Shofner, W.P. & Yost, W.A. (1997). Discrimination of rippled-spectrum noise from flat-

756     spectrum noise by chinchillas: evidence for a spectral dominance region. *Hearing Research*, 110,

757     15-24. https://doi.org/10.1016/S0378-5955(97)00063-4

758

759     Shofner, W.P., Yost, W.A. & Sheft, S. (1993). Increment detection of bandlimited noises in the

760     chinchilla. *Hearing Research*, 66, 67-80. https://doi.org/10.1016/0378-5955(93)90261-X

761

762     Teoh, S.W., Pisoni, D.B. & Miyamoto, R.T. (2004a). Cochlear implantation in adults with

763     prelingual deafness. Part I. Clinical results. *Laryngoscope*, 114, 1536–1540.

764     https://doi.org/10.1097/00005537-200409000-00006

765

766     Teoh, S.W., Pisoni, D.B. & Miyamoto, R.T. (2004b). Cochlear implantation in adults with

767  prelingual deafness. Part II. Underlying constraints that affect audiological outcomes.

768  *Laryngoscope*, 114, 1714-1719. https://doi.org/10.1097/00005537-200410000-00007

769

770  Trout, J.D. (2001). The biological basis of speech: What to infer from talking to animals.

771  *Psychological Review,* 3, 523-549.  http://dx.doi.org/10.1037/0033-295X.108.3.523

772

773  Uddin, M., Wildman, D.E., Lui, G., Xu, W., Johnson, R.M., Hof, P.R., Kapatos, G., Grossman,

774  L.I., & Goodman, M. (2004)  Sister grouping of chimpanzees and humans as revealed by

775  genome-wide phylogenetic analysis of brane gene expression profiles.  *Proceedings of the*

776  *National Academy of Sciences, 101,* 2957-2962. http://dx.doi.org/10.1073/pnas.0308725100

777

778  Wang, G-d., Zhai, W., Yang, H-c., Fan, R-x., Cao, A., Zhong, L., Wang, L., Liu, F., Wu, H.,

779  Cheng, L-g., Poyarkov, A.D., Poyarlov,. N.A., Tang, S-s., Zhao, W-m., Gao., Y., Lv, X-m., Irwin,

780  D.M., Savolainen, P., Wu, C-I., & Zhang, Y-p. (2013). The genomics of selection in dogs and the

781  parallel evolution between dogs and humans.  *Nature Communications,* 4, 1-9. DOI:

782  10.1038/ncomms2814

783

784    **Figure Legends**

785

786    **Figure 1.** Time-domain waveforms (top) and narrowband spectrograms (bottom) are illustrated

787    for the word "sit". Panels show the naturally-spoken word and noise-vocoded versions based on

788    2, 8 and 64 vocoder channels. The bands in the spectrograms corresponding to the first, second

789    and third formant frequencies of the vowel are indicated by F1, F2 and F3, respectively.

790

791    **Figure 2.** Examples of temporal envelopes are illustrated for "cut" (A) and "sit" (B). Black

792    lines illustrate the 2-channel NV words; red lines illustrate the naturally-spoken words.

793    Envelopes were extract using Praat by half-wave rectification and low-pass filtering with an

794    upper cut-off frequency of 100 Hz.

795

796    **Figure 3.** Behavioral responses obtained for individual chinchillas in the stimulus generalization

797    task are shown for each of the 6 test words. Different symbols represent different test words.

798    Labels along the abscissa indicate the naturally-spoken word (nat) and the number of channels

799    for vocoded words. Percent responses are the number of positive responses divided by the

800    number of trials multiplied by 100.

801

802    **Figure 4.** A: Mean responses obtained from chinchillas (red solid line & red filled circles) and

803    human listeners (blue dashed line & blue filled squares) are compared. Naturally-spoken words

804    were used as the signal in the stimulus generalization paradigm. Responses are averaged across

805    individual listeners and across the 6 words. Error bars indicate $\pm$ 1 standard deviation of the

806    mean. B: Mean RAUs obtained from chinchillas (red solid line & red filled squares) for C12,

807    C24, C36 and C47.  Data for C15 were incomplete and not included in the ANOVA.  Open blue

808    triangles show mean RAUs when the data for C15 were also included.  Error bars indicate 95%

809    confidence intervals.

810

811    **Figure 5.**  Responses obtained when 64-channel NV words were used as signals in the

812    generalization task.  A.  Responses of 3 chinchillas averaged across animals and words are

813    illustrated by the solid red line and red open circles.  Error bars show $\pm$ 1 standard deviation.  For

814    comparison, the average responses from 16 human listeners are shown by the blue dashed line

815    and blue filled squares. RAUs averaged across words are shown for C12 (B), C24 (C) and C47

816    (D).  Red filled squares and red solid lines show responses obtained when the naturally-spoken

817    words were used as signals in the generalization task.  Blue filled circles and blue dashed lines

818    show responses obtained when the 64-channel NV words were used as signals.

819

820    **Figure 6.**  Example waveforms are illustrated for time-normal "cut" for 64-channel (A) and 2-

821    channel (B) noise-vocoded versions and for 64-channel time-reversed "cut" (C).  Corresponding

822    envelopes (D) are shown for the 2-channel time-normal word (2-ch) by black solid line, for the

823    64-channel time-normal version (64-ch) by the red solid line and for the 64-channel time-

824    reversed version (64r-ch) by the blue solid line.

825

826    **Figure 7.**  Spectrograms are illustrated for time-normal (A) and time-reversed (B) 64-channel

827    NV "cut".  F1, F2, and F3 illustrate formant frequencies of the vowel.

828

829    **Figure 8.**  Responses of 8 humans (left-hand column) and 3 chinchillas (right-hand column) to

38

830    time-normal and time-reversed noise-vocoded words.  X-axis shows the number of vocoder

831    channels.  Red filled triangles and red dashed lines show responses to time-normal words; blue

832    filled circles and blue solid line show response to time-reversed words.  The gray dotted line in

833    the right-hand column shows the responses obtained for the 2-channel standard in the

834    generalization task.

835

836    **Figure 9.**  (A) Bar graph illustrating responses of 3 chinchillas in the generalization task when

837    the naturally-spoken "cap" was the signal.  Stimuli presented are shown on the X-axis.  The

838    standard was the 2-channel NV "cap" (2-chan cap) and the signal was the naturally-spoken "cap"

839    (cap) as illustrated by the diagonally striped bars.  The responses to other naturally-spoken test

840    words and two musical instruments are indicated in the graph by the filled bars: C12 in blue; C47

841    in green; C24 in red.  (B) Bar graph illustrating mean RAUs from 3 chinchillas in the

842    generalization task when the naturally-spoken "cap" was the signal.  The standard was the 2-

843    channel NV "cap" (2-chan cap) and the signal was the naturally-spoken "cap" (cap) as illustrated

844    by the blue diagonally striped bars.  The responses to other naturally-spoken test words and two

845    musical instruments are indicated in the graph by the blue filled bars.  Error bars show 95%

846    confidence intervals.

847 **Table 1.**

848 *Formant and fundamental frequencies of the unmodified, non-vocoded sounds*

849

|  | F1 (Hz) | F2 (Hz) | F3 (Hz) | F4 (Hz) | F0 (Hz) | 1/F0 (ms) |
|---|---|---|---|---|---|---|
| 'ball' | 846 | 1156 | 3025 | 4193 | 180.8 | 5.53 |
| 'cap' | 944 | 1862 | 2034 | 4889 | 193.2 | 5.18 |
| 'cut' | 800 | 2144 | 3223 | 4441 | 186.5 | 5.36 |
| 'hot' | 963 | 1655 | 2758 | 4310 | 184.5 | 5.42 |
| 'meat' | 395 | 3029 | 3099 | 4856 | 192.5 | 5.19 |
| 'sit' | 610 | 1203 | 2325 | 3286 | 200.5 | 4.99 |
| 'wet' | 1067 | 2407 | 3470 | 4410 | 197.8 | 5.06 |
| Piano $G^b3$ | 753 | 1327 | 1871 | 2485 | 186.2 | 5.37 |
| Cello $G^b3$ | 951 | 1960 | 2574 | 3889 | 187.7 | 5.33 |

850
851

852    **Table 2.**
853    *2-factor analysis of variance for the effect of signal for C12, C24, C47*
854

|  | Factor a<br>Effect of signal<br>(natural vs. 64-channel) | Factor b<br>Effect of stimulus | Interaction<br>a x b |
|---|---|---|---|
| C12 | $F(1, 60) = 1158.0$<br>$p < 0.001$<br>$\eta^2 = 0.54$ | $F(1, 60) = 69.6$<br>$p < 0.001$<br>$\eta^2 = 0.23$ | $F(1, 60) = 61.4$<br>$p = 1.52588E{-}05$<br>$\eta^2 = 0.199$ |
| C24 | $F(1, 60) = 532.4$<br>$p < 0.001$<br>$\eta^2 = 0.35$ | $F(1, 60) = 99.9$<br>$p < 0.001$<br>$\eta^2 = 0.47$ | $F(1, 60) = 27.1$<br>$p = 1.52588E{-}05$<br>$\eta^2 = 0.126$ |
| C47 | $F(1, 36) = 1256.7$<br>$p < 0.001$<br>$\eta^2 = 0.55$ | $F(1, 36) = 65.9$<br>$p < 0.001$<br>$\eta^2 = 0.202$ | $F(1, 36) = 73.2$<br>$p = 1.52588E{-}05$<br>$\eta^2 = 0.225$ |

855
856

857 **Table 3.**

858 *2-factor analysis of variance for time-reversed noise-vocoded words*

859

|  | Factor a<br>Effect of time reversing | Factor b<br>Effect of number of channels | Interaction<br>a x b |
|---|---|---|---|
| 'ball' | $F(1,16) = 4.73$<br>$p = 0.04$<br>$\eta^2 = 0.178$ | $F(1,16) = 1.14$<br>$p = 0.30$<br>$\eta^2 = 0.128$ | $F(1,16) = 0.81$<br>$p = 0.38$<br>$\eta^2 = 0.091$ |
| 'cut' | $F(1,16) = 47.9$<br>$p < 0.001$<br>$\eta^2 = 0.584$ | $F(1,16) = 4.15$<br>$p = 0.06$<br>$\eta^2 = 0.152$ | $(F(1,16) = 1.89$<br>$p = 0.19$<br>$\eta^2 = 0.069$ |
| 'hot' | $F(1,16) = 38.2$<br>$p < 0.001$<br>$\eta^2 = 0.638$ | $F(1,16) = 0.54$<br>$p = 0.48$<br>$\eta^2 = 0.027$ | $F(1,16) = 1.36$<br>$p = 0.26$<br>$\eta^2 = 0.068$ |

860

861 **Table 4.**

862 *One-tailed t-test for time-reversed noise-vocoded words*

863

| | | |
|---|---|---|
| 'ball' | $t = 9.20$, $p < 0.001$ | Cohen's d = 6.00 ($\eta^2 = 0.9$) |
| 'cut' | $t = 6.90$, $p < 0.001$ | Cohen's d = 4.47 ($\eta^2 = 0.8332$) |
| 'hot' | $t = 7.10$, $p < 0.001$ | Cohen's d = 4.58 ($\eta^2 = 0.8398$) |

864

865

# "sit"

(A) "cut"

(B) "sit"

868

869

870

871

46

**(A) Grand Averages ± S.D.**

**(B) Grand Averages ± 95% C.I.**

872

873

**(A)   Grand Averages**

human

chinchilla

% Response

Stimulus

**(C)   C24**

% RAU

64-chan signal

nat signal

Stimulus

**(B)   C12**

% RAU

64-chan signal

nat signal

Stimulus

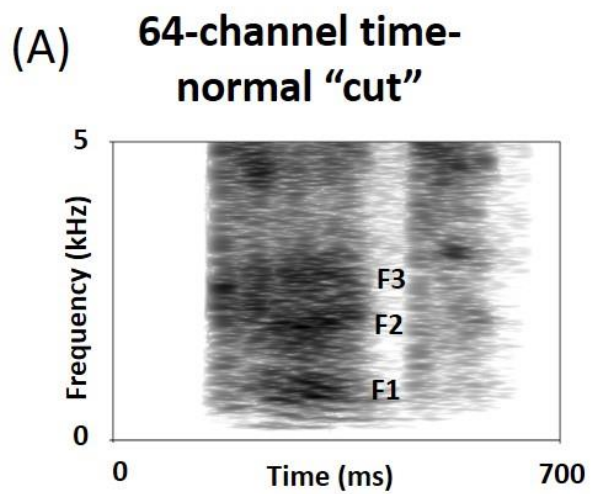**(D)   C47**

% RAU

64-chan signal

nat signal

Stimulus

874

875

48

(A) **64-channel time-normal "cut"**

(B) **2-channel time-normal "cut"**

(C) **64-channel time-reversed "cut"**

(D)

876

877

49

(A) **64-channel time-normal "cut"**

(B) **64-channel time-reversed "cut"**

878

879

## Human

### 'ball'

% Response

time-normal

16    32    64    128
Channels

### 'cut'

% Response

16    32    64    128
Channels

### 'hot'

% Response

16    32    64    128
Channels

## Chinchilla

### 'ball'

% Response

time-reversed

2-channel standard

16    32    64    128
Channels

### 'cut'

% Response

16    32    64    128
Channels

### 'hot'

% Response

16    32    64    128
Channels

880

881

(A)



(B)

882