1    **Does passive sound attenuation affect responses to pitch-shifted auditory feedback?**

2

3    Matthias K. Franken[1], Robert J. Hartsuiker

4    Experimental Psychology, Ghent University, Henri Dunantlaan 2, 9000 Ghent, Belgium

5

6    Petter Johansson, Lars Hall

7    Lund University Cognitive Science, Department of Philosophy, Lund University, Box
8    192, 221 00 Lund, Sweden
9

10

11    Tijmen Wartenberg

12    Hearing Technology @ WAVES, Information Technology, Ghent University,

13    Technologiepark-Zwijnaarde 126, 9052 Ghent, Belgium

14

15    Andreas Lind

16    Lund University Cognitive Science, Department of Philosophy, Lund University, Box
17    192, 221 00 Lund, Sweden
18

19

20

[1] Electronic mail: matthias.franken@ugent.be

24

25      **Abstract**

26      The role of auditory feedback in vocal production has mainly been investigated by altering

27      auditory feedback (AAF) in real time. In response, speakers compensate by shifting their

28      speech output in the opposite direction. Current theory suggests this is caused by a

29      mismatch between expected and observed feedback. A methodological issue is the difficulty

30      to fully isolate the speaker's hearing so that only AAF is presented to their ears. As a result,

31      participants may be presented with two simultaneous signals. If this is true, an alternative

32      explanation is that responses to AAF depend on the contrast between the manipulated and

33      the non-manipulated feedback. This hypothesis was tested by varying the passive sound

34      attenuation (PSA). Participants vocalized while auditory feedback was unexpectedly pitch-

35      shifted. The feedback was played through three pairs of headphones, with varying amounts

36      of PSA. The participants' responses were not affected by the different levels of PSA. This

37      suggests that across all three headphones, PSA is either good enough to make the

38      manipulated feedback dominant, or that differences in PSA are too small to affect the

39      contribution of non-manipulated feedback. Overall, the results suggest that it is important to

40      realize that non-manipulated auditory feedback could affect responses to AAF.

41

42

43    **I. INTRODUCTION**

44    An influential technique for investigating the interplay between speech and auditory

45    feedback is to alter auditory feedback in real time so that speakers hear their productions

46    perturbed in various ways (e.g., in pitch or formants). The dominant view in the field holds

47    that speakers usually compensate for feedback perturbations of pitch and formants because

48    they try to minimize the discrepancy between an internal representation of the sensory

49    speech target and the perceived auditory feedback (Hain et al., 2000; Liu and Larson, 2007).

50    This view, however, ignores a methodological issue associated with the altered auditory

51    feedback (AAF) technique: it is very difficult to completely rule out that speakers still

52    perceive their original, unperturbed feedback in addition to the manipulated signal. Thus, it

53    is possible that speakers receive conflicting evidence of what they are producing: their

54    actual, unperturbed, auditory feedback, and the altered auditory feedback provided by the

55    researchers. If so, an alternative explanation for compensation responses is that

56    compensatory responses depend on the conflict between two simultaneous auditory

57    feedback signals. The speaker, in the assumption that the dominant, manipulated, feedback

58    is self-produced, tries to minimize the discrepancy between the manipulated and the original

59    feedback, which leaks through the headphones and is considered as an external reference in

60    this scenario. The current study aims to test this alternative hypothesis.

61    Speakers receive both somatosensory as well as auditory feedback during speech

62    production. Auditory feedback is composed of both air-conducted and bone-conducted

63    feedback. While it is important to acknowledge the contribution of somatosensory feedback

64    and bone-conducted auditory feedback during speech production, the current study focuses

65    explicitly on air-conducted auditory feedback. The study and manipulation of (air-conducted)

66    auditory feedback through AAF has strongly advanced the field of speech motor control.

67    Studies using this technique have led to several theoretical frameworks for speech motor

68    control (Guenther, 2016; Houde and Nagarajan, 2011). In experiments that make use of AAF,

69    participants are instructed to speak, while their speech is being recorded with a microphone

70    and played back to them, near-simultaneously, through headphones. The experimenters

71    take control of the auditory feedback by manipulating it in real time, creating a discrepancy

72    between speech intent and the observed auditory signal. The type of manipulations that

73    have been applied include shifting the pitch (Burnett et al., 1998; Elman, 1981), formant

74    values (Houde and Jordan, 1998; Purcell and Munhall, 2006), or fricative noise (Casserly,

75    2011; Shiller et al., 2009) of the speech signal.

76         Most of these studies use one of two common paradigms. The first paradigm

77    ('adaptation') focuses on how speech production is affected after being exposed to altered

78    auditory feedback that is consistently altered in a specific manner. For example, when the

79    value of the first formant (F1) in the auditory feedback was gradually shifted upwards over

80    the course of an experiment, speakers responded by shifting the F1 in their speech in the

81    opposite way (i.e., downwards), and vice versa (Houde and Jordan, 1998; Jones and Munhall,

82    2000; Purcell and Munhall, 2006). These studies suggest that over time, speakers adapted to

83    consistently altered auditory feedback by changing their feedforward speech motor

84    commands (Franken et al., 2019). The second paradigm ('compensation') is aimed at

85    investigating how speakers respond to brief, unexpected changes in auditory feedback

86    during speech production. The present study makes use of this second AAF paradigm, in

87    order to investigate the effect of passive sound attenuation on immediate responses to

88    unexpected auditory feedback. This also allows us to investigate responses to feedback

89    perturbations of different magnitudes and directions. This is in contrast with the earlier

90    study by Mitsuya & Purcell (2016), where the 'adaptation' paradigm was used to investigate

4

91  adaptation to formant manipulations with either insert earphones or circumaural

92  headphones. While the authors concluded that the headphone type did not affect the

93  adaptation results, it is possible that headphone types will affect immediate responses. This

94  is viable since some recent studies have argued that compensation and adaptation are in

95  fact distinct processes (Franken et al., 2019; Parrell et al., 2017).

96      In the 'compensation' paradigm, speakers usually compensate for the altered feedback

97  by shifting their speech production in the opposite direction (Burnett et al., 1998; Hain et al.,

98  2000). For example, when pitch in the auditory feedback was shifted up, participants

99  responded by lowering their pitch, or vice versa. Interestingly, sometimes speakers may also

100 follow the feedback, by changing their speech in the same direction as the feedback

101 manipulation (Behroozmand et al., 2012; Franken et al., 2018a; Patel et al., 2014). Currently,

102 it is unclear what causes following responses, but multiple factors may play a role. Some

103 authors have suggested that following responses indicate that the feedback manipulation is

104 not considered to be self-generated but treated as an external referent, similar to a singer

105 trying to match the pitch of, for example, an accompanying piano (Hain et al., 2000; Patel et

106 al., 2014). Others have suggested following responses might have to do with the velocity of

107 the pitch shift (Guenther, 2016). A recent study has shown that the current state of the

108 speech system (i.e., ongoing pitch fluctuations) may affect whether a speaker opposes or

109 follows a pitch shift (Franken et al., 2018a). The neural correlates of following responses are

110 poorly understood, but recent studies claim that different neural mechanisms may underlie

111 following and opposing responses (Franken et al., 2018b; Li et al., 2013).

112      A methodological issue with AAF is that it is very difficult to fully isolate the speaker's

113 hearing so that only the altered feedback is presented to their ears. Many research groups

114 make use of commercial headphones (see Table 1 for a few examples) and these vary in how

115    much passive sound isolation they offer. As a result, the speaker may be presented with two

116    simultaneous auditory feedback signals: a) What they are actually uttering (the original

117    speech signal), and b) what is relayed through the headphones (the manipulated signal).

118    While two simultaneous auditory signals can be perceived as a single blended signal (Alain,

119    2007), small discrepancies, for example in pitch, may lead to a perception of two separate

120    signals. For example, perception of two simultaneous vowels is aided by small pitch

121    differences between the two vowels (Darwin, 1997; Darwin et al., 2003).

122

123    **Table 1.** Overview of different headphones used in published perturbation studies.

| Headphones | Type | Attenuation | Example studies |
| --- | --- | --- | --- |
| AKG boomset (K 270 H/C) | Circumaural | NA | Hain et al. (2000), Exp. Brain Research; Liu et al. (2010b), JASA |
| Etymotic Research ER | Insert Earphones | >30 dB | Cai et al. (2010), JASA ; Behroozmand et al. (2012), JASA |
| Sennheiser HD 280 Pro | Circumaural | up to 32dB | Franken et al. (2018a), PBR; Keough et al. (2009), JASA |
| BeyerDynamic DT 770 Pro | Circumaural | 18 dBA | Schuerman et al. (2017), JASA; |
| Stax SR001-MK2 | Insert Earphones | NA | Lametti et al. (2012, 2014), J of Neuroscience |

| Koss ESP950 | Circumaural | NA | Flagmeier et al. (2014), Brain and Language; Behroozmand et al. (2015), NeuroImage |

124

125

126    Thus, with low passive sound attenuation, the speaker could receive two conflicting

127    sets of evidence about what they are saying. Most studies increase the volume or add noise

128    to the manipulated signal to make it dominant over the original signal. However, it is very

129    difficult to completely rule out that participants still hear their original speech output. As the

130    presence of the original speech signal is often ignored, it is unclear how its potential

131    interaction with the manipulated signal may have affected the results in many of these

132    studies.

133    The present study aims at investigating the impact of sound attenuation observed in

134    typical headphones used in AAF experiments. Note that while we acknowledge that sound

135    attenuation has no impact on the contribution of bone-conducted auditory feedback, it will

136    affect the level of air-conducted auditory feedback leaking through the headphones, and

137    thus the overall level of non-manipulated auditory feedback. We first established the passive

138    sound attenuation offered by a number of different headphones (Experiment 1), and then

139    carried out an altered auditory feedback experiment (Experiment 2). This allowed us to

140    evaluate the effect of varying sound attenuation degree of different headphones on the

141    response to unexpectedly altered auditory feedback. Specifically, we aim to answer two

142    questions: (1) will passive sound attenuation affect the magnitude of the compensatory

143    response, and (2) will passive sound attenuation affect the likelihood of opposing responses?

144     The dominant view in the literature is reflected in models that suggest that

145     compensatory responses to altered feedback arise in order to minimize the discrepancy

146     between the intended speech target and the observed feedback signal (Guenther, 2016;

147     Hain et al., 2000; Houde and Nagarajan, 2011). If both the manipulated and the original

148     feedback signals are present, increased passive sound attenuation would make the

149     manipulated signal more dominant compared to the original feedback and thus make the

150     discrepancy between intended pitch and manipulated pitch more salient. Therefore, based

151     on the dominant theoretical framework, we would expect that increased sound attenuation

152     leads to stronger or more compensatory responses.

153     Alternatively, compensation could depend on the speaker hearing not only the

154     perturbed feedback, but also the (non-perturbed) normal feedback leaking through the

155     headphones. In other words, instead of an internal pitch target as the referent, the non-

156     manipulated auditory signal is considered the referent, as it is the louder of the two auditory

157     feedback signals. Compensatory behaviour could thus be a consequence of the perceived

158     mismatch between the two conflicting auditory signals that the speaker receives. This

159     hypothesis assumes that speakers consider the manipulated feedback as self-produced, and

160     thus try to minimize the mismatch by bringing this signal closer to the original ("actual")

161     feedback that leaks through the headphones. This would suggest that the intended speech

162     target (or an internal forward model prediction) plays a smaller role than often assumed, in

163     line with views that speech production targets are less well defined than most models

164     hypothesize, as it has been proposed for semantic aspects of language production by

165     inferential models (Lind et al., 2014). Note that we don't claim that speakers should be

166     consciously aware of the presence of two simultaneous auditory signals. Previous studies

167     have shown that responses to pitch-shifted feedback occurs automatically, even when

8

168 instructed not to (Hain et al., 2000). If this alternative hypothesis is true, increased passive

169 sound attenuation should lead to smaller compensations because this would decrease the

170 saliency of the conflict between the two auditory signals. With increased sound attenuation,

171 there would be less sound leaking through the headphones and thus the original, non-

172 manipulated feedback would be less salient, thereby reducing the conflict between two

173 simultaneous feedback signals.

174      Recent studies have shown that participants sometimes follow and sometimes oppose

175 pitch-shifted feedback (Behroozmand et al., 2012; Franken et al., 2018a). The alternative

176 hypothesis proposed here also provides a more straightforward account for the presence of

177 both opposing and following responses: if two simultaneous signals are perceived, the

178 response direction may depend on which of the two signals is considered by the participant

179 as under their control. To test this hypothesis, we ask whether the proportion of opposing

180 responses might be affected by sound attenuation. While different explanations have been

181 offered to explain following responses, an explanation based on source monitoring of the

182 auditory input as presented here is similar to the account by Hain et al. (2000), who suggest

183 that following might be appropriate when the speaker considers the incoming auditory

184 signal as externally generated, instead of being self-produced (Patel et al., 2014). If this is the

185 case, low passive sound attenuation and thus the presence of two simultaneous auditory

186 signals could lower the probability that the participant will consider the manipulated

187 feedback signal as self-produced, and thus increase the likelihood of a following rather than

188 an opposing response. Thus, this view would predict that increased sound attenuation would

189 lead to a higher proportion of opposing responses.

190      A recent study investigated a related, but different question (Mitsuya and Purcell,

191 2016). In order to investigate the role of the occlusion effect, the authors compared insert

192    earphones with circumaural headphones in an adaptation paradigm, and found no effect of

193    headphones on F1 adaptation. In other words, adaptation over time to a consistent

194    manipulation of F1 was not affected by the type of headphones. There is evidence ,

195    however, for the hypothesis that longer-term adaptation and immediate compensation to

196    unexpected feedback perturbations may be supported by two different mechanisms

197    (Franken et al., 2019; Parrell et al., 2017). If that hypothesis is correct, the type of

198    headphones could affect these processes in different ways. The current study will focus on

199    real-time compensation responses to unexpected feedback perturbations. In addition, the

200    earlier study compared two headphones in order to examine the role of the occlusion effect.

201    Although it is likely that the headphones used also differed in sound attenuation, the current

202    study will look at the effect of sound attenuation specifically in a pitch-shift compensation

203    paradigm.

204    **II. EXPERIMENT 1: HEADPHONES MEASUREMENTS**

205       In Experiment 1, we investigated the passive sound attenuation of one set of hearing

206    protection ear muffs and seven pairs of headphones. The goal was to have a comparable

207    measure of passive sound attenuation for each pair of headphones, in order to be able to

208    investigate its effect on responses to altered auditory feedback in Experiment 2. Although

209    headphones manufacturers provide sound attenuation measures, it is unclear what method

210    different manufacturers use and thus how these numbers could be compared across

211    headphones. In addition, we measured each headphones' frequency response to make sure

212    differences between the headphones' frequency responses were not a contributing factor to

213    the behavioural differences in Experiment 2.

214

215    **A. Methods**

216     *1. Headphones*

217     Four pairs of commercially available headphones were selected, as well as one pair of

218     hearing protection ear muffs. The headphones were chosen as they were all designed to

219     have high sound attenuation, and reflect the range of headphones commonly used for

220     speech manipulation research. The headphones included three closed-back circumaural

221     headphones, designed to have high passive sound attenuation, as well as the ER-3C insert

222     earphones (Etymotic Research), designed for research. The headphones are listed in Table 2,

223     along with the average attenuation magnitude as specified by the manufacturer. Many of

224     the headphones selected have been used in altered auditory feedback studies (see also

225     Table 1).

226

227     **Table 2.** Attenuation of headphones/ear muffs used in Experiment 1.

| Name | Type | Attenuation (according to manufacturer) |
|---|---|---|
| Peltor X5A | Hearing protection | 37 dB |
| Beyerdynamic DT 770 Pro | Closed-back circumaural | 18 dBA |
| Sennheiser HD 280 Pro | Closed-back circumaural | up to 32 dB |
| Vic Firth SIH1 | Closed-back circumaural | 24 dB |
| Etymotic ER-3C | Insert earphones | over 30 dB |
| Custom-built no. 1 | Closed-back circumaural | - |
| Custom-built no. 2 | Closed-back circumaural | - |
| Custom-built no. 3 | Closed-back circumaural | - |

228
229

230    In addition to the commercially available headphones, we custom-built headphones by

231    placing the loudspeakers (including their plastic casings) of Sennheiser HD 280 Pro

232    headphones into Peltor 3M X5A hearing protection ear muffs[1] (see supplementary

233    materials[2]). Since these custom-built headsets are not standardized, we built three copies of

234    the same design in order to see how they compare to each other. We included the hearing

235    protection ear muffs in our attenuation measurements to check how the construction of the

236    custom-built headphones affected the passive sound attenuation of the ear muffs. The

237    custom-built headphones were created in order to maximize sound attenuation with

238    circumaural headphones. While insert earphones could lead to better sound attenuation

239    still, circumaural headphones avoid an occlusion effect (Mitsuya and Purcell, 2016) and are

240    easier to use.

241

242    **_2. Equipment_**

243    For this study we used a Head And Torso Simulator (HATS) Type 4128-C placed in a

244    near-anechoic chamber (only the floor is not anechoic). The HATS is a model of a head and

245    torso designed for in-situ electroacoustic tests. It has models for the human pinnae.

246    However, in the current study the pinnae models were not used as they might interact with

247    some of the circumaural headphones (except for the measurements of the Etymotic

248    Research ER-3C insert earphones, where the pinnae were used). The HATS contains ear

249    simulators with ½'' microphones, which allows the researcher to record the sound reaching

250    the ears. For the attenuation measurements, acoustic stimuli were played from a single

251    ADAM S1X Active Studio Monitor placed at 1.5m in front of the HATS. At about 2.5cm in

252    front of the mouth of the HATS, a reference microphone (1/2" preprolarized free-field

253    microphone, Bruel & Kjaer type 4189) was placed. Microphones and speakers were

254    connected to a Bruel & Kjaer Input/Output Module (Type 3109).

255

256    ***3. Sound materials***

257    For the measurements of passive sound attenuation, a white noise stimulus was

258    created using Praat (Boersma and Weenink, 2017). In addition, for the frequency response

259    measurements we created stimuli with a male and a female speech-weighted speech

260    spectrum by taking the male-weighted and female-weighted speech-modulated noises from

261    the ICRA project (Dreschler et al., 2001) and randomly shifting the phases in Matlab

262    (Mathworks. Inc., R2016b). All stimuli had a duration of at least 25s.

263

264    ***4. Procedure***

265    In order to measure passive sound attenuation, each pair of headphones was placed

266    on the HATS, while white noise was played at 80 dB SPL through the studio monitor

267    (measured at the reference microphone in front of the HATS mouth). PULSE LabShop (Bruel

268    & Kjaer, v. 15.1.0) was used to control stimulus playback and to record signals from the in-

269    ear microphones as well as from the reference microphone in front of the HATS mouth.

270    Every measurement with the headphones was carried out twice, with repositioning the

271    headphones in between measurements to check for accuracy. The signals were transformed

272    to power spectra (in $Pa^2$) with 1/3 octave filter bands by averaging over a 20s time window.

273    Before every measurement with headphones, a measurement was carried out without

274    headphones to serve as a baseline measurement. The reference microphone in front of the

275    HATS mouth was used to control the stimulus volume across measurements online. In

276    addition, an offline analysis confirmed that the reference signal was not affected by the

277    presence or absence of headphones on the HATS.

278         For the measurements of the frequency responses of the headphones, acoustic stimuli

279    were played through each pair of headphones after it was placed on the HATS. Before each

280    measurement, it was made sure that the overall intensity level reaching the in-ear

281    microphones when the headphone was not mounted on the head was 80 dB SPL. Every

282    measurement was carried out twice with headphones repositioning in between. These

283    measurements were repeated with the white noise and the two speech-weighted stimuli.

284

285    ***5. Analysis***

286         The data and analysis scripts are publicly accessible via https://osf.io/vm84u/. All

287    further analyses were done in R (R Core Team, 2018) and focused on frequency bands

288    ranging from 100Hz to 8kHz, which include the frequencies most relevant for speech. The

289    power spectra were expressed in dBA. In order to calculate the attenuation in each

290    frequency band, the intensity in the corresponding frequency band in the baseline

291    measurement without headphones was subtracted from the intensity in the measurements

292    with headphones. This was done for both headphones measurements, after which the

293    results were averaged.

294         In order to quantitatively compare frequency responses to each other, two metrics

295    were used: spectral flatness, and the average root-mean-square error (Breebaart, 2017).

296    Spectral flatness was quantified as the dB-scaled ratio between the geometric and the

297    arithmetic mean of the power spectrum (Johnston, 1988):

298
$$Spectral\ Flatness = 10 \cdot \log_{10}\left(\frac{\sqrt[N]{\prod_{n=1}^{N} x(n)}}{\frac{1}{N}\sum_{n=1}^{N} x(n)}\right)$$

299    Where N is the number of frequency bands, and x(n) the power in frequency band n.

300    The spectral flatness measure has been used to quantify how flat (or noise-like) a spectrum

301    is. It is bounded between -∞ and 0. Given white noise as an input signal, a higher spectral

302    flatness score would thus indicate a frequency response that is closer to the input signal.

303    Only with a perfectly flat spectrum is the geometric mean equal to the arithmetic mean and

304    thus the spectral flatness score 0. Furthermore, frequency responses can be compared to

305    each other by looking at the root-mean-square error (RMSE) between two frequency

306    responses. The RMSE was calculated as follows:

307

308    $$RMSE = \sqrt{\frac{1}{N}\sum_{n=1}^{N}(x_1(n) - x_2(n))^2}$$

309    Where $x_i(n)$ indicates the power at frequency band n for headphones i. Both the

310    correlation coefficient and the RMSE value were calculated for each possible pair of

311    headphones, and averaged per headphone. The resulting average values indicate how well a

312    pair of headphones' frequency response compares on average to all the other pairs of

313    headphones. A low value indicates that the frequency response of the headphones is very

314    similar to the other headphones' frequency responses.

315

316    **B. Results**

317    Figure 1 shows the passive sound attenuation over the frequency range 100-8000Hz

318    for each pair of headphones for both the left and right ears. It is clear from the figure that

319    the passive sound attenuation varies across the frequency spectrum as well as headphones.

320    Note that the Etymotic ER insert earphones seem to be the most attenuating below 300Hz

321    and above 3000Hz, while the hearing protection ear muffs are the most attenuating

15

322    between 300 and 1600Hz. The different shape of the ER attenuation spectrum compared to

323    the other headphones could be due to the fact that these are the only insert earphones

324    compared to the other (circumaural) headphones, possibly leading to different in-ear

325    resonance frequencies. The fact that we used the HATS' pinna models for the ER

326    measurements but not for the circumaural headphones measurements could be an

327    additional contributing factor. However, for measurements conducted without headphones,

328    the addition of the pinnae models only lead to a slight amplitude increase between 2000 and

329    5000Hz, suggesting that the pinnae were not a major contributing factor to the spectral

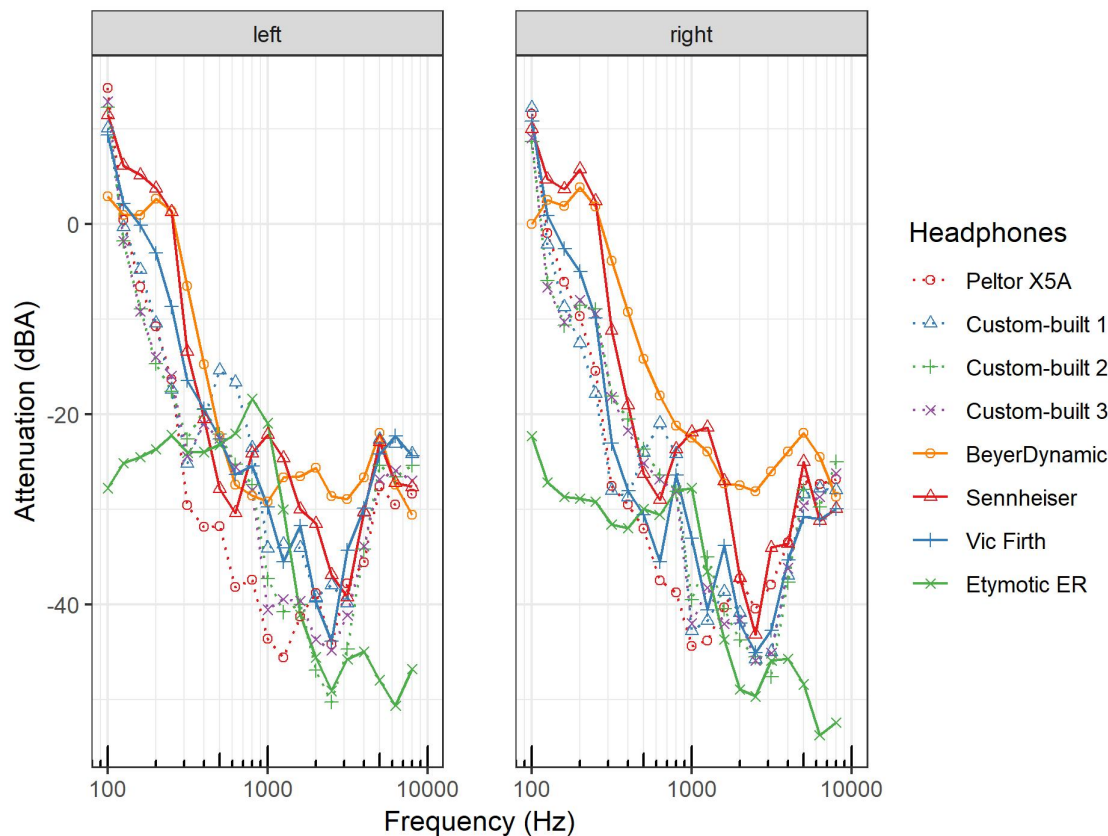330    differences observed between the ER and the other headphones.

331



332    **Fig. 1.** The measured passive sound attenuation over the frequency spectrum of 100-8000Hz

333    for both left and right ears.

334

335        Figure 2 shows the same data, this time averaged across the frequency range, which

336    allows for an overall measure of passive sound attenuation in speech-relevant frequencies. It

337    can be seen from both figures that passive sound attenuation varies across headphones,

338    from the pair of BeyerDynamic and Sennheiser headphones with relatively low attenuation

339    to the most attenuation in the hearing protection ear muffs (Peltor X5A) and the Etymotic ER

340    insert earphones. These values do not precisely correspond to the values provided by the

341    manufacturers as shown in Table 2. A comparison between headphones based on the

342    manufacturer-provided values is difficult, as manufacturers do not disclose how they arrived

343    at these values, and different manufacturers may use different measuring methodologies.
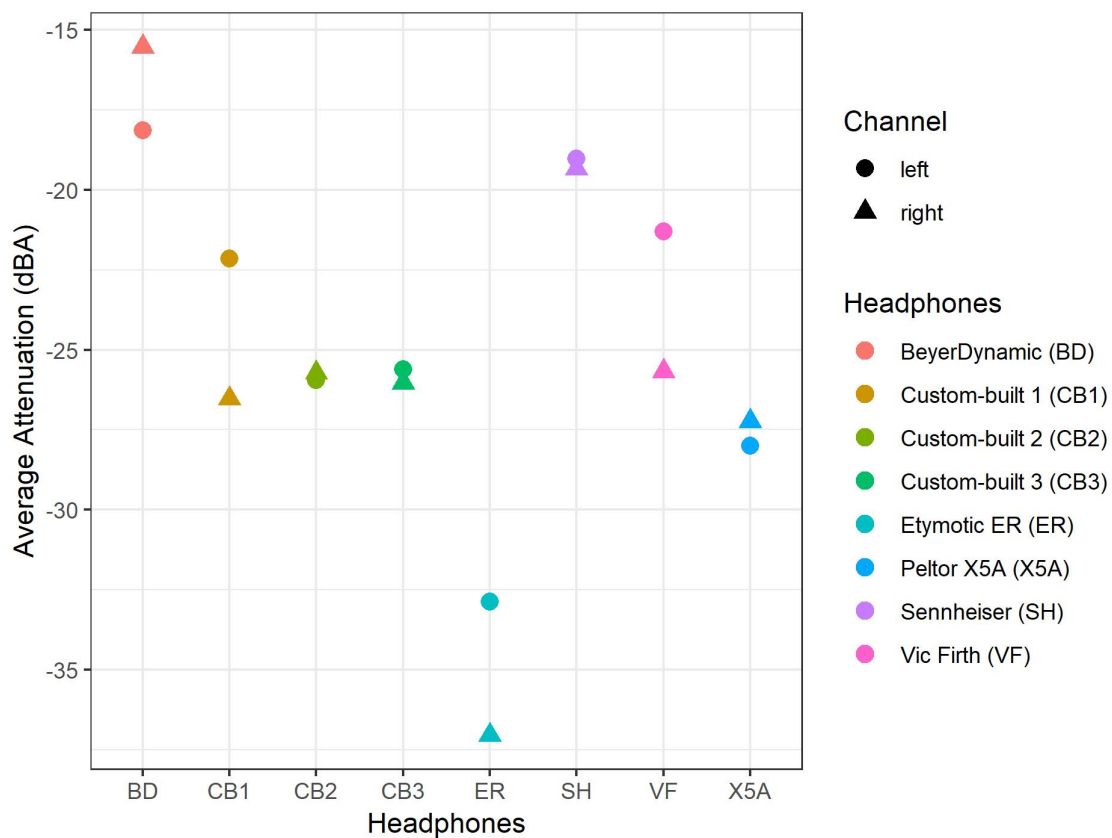


344

345    **Fig. 2.** Average sound attenuation for each pair of headphones, averaged across the 100-

346    8000Hz frequency range.

347

17

348        In order to make sure that any headphones-specific differences in passive sound

349    attenuation are not confounded with headphones-specific frequency response

350    characteristics, the frequency spectrum was quantified for each pair of headphones. First,

351    the spectral flatness of the frequency response to white noise input was quantified for each

352    pair of headphones, shown in Figure 3. Judging from the figure, there is a clear difference in

353    spectral flatness between the Vic Firth headphones and the other headphones. In addition,

354    we see a smaller difference for the left channel between the custom-built pair of

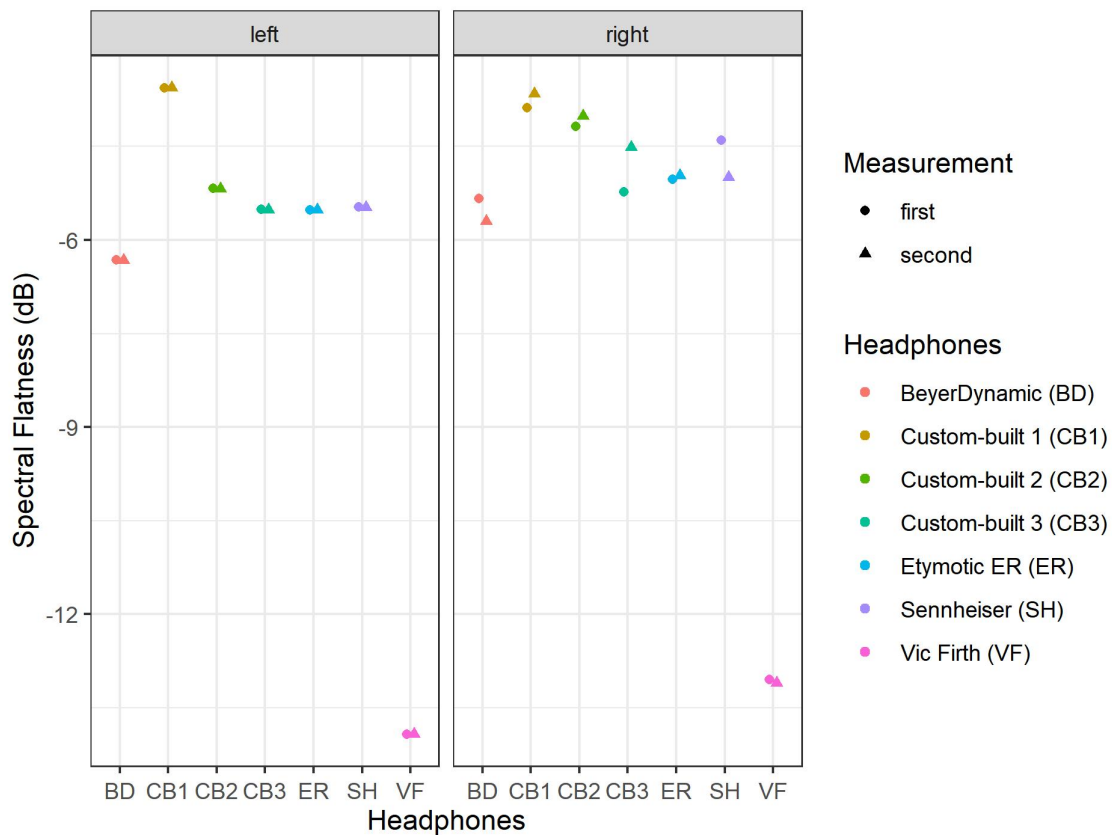355    headphones nr. 1 and the other custom-built pair of headphones.



356

357    **Fig. 3.** Spectral flatness of the headphones' frequency response to a white noise (WN) input

358    signal.

359

360    A second way to evaluate the differences between headphones' frequency responses

361    is to quantify the average RMSE between a headphones' frequency response and the

362    response of every other pair of headphones. This was done for the frequency response to a

363    white noise input signal, as well as for responses to male (ICRA4) and female (ICRA5) speech-

364    weighted noises, shown in Figure 4. The figure suggests that the Etymotic ER, the custom-

365    built headphones nr. 1, and the Vic Firth headphones show a frequency response which is
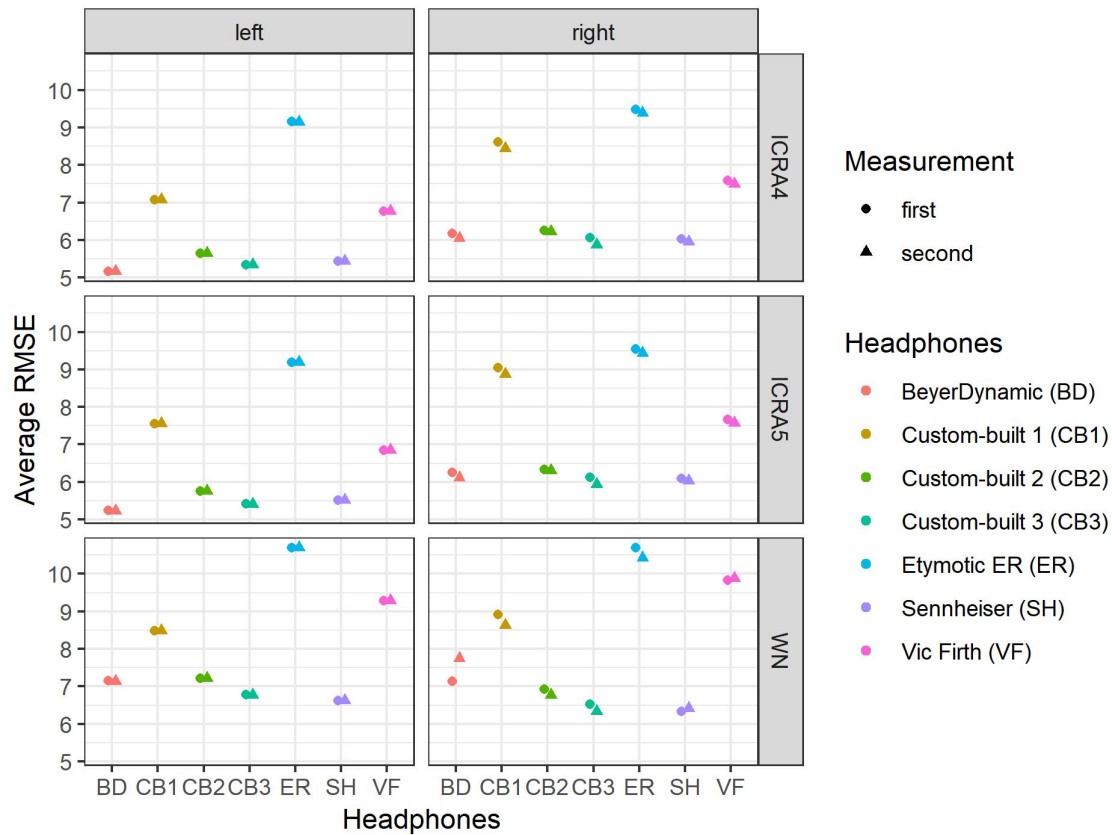
366    considerably different from the other headphones.

367

**Fig. 4.** Average RMSE of every headphones' frequency response.

369

370    **C. Discussion**

371    Overall, it can be concluded that the pairs of headphones tested in Experiment 1 show

372    variable passive sound attenuation. The BeyerDynamic and Sennheiser headphones show

373    the least sound attenuation, while the Etymotic insert earphones show the highest sound

374    attenuation. The custom-built headphones show medium sound attenuation, approaching

375    the values of hearing protection ear muffs. It should be noted that the high sound

376    attenuation in the Etymotic insert earphones is visible especially for low (<300Hz) and higher

377    (>3000Hz) frequencies. This is interesting in light of the evidence that in speech perception

378    important phonetic cues are conveyed between 100 and about 2000Hz (Epstein et al., 1968;

379    Warren et al., 1995), although higher frequencies also convey speech-relevant

380    information(e.g., for the perception of sibilants). So this result shows that in the ER

381    earphones, attenuation is especially strong in frequencies that are less relevant for many

382    speech sounds.

383        In order to maximize the range of passive sound attenuation while limiting the number

384    of headphones used, three pairs of headphones spanning the attenuation scale were

385    selected for use in Experiment 2: the BeyerDynamic headphones, a pair of custom-built

386    headphones (Nr. 3), and the Etymotic ER insert earphones. The BeyerDynamic are the least

387    sound-attenuating, the custom-built headphones offer intermediate attenuation, and the ER

388    offer most sound attenuation. This will allow us to interpret differences between

389    headphones in Experiment 2 as a function of passive sound attenuation. Both the

390    BeyerDynamic headphones and the Etymotic insert earphones have been used for altered

391    auditory feedback research in the past (see Table 1). It should be noted that the Etymotic ER

392    insert earphones are somewhat different from the other two, both in type (insert earphones

393    vs. circumaural headphones) as well as in the measured frequency responses (Figure 4). The

394    different frequency response for the ER could affect both the air-conducted auditory

395    feedback, as well as the relative contributions of air-conducted and bone-conducted

396    feedback to the overall auditory feedback, as these differ across the frequency range

397 (Pörschmann, 2000). This suggests we should take caution interpreting differences between

398 condition with ER earphones and the other two pairs of headphones in Experiment 2 as

399 being solely due to passive sound attenuation.

400     Finally, Experiment 1 shows that the construction of custom-built headphones by

401 placing Sennheiser headphones speakers into Peltor X5A hearing protection ear muffs was

402 successful, especially for pairs nr. 2 and nr. 3. They showed passive sound attenuation that

403 was not far from the attenuation measured for the Peltor X5A ear muffs, and their frequency

404 response measures were similar to the frequency response of the Sennheiser headphones

405 from which they were constructed.

406

407

408     **III. EXPERIMENT 2**

409     In order to investigate whether passive sound attenuation has an effect on speakers'

410 behaviour in a feedback perturbation experiment, three pairs of headphones were selected

411 based on their sound attenuation properties as measured in Experiment 1. Participants took

412 part in a pitch perturbation experiment with three blocks, one for each pair of headphones.

413 If responses to pitch perturbations depend on a comparison between the manipulated

414 feedback and an internal target representation, increased sound attenuation should lead to

415 stronger opposing responses (compared to weaker sound attenuation). If, on the other

416 hand, responses depend on a comparison between two simultaneous auditory signals,

417 increased sound attenuation should lead to smaller responses and/or more opposing

418 responses.

419

420     **A. Method**

421     *1. Participants*

422         49 native speakers of Dutch participated in the experiment in exchange for course

423     credit. All participants were students at Ghent University (41 female and 8 male, mean age =

424     19.4). None of them had any history of speech, hearing, or language impairments. The study

425     was approved by the ethics committee of the Ghent University faculty of psychology and

426     educational sciences.

427

428     *2. Procedure*

429         Participants were fitted with a pair of headphones and a head-mounted microphone.

430     On each trial, the appearance of the letters "EE" (pronounced in Dutch as [e]) on a laptop

431     screen provided a signal for participants to start vocalizing the vowel [e] and to hold the

432     vowel until the letters disappeared after 4s. Participants were instructed to try to keep the

433     volume, pitch, and articulation of the vowel constant. During vocalization, participants

434     received auditory feedback via the headphones. During each vocalization, pitch was shifted

435     for 200ms by either -25 cents, +25 cents, -100 cents, +100 cents, or 0 cents. This happened

436     three times during every vocalization. The addition of 0 cents shifts (null shifts) has two

437     advantages. First, they allowed us to represent responses to pitch shifts not just as

438     deviations from a pre-shift baseline pitch as in previous studies (Bauer and Larson, 2003;

439     Larson et al., 2007; Liu et al., 2012; Liu and Larson, 2007), but also as deviations from

440     "responses" to a null shifts. in this way, a constant pitch drift common to pitch contours in all

441     conditions cannot affect estimations of response direction and magnitude. Second, the

442     presence of null shifts means it was not predictable for participants how many pitch shifts

443     would occur within one vocalization, thus avoiding any anticipation effects. The shifts were

444     separated from each other and from speech onset by a jittered interval of 600 to 800ms. The

445     pitch shifts were randomized within each experimental block in such a way that each set of

446     two consecutive trials contained all four shifts and two null shifts. An experimental block

447     consisted of 80 trials and thus of 240 shifts, including 40 shifts of each perturbation type as

448     well as 80 null shifts. Before each block, participants produced 10 practice vocalizations to

449     get acquainted with the task and the sound of their voice played via the headphones. After

450     each experimental block, participants got a short break during which they changed the

451     headphones. The order of headphones was counterbalanced across all participants.

452

453     ***3. Equipment***

454     Three pairs of headphones were used: the BeyerDynamic DT 770 Pro (hereafter BD),

455     the custom-built headphones (pair nr. 3, hereafter CB), and the Etymotic ER-3C insert

456     earphones (hereafter the ER). Speech was recorded with a head-mounted microphone (DPA

457     4088-B) positioned at about 2cm from the participant's mouth. The microphone was

458     connected to a Xenyx 802 audio mixer, which sent the signal to an Eventide Eclipse multi-

459     effects processor, which generated the pitch manipulations. The pitch manipulations were

460     controlled via MIDI by a custom PureData (Puckette, 1996) program written by the first

461     author. The output signal from the multi-effects processor was sent, via a different channel

462     on the Xenyx 802 audio mixer, to an Aphex HeadPod 4 headphones amplifier, which

463     connected to the headphones. At the same time, both the microphone signal and the

464     manipulated audio signal were sent to a MicroBook IIc audio interface connected to the

465     laptop in order to store them for offline analysis. All signals were stored at a 44.1kHz

466     sampling rate.

467    In accordance with previous studies, the volume of the auditory feedback was set 10dB

468    above the signal picked up by the microphone (Behroozmand et al., 2014; Hawco et al.,

469    2009; Liu et al., 2011). Any volume differences between headphones were compensated for

470    by adjusting the output gain on the Eclipse Eventide processor. The output gain values used

471    were -16dB, -5dB and -1dB for the CB, the ER and the BD, respectively. These were

472    determined beforehand during a session in which the output volume of each headphone

473    pair was measured with an oscilloscope. The output gain on the Eclipse Eventide processor

474    was adjusted such that all headphones would show a 10dB increase compared to the input

475    volume at the microphone. The delay between microphone input and the auditory feedback

476    output was on average 14.3ms (SD = 5.3ms).

477

478    *4. Analysis*

479    All data and analysis scripts are publicly available via https://osf.io/vm84u/. The data

480    from three participants was not further analysed, because the ER insert earphones did not fit

481    well and thus could have led to a different feedback volume compared to the other

482    headphones. For one of these three participants, the ER earphones fell out during the

483    experiment. The other two participants reported after the experiment that they felt like the

484    earphones were about to fall out, and had difficulty fitting the earphones before the

485    experiment. For the remaining participants, sometimes vocalization was too soft or initially

486    too soft to trigger the pitch shifts in time. This sometimes led to mistiming of the pitch shifts.

487    As long as the pitch shifts were applied during vocalization with ample time of vocalization

488    around (200ms before and 700ms after shift onset), the data was included in the analysis.

489    A pitch estimation algorithm based on autocorrelation in Praat (Boersma and Weenink,

490    2017) was used to estimate pitch in Hertz in every vocalization with a 1ms resolution. The

491    resulting pitch contours were exported to Matlab (R2016b). From every perturbation's pitch

492    contour (including null perturbations), epochs were extracted from 200ms preceding to

493    700ms following the perturbation onset. Pitch was converted to the cents scale as follows:

494

495    $$pitch_{cents} = 1200 \cdot \log_2\left(\frac{pitch_{Hz}}{baseline_{Hz}}\right)$$

496

497        Here, $baseline_{Hz}$ is the mean pitch in Hertz over the 100ms preceding the perturbation

498    onset. The pitch contours for all epochs were visually inspected for pitch estimation errors.

499    As a result of visual inspection, epochs with sharp discontinuities or unusually high variability

500    were discarded. Epochs where more than 10% of the pitch contour was undefined (due to a

501    pitch estimation failure) were discarded as well. On average, about 77.3 epochs (i.e., about

502    10.7 % of the maximum of 720 epochs) were discarded per participant. This includes epochs

503    that displayed pitch tracking errors as well as epochs containing a mistimed pitch shift. The

504    maximal number of epochs discarded for a single participant was 264, with only four

505    participants having more than 200 discarded epochs. Undefined stretches in remaining

506    epochs' pitch contour were linearly interpolated from neighbouring samples.

507        For each participant, headphones, and perturbation condition, the average pitch

508    response contour was calculated by averaging across epochs, as in previous studies with this

509    paradigm (Bauer and Larson, 2003; Larson et al., 2007; Liu et al., 2012; Liu and Larson, 2007).

510    For each participant, only conditions (i.e., a specific headphones by perturbation

511    combination) in which there were at least 20 epochs were included in further analysis. This

512    resulted in the rejection of the data of two additional participants (they each had no

513    condition with over 20 epochs for two of the three headphones) as well as the rejection of

514    data for three conditions in one participant and one condition in another. The resulting 656

515    average pitch contours were derived from on average 35.2 epochs (ranging from 20 to 40

516    out of maximally 40) in the non-null perturbation conditions and from 70.2 epochs (ranging

517    from 33 to 80 out of maximally 80) for the null perturbation conditions. To ensure that

518    response magnitude estimations were not affected by gradual drifts, difference contours

519    were calculated for each participant by subtracting the average for the null perturbation

520    from the average of the corresponding non-null perturbations. The sign of the difference

521    contours for the upward perturbations was flipped such that positive values indicate

522    opposing responses while negative values indicate following responses.

523        For every pair of headphones and perturbation condition, the compensation response

524    magnitude was estimated as the maximal value after 60ms after the perturbation onset.

525        In addition, we used a response classification method to classify every single epoch as

526    containing either an opposing or a following response. The epochs were classified based on

527    the slope of the pitch contour over the time window of 60 to 260ms after perturbation onset

528    (Franken et al., 2018a). As 60ms is considered the minimal time that is necessary to respond

529    to a pitch shift (Chen et al., 2007; Larson et al., 2001), this is presumably the window

530    containing possible responses to the pitch shift onset but not (yet) responses to the pitch

531    shift offset. If the slope was positive, the response was labelled as an upward response (i.e.,

532    an opposing response for downward perturbations and a following response for upward

533    perturbations). The response classification was run on the difference pitch contours, in order

534    to avoid a bias due to an overall pitch drift that is unrelated to the specific condition's pitch

535    perturbation.

536

537    *5. Statistical inference*

538    In order to assess whether pitch shifts led to a general response, each condition's

539    response contour was compared to the response contour for the null shift with the same

540    headphones. This comparison was carried out using a cluster-based permutation test (Maris

541    and Oostenveld, 2007). For every condition, a t-value was calculated at each time sample,

542    and neighbouring time points that exceeded a value corresponding to an uncorrected p-

543    value of 0.05 were clustered. The summed t-value was calculated per cluster and the largest

544    sum was used as the statistic of interest. The same was done after permuting condition

545    labels randomly, arriving at a permutation distribution against which the original statistic

546    value was tested.

547    All further statistical tests were carried out in R (R Core Team, 2018). Response

548    magnitudes were entered in linear mixed effects models, with headphone type, perturbation

549    magnitude, and perturbation direction as fixed effects (main effects and all pairwise

550    interactions as well as a three-way interaction) and random intercepts across subjects. The

551    factor headphone type was dummy-coded with the BD as the reference level, while the

552    perturbation direction and the perturbation magnitude were contrast coded. If model

553    convergence allowed it, random slopes across subjects for headphone type, perturbation

554    magnitude, and perturbation direction were added as well (but no random slopes for

555    interaction effects). Reported p values are calculated using Satterthwaite's methods for

556    estimating degrees of freedom. The omnibus results shown are a type-III table of variance

557    calculated using the anova() function in R, while the pairwise comparisons are calculated

558    using the 'emmeans' package, with Tukey-adjusted p-values if appropriate.

559    The response classification results as either opposing or following responses were

560    entered in a logistic mixed effects model. Reported p values were calculated using Laplace

561    approximation. Omnibus results were derived from type-III Wald $\chi^2$-tests from the 'Anova()'

27

562    function in the 'car' package (Fox and Weisberg, 2019), while pairwise comparisons were

563    calculated with the 'emmeans' package (Lenth, 2019) as before. All mixed effects modelling

564    was performed using the R packages 'lme4' (Bates et al., 2015) and 'lmerTest' (Kuznetsova et

565    al., 2016).

566

567    **B. Results**

568

569        Figure 5 shows the grand average pitch compensation responses, as a function of

570    headphone type, perturbation direction, and perturbation magnitude. These responses

571    show the difference between responses in each condition and the response to a null shift

572    with the same headphones. As expected, in all conditions the grand average pitch contour

573    shows a compensatory response, which starts around 100ms after the perturbation onset

574    and peaks around 250ms after the perturbation onset. At first sight, there seems to be little

575    difference between the responses to the different headphones. Cluster-based permutation

576    tests revealed that for each condition, the response contour differed from the corresponding

577    pitch contour for a null pitch shift (Table 3).
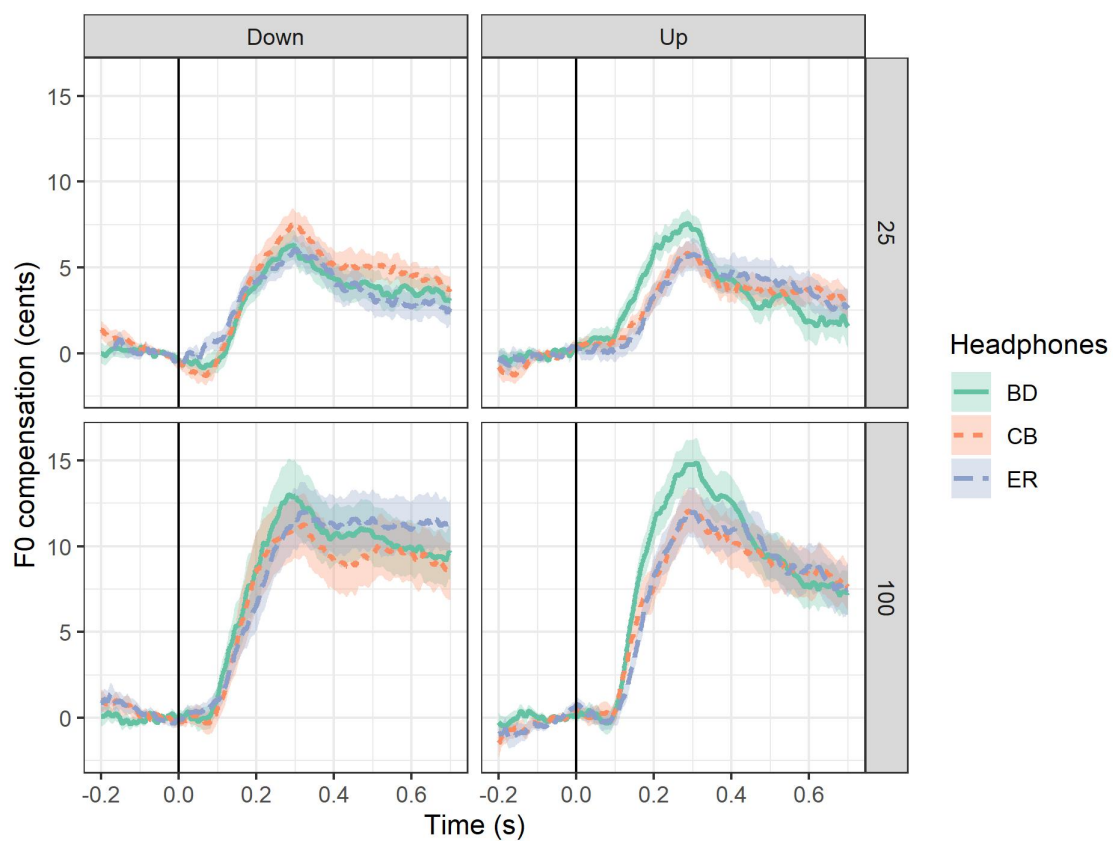
578

579    **Table 3.** Results of the clustered based permutation tests. The reported condition is

580    compared to the corresponding null shift condition in each case. Along with the p-value, the

581    onset time of the largest cluster responsible for the statistical difference is shown.

| Headphones | Pitch shift (cents) | Onset Largest Cluster (ms) | p |
|------------|---------------------|----------------------------|------|
| BD | 25 | 101 | < .001 |
| BD | -25 | 206 | .018 |

| | | | |
|---|---|---|---|
| BD | 100 | 108 | < .001 |
| BD | -100 | 129 | < .001 |
| ER | 25 | 162 | < .001 |
| ER | -25 | 152 | .022 |
| ER | 100 | 119 | < .001 |
| ER | -100 | 135 | .0012 |
| CB | 25 | 130 | < .001 |
| CB | -25 | 159 | .0028 |
| CB | 100 | 97 | < .001 |
| CB | -100 | 152 | < .001 |

582

583



584

585     **Fig. 5.** Grand average pitch compensation contours as a function of Headphones,

586     Perturbation Magnitude, and Perturbation Direction. These contours reflect the difference

587     between the response in each condition and the response to a null shift with the same

588     headphones. The top row displays responses to perturbation with an absolute magnitude of

589     25 cents, the bottom row displays responses to perturbations with an absolute magnitude of

590     100 cents. The left column shows responses to downward pitch shifts (i.e., pitch decreases),

591     while the right column shows responses to upward pitch shifts. The sign of the responses to

592     upward pitch shifts was flipped, so positive values indicate an opposing response. Shaded

593     areas around the contours indicate the standard error of the mean. The vertical black lines

594     indicate the perturbation onset.

595

596         *1. Response Magnitude*

597         The results of a linear mixed effects model of the response magnitude estimates,

598     reported in Table 4, indicate a significant effect of Perturbation Magnitude, showing that

599     responses to 100 cents perturbations were larger than to 25 cents perturbations (contrast =

600     8.36, SE = 0.67, $t(451)$ = 12.46, $p$ < .001). Contrary to our expectations, the response

601     magnitude did not vary as a function of headphone type. This suggests that the amount of

602     passive sound attenuation associated with the different headphones did not affect response

603     magnitude. Response magnitude was also not affected by perturbation direction, or any of

604     the 2-way or 3-way interactions between the three factors. The results are visualized in

605     Figure 6.

606

607     **Table 4.** Omnibus fixed effects on the overall response magnitude. The factors pertMag and

608     pertDir refer to perturbation magnitude and perturbation direction respectively. Colons

30

609    indicate interaction terms (e.g., Headphones:pertMag refers to the two-way interaction

610    between headphone type and perturbation magnitude).

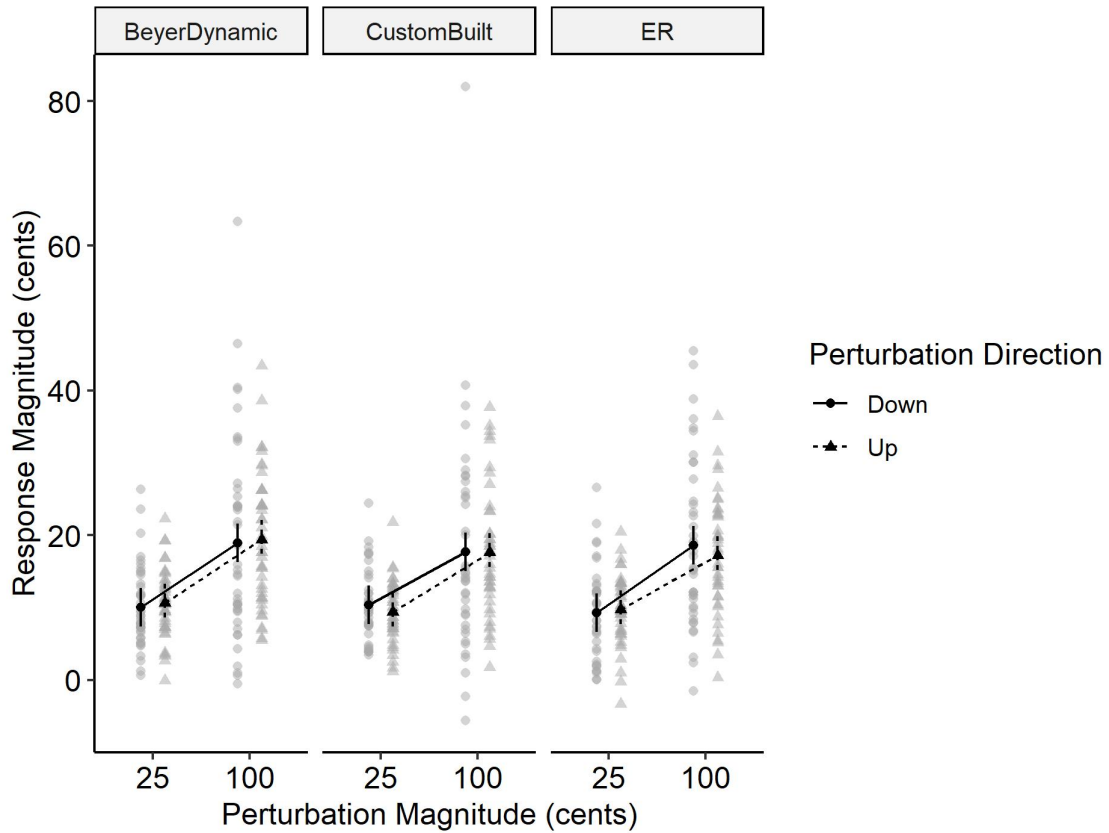|  | SS | df | F | p |
|---|---|---|---|---|
| Headphones | 117.26 | 2, 451 | 1.03 | .36 |
| pertMag | 8797.00 | 1, 451 | 155.13 | **<.001*** |
| pertDir | 2.93 | 1, 451 | 0.052 | .82 |
| Headphones:pertMag | 23.66 | 2, 451 | 0.21 | .81 |
| Headphones:pertDir | 31.04 | 2, 451 | 0.27 | .76 |
| pertMag:pertDir | 4.54 | 1, 451 | 0.080 | .78 |
| Headphones:pertMag:pertDir | 39.75 | 2, 451 | 0.35 | .70 |

611

612

613

**Fig. 6.** Response magnitude as a function of Perturbation Magnitude, Perturbation Direction, and Headphones. In grey, the data for individual subjects is plotted. In black, the fitted values from the mixed effects model are plotted. The error bars indicate 95% confidence intervals.

In a second analysis, the response magnitude was also quantified using only the epochs classified as having opposing responses. Again, only the perturbation magnitude affected response magnitude (contrast = 7.45, $t(412) = 11.31$, $p < .001$). None of the other main effects or interactions yielded significant results.
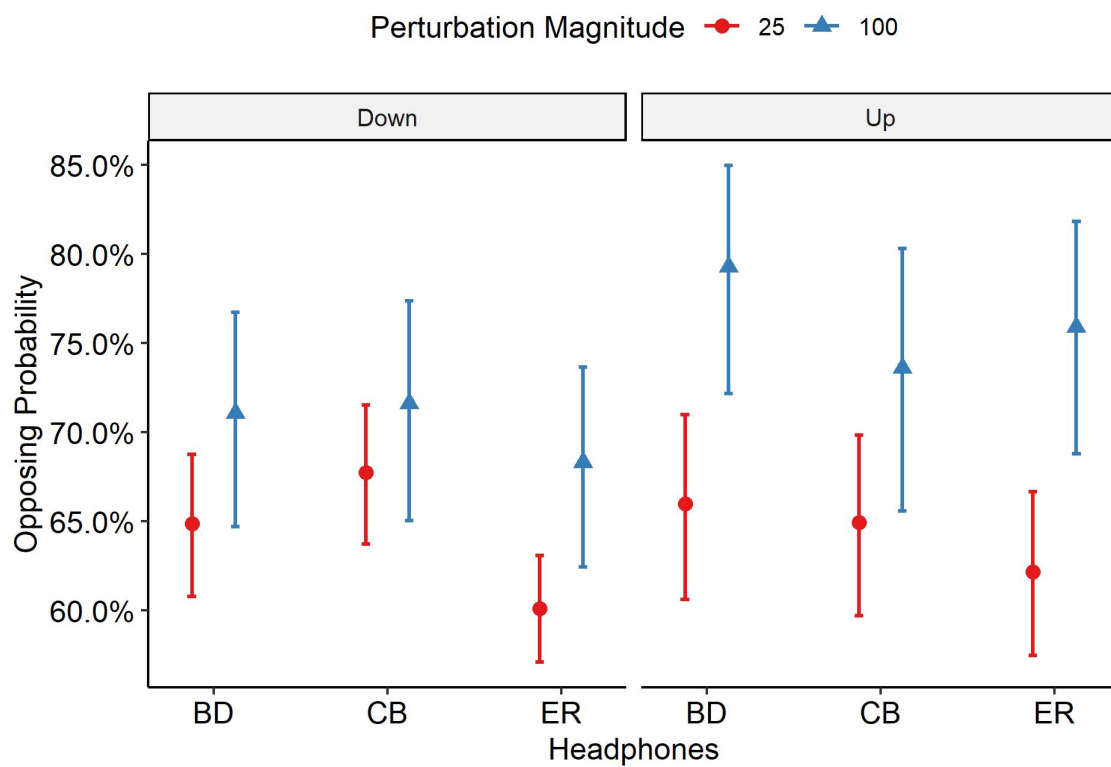
*2. Proportion of opposing responses*

626    Next, epochs were classified as containing either an opposing or a following response.

627    Out of a total of 19857 analysed epochs across all participants and conditions (the null shift

628    excluded), 13377 (about 67%) were classified as opposing and 6480 as following. The

629    probability of an opposing response ('opposing probability') was modelled as a function of

630    the perturbation magnitude, perturbation direction, and type of headphones in a logistic

631    mixed effects model. The results are visualized in Figure 7 and the omnibus effects are

632    shown in Table 5. The results suggest a main effect of headphones ($\chi^2(2) = 7.39$, $p = .025$), a

633    main effect of perturbation magnitude ($\chi^2(1) = 35.45$, $p < .001$), a marginally significant main

634    effect of perturbation direction ($\chi^2(1) = 3.12$, $p = .077$), as well as significant two-way

635    interactions between these factors. Closer examination of the interaction between

636    headphone type and perturbation magnitude suggested that for 25 cents perturbations, the

637    ER led to a lower opposing probability compared to the BD (est. = 0.18, $z = 2.49$, $p = .034$) as

638    well as compared to the CB (est. = 0.23, $z = 2.84$, $p = .013$). For the 100 cents perturbations,

639    there were no significant pairwise contrasts. A similar pattern is visible for the interaction

640    between headphone type and perturbation direction. For downward perturbations, the ER

641    show a lower opposing probability compared to CB (est. = 0.24, $z = 3.01$, $p = .0074$) and a

642    trend towards a lower opposing probability when compared to BD (est. = 0.17, $z = 2.21$, $p =$

643    .070). For upward perturbations, there are no significant contrasts, although there is a trend

644    of BD showing higher opposing probability compared to either CB (est. = 0.18, $z = 2.23$, $p =$

645    .066) or ER (est. = 0.18, $z = 2.32$, $p = .053$).

646

647    **Table 5.** Omnibus (Type-III Wald $\chi^2$ tests) results for opposing probability.

| | $\chi^2$ | df | p |
|---|---|---|---|
| | | | |

| | | | |
|---|---|---|---|
| (Intercept) | 124.50 | 1 | **< .001*** |
| Headphones | 7.39 | 2 | **.025*** |
| pertMag | 35.45 | 1 | **< .001*** |
| pertDir | 3.12 | 1 | .077 |
| Headphones:pertMag | 7.92 | 2 | **.019*** |
| Headphones:pertDir | 12.82 | 2 | **.0016*** |
| pertMag:pertDir | 20.94 | 1 | **< .001*** |
| Headphones:pertMag:pertDir | 1.08 | 2 | .58 |

648
649



650

651 **Fig. 7.** The probability of opposing responses as a function of Perturbation Magnitude,

652 Perturbation Direction, and Headphones. The error bars reflect the 95% confidence intervals

653 of the model's fixed effect estimates.

654

655       In addition, the effect of perturbation magnitude interacted with perturbation

656 direction, suggesting that the difference in probability of opposing responses for 25 cents

657 and 100 cents perturbation was larger for upward (est. = 0.58, $z$ = 7.29, $p$ < .001) than for

658 downward perturbations (est. = 0.28, $z$ = 3.50, $p$ < .001). The three-way interaction between

659 headphone type, perturbation magnitude, and perturbation direction was not significant.

660

661       **IV. GENERAL DISCUSSION**

662       The current study investigated speakers' responses to pitch perturbations with three

663 different headphones varying in amount of passive sound attenuation. Our main research

664 question was whether there are sound attenuation-related differences in the responses to

665 pitch shifts. If responses to pitch shifts are driven by a comparison between an internal pitch

666 target and perceived auditory feedback, increased passive sound attenuation would make

667 the discrepancy between target and feedback more salient, leading to larger responses for

668 more attenuating headphones (like ER). On the other hand, if responses are driven by a

669 comparison between the manipulated signal and the original feedback leaking through the

670 headphones, increased attenuation would make the discrepancy less salient and thus more

671 attenuating headphones should lead to smaller responses. Similarly, if increased attenuation

672 makes it more likely that the manipulated feedback is considered by the speaker as self-

673 generated, we expect more sound attenuation to lead to more opposing responses. In terms

674 of response magnitudes, there were no differences between headphones, in contrast with

675    our hypotheses. This null result suggests that in terms of response magnitude, it does not

676    matter which of the three headphones was used.

677        In terms of the response type (i.e., proportion of following vs. opposing responses), the

678    current analysis revealed no clear overall association between sound attenuation and

679    response type, although there was an interaction of headphone type with both perturbation

680    direction and perturbation magnitude. This pattern of results, although not consistent,

681    suggests that the type of headphones does play a role in this paradigm, although it is hard to

682    pinpoint what role precisely. Specific contrasts showed that response types were only

683    affected by headphone type in some conditions. Although these contrasts may tentatively

684    suggest that higher attenuation (as in ER) is associated with fewer opposing responses, this

685    should be interpreted with caution since a number of the examined contrasts are not strictly

686    significant, and it is unclear how the interactions with perturbation direction and magnitude

687    should be interpreted. In addition, the results of Experiment 1 suggested that the ER show a

688    somewhat different frequency response compared to the circumaural headphones,

689    suggesting that differences that only affect the ER without a difference between CB and BD

690    could be driven by either sound attenuation or the different frequency responses. If,

691    however, future work would corroborate a link between more sound attenuation and less

692    opposing responses, this suggests that passive sound attenuation affects response type but

693    not response magnitude. This is in contrast with our hypotheses. Previous studies showed

694    that response magnitude varies with perturbation magnitude (Chen et al., 2007; Hawco et

695    al., 2009; Liu and Larson, 2007), suggesting that response magnitude can be treated as an

696    index of the conflict introduced by the feedback perturbation. While the causes of following

697    responses are unclear, several authors have proposed that one contributing factor may be

698    that participants treat the feedback signal itself as a referent (Hain et al., 2000; Patel et al.,

699   2014), rather than as self-generated auditory feedback. If so, the current (tentative) results

700   suggest that the sound attenuation of the headphones affects the participants' source

701   monitoring of the auditory signal, but not the magnitude of the feedback mismatch itself.

702          Why would sound attenuation affect response type but not response

703   magnitude? The auditory feedback in the current study was set at 10dBA louder than the

704   signal picked by the microphone, which leads to a quite loud feedback signal. In fact, some

705   of the participants in the current study spontaneously noted that the feedback was very

706   loud. Increasing the loudness of the feedback signal, like many studies do to try and drown

707   out bone-conducted auditory feedback (Behroozmand et al., 2014; Chen et al., 2007; Liu et

708   al., 2011), may create an atypical situation as speakers do not usually hear themselves so

709   loud. If the volume-induced unnaturalness of this feedback signal is exacerbated by the high

710   passive sound attenuation in the ER, it could lead participants to treat the signal as an

711   external referent and to follow the feedback. Just as very large perturbation magnitudes are

712   considered to lead to following because they are unlikely to be self-generated, very high

713   sound attenuation combined with louder than usual auditory feedback could lead to an

714   intrusive auditory signal, which is unlikely to be self-generated. However, the weak statistical

715   evidence in the current study as well as the absence of an effect on the magnitude of the

716   (opposing) responses shows that additional work is necessary to disentangle these

717   possibilities.

718          If the differences between headphones in the current study turn out to be false

719   positives, in line with the absence of a clear overall association between attenuation and

720   response magnitude and types, this would suggest that the differences in passive sound

721   attenuation in the current study did not play a significant role in responses to pitch-shifted

722   feedback, as suggested also by the null effect for the response magnitude. This suggests that

723    compensatory responses across different pitch perturbation studies should be comparable

724    regardless of the headphones used. Although it is possible that varying sound attenuation

725    could have different results if other speech features (e.g., formant values) were

726    manipulated, the current result seems to be in line with a previous study perturbing the first

727    formant (Mitsuya and Purcell, 2016). However, we should be cautious with this conclusion

728    given that there are possible alternative explanations for the absence of an effect of sound

729    attenuation. This result may indicate that the passive sound attenuation of all three

730    headphones in the current study is either good enough to make the manipulated feedback

731    dominant over any leaking original feedback or that participants simply treat only the

732    loudest auditory input as their feedback signal, but it is similarly possible that the differences

733    in attenuation are not large enough, and that therefore non-manipulated auditory feedback

734    plays a similar role in all three headphones. We speculate that this may be caused, in part,

735    by bone-conducted auditory feedback which is potentially not affected by the attenuation

736    properties of the different headphones[3]. Thus, strictly speaking, the current results are not

737    able to distinguish between the dominant hypothesis suggesting that responses are

738    dependent on a comparison between an internal pitch representation and the manipulated

739    feedback, and the alternative hypothesis where responses are dependent on the contrast

740    between the manipulated and non-manipulated auditory feedback.

741        For future studies, an interesting way to address this issue is to mask bone-conducted

742    auditory feedback in order to limit its role and isolate air-conducted feedback that would be

743    affected by headphones' sound attenuation. One could attempt this by playing speech-

744    shaped noise through bone conduction headphones while manipulated feedback is played

745    through normal headphones as in the current experiment. Another way forward may be to

746    take more control over the relative level of the normal and manipulated feedback signals by

747    playing both normal and manipulated feedback through a single pair of headphones while

748    varying the relative levels of both signals. In most altered auditory feedback experiments, it

749    is common to amplify the auditory feedback (as in the current study), in an attempt to make

750    it more salient than potentially conflicting feedback signals. An experiment that compares

751    different relative loudness levels of non-manipulated and manipulated feedback would yield

752    more insight into the role of the relative weighting of conflicting feedback signals.

753        As expected, speakers in the current study show stronger compensation responses to

754    larger perturbations, in line with previous studies (Chen et al., 2007; Hafke, 2008; Hawco et

755    al., 2009; Liu and Larson, 2007), although others have failed to find such an effect (Burnett et

756    al., 1998; Liu et al., 2010b). Interestingly, some studies have shown that this relationship

757    between perturbation magnitude and response magnitude holds only for relatively small

758    perturbations (i.e., up to 200/250 cents), with the compensation response decreasing again

759    for larger responses (Behroozmand et al., 2012; Scheerer et al., 2013), because very large

760    perturbations are unlikely to be considered to be self-generated by the speaker. In addition,

761    we find that 100 cents perturbations in the current study led to a higher probability of

762    opposing responses compared to 25 cents perturbations. Although this does not speak to

763    the influence of leaking non-manipulated auditory feedback, this result is in contrast with

764    some previous studies finding more following responses with larger magnitudes (Burnett et

765    al., 1998; Liu et al., 2010a, 2011, 2010b). It is important to note here that the perturbation

766    magnitudes used in these previous studies were larger than in the present study: for

767    example, most of these studies (Liu et al., 2010a, 2010b) found more following responses to

768    200 or 500 cents perturbations compared to smaller (50 cents and 100 cents) perturbations.

769    As with the findings of differences in response magnitude, based on the findings in the

770    current and in previous studies, we propose that the response type (following vs. opposing)

771  may also show an (inverted) U-shaped relationship with the perturbation magnitude. On the

772  one hand, literature suggests that following responses occur more frequently with very large

773  pitch shifts (200, 500 cents), which may be due to large perturbations being less likely to be

774  recognized as self-produced by the speaker and therefore participants follow it as an

775  external pitch referent. This is in line with suggestions that following responses are observed

776  when the speaker does not consider the presented auditory feedback signal as self-produced

777  speech (Hain et al., 2000; Patel et al., 2014). On the other hand, while very small

778  perturbations (e.g., 25 cents in the current study) are highly likely to be self-generated, they

779  lead to fewer opposing responses as the shifts are less salient compared to slightly larger

780  shifts that are still considered to be self-generated (e.g., 100 cents). In the same vein, a 25

781  cents shift leads to a smaller compensation response than a 100 cents shift because it is less

782  salient, while a 500 cents shift leads to a smaller response compared to 100 cents shifts

783  because it is no longer considered to be self-generated.

784          In addition, it is important to note that most of the previous studies identified

785  the response type (following or opposing) at the average level: epochs were averaged for

786  every condition and participant, and it was identified whether this average response was

787  either following or opposing. Given that recent studies suggest that speakers generally both

788  oppose and follow the pitch shift even within the same condition (Behroozmand et al., 2012;

789  Franken et al., 2018a), the current study classified responses at the single epoch level.

790  Behroozmand et al. (2012) did the same but found no effect of perturbation magnitude on

791  the amount of following/opposing responses (they used perturbations of 50, 100, and 200

792  cents). Given the variability at the single epoch level, it may be the case that correct

793  response classification is harder for 25 cents shifts, as the response magnitudes are smaller

794  and therefore have a lower signal-to-noise ratio. A more precise characterization of the

795     effect of perturbation magnitude on the frequency of opposing and following responses

796     deserves further investigation.

797         With respect to the probability of opposing responses, the current results showed an

798     interaction between perturbation magnitude and perturbation direction, with a stronger

799     effect of perturbation magnitude on opposing probability in the upward shifts compared to

800     the downward shifts. To our knowledge, this is the first paper reporting that the response

801     type may vary as a function of direction, but previous studies have shown directionality

802     effects on response magnitude. The current results seem in contrast with some studies

803     showing larger responses to downward shifts compared to upward shifts (Liu et al., 2011; Liu

804     and Larson, 2007; Sturgeon et al., 2015), while others have found no effect of perturbation

805     direction (Larson et al., 2001, 2008). In the current study we find no effect of perturbation

806     direction on response magnitude, but only on opposing probability. Instead of a

807     directionality effect, this could also suggest an overall bias, in our sample, for downward

808     responses, which would be opposing in response to upward shifts and following in response

809     to downward shifts. We suggest further investigation is needed to investigate the effect of

810     the direction on pitch response types.

811     Overall, the current results suggest that passive sound attenuation has no effect on response

812     magnitudes to unexpected pitch shifts in online auditory feedback. In addition, sound

813     attenuation did also not have a clear effect on the response type. While response type may

814     be affected by multiple factors, including pitch fluctuations before the perturbation onset

815     (Franken et al., 2018a) and properties of the pitch manipulation (Burnett et al., 1998; Liu et

816     al., 2010b), we suggest that in the current study it may be treated as an index of source

817     monitoring. In other words, response type could reflect whether participants attribute the

818     pitch shift to their own production or to an external source. Although we should be cautious

819    to interpret the weak evidence in the current study, we propose that it is important to take

820    into account that non-manipulated auditory feedback may not be completely masked in

821    pitch-shift studies. Responses to pitch-shifted feedback are not only driven by the mismatch

822    between an internal speech target and the manipulated auditory signal, but potentially also

823    by the source attributed to the auditory signal by the speaker. We suggest that it would be

824    interesting for future studies to measure both the response magnitude as well as the

825    response types at an epoch by epoch level. In addition, we have suggested that both

826    response magnitude and response type show an inverted U-shaped relationship with pitch-

827    shift magnitude: for small perturbations, which are likely to be treated as self-generated,

828    larger perturbations lead to larger responses and more opposing responses. Other studies

829    have suggested, in addition, that very large perturbations lead to a decrease in response

830    magnitude and an increase in following responses. Both error-monitoring in speech

831    production as well as source monitoring are functions that have been associated with

832    auditory feedback processing previously (Hain et al., 2000; Korzyukov et al., 2017;

833    Subramaniam et al., 2018). In future studies, it will be important to further investigate the

834    interplay between these two processes.

835

836    **ACKNOWLEDGEMENTS**

842     Foundation. TW was supported by a European Research Council starting grant (ERC-StG-

843     678120 ('RobSpear').

844

845     **FOOTNOTES**

846     1. The headphones were conceived and constructed by the last author.

847     2. See supplementary material at [URL will be inserted by AIP] for details on the

848     construction of the custom-built pair of headphones.

849     3. With current technology, it is difficult to know for sure whether bone-conducted

850     feedback is masked by manipulated auditory feedback or not.

851

852     **REFERENCES**

853 Alain, C. (**2007**). "Breaking the wave: Effects of attention and learning on concurrent sound

854     perception," Hear. Res., **229**, 225–236. doi:10.1016/j.heares.2007.01.011

855 Bates, D., Mächler, M., Bolker, B., and Walker, S. (**2015**). "Fitting Linear Mixed-Effects

856     Models Using lme4," J. Stat. Softw., , doi: 10.18637/jss.v067.i01.

857     doi:10.18637/jss.v067.i01

858 Bauer, J. J., and Larson, C. R. (**2003**). "Audio-vocal responses to repetitive pitch-shift

859     stimulation during a sustained vocalization: Improvements in methodology for the

860     pitch-shifting technique," J. Acoust. Soc. Am., **114**, 1048. doi:10.1121/1.1592161

861 Behroozmand, R., Ibrahim, N., Korzyukov, O., Robin, D. A., and Larson, C. R. (**2014**). "Left-

862     hemisphere activation is associated with enhanced vocal pitch error detection in

863     musicians with absolute pitch," Brain Cogn., **84**, 97–108.

864     doi:10.1016/j.bandc.2013.11.007

865 Behroozmand, R., Korzyukov, O., Sattler, L., and Larson, C. R. (**2012**). "Opposing and

866       following vocal responses to pitch-shifted auditory feedback: evidence for different

867       mechanisms of voice pitch control," J. Acoust. Soc. Am., **132**, 2468–77.

868       doi:10.1121/1.4746984

869  Behroozmand, R., Shebek, R., Hansen, D. R., Oya, H., Robin, D. A., Howard, M. A., and

870       Greenlee, J. D. W. (**2015**). "Sensory-motor networks involved in speech production and

871       motor control: an fMRI study," Neuroimage, **109**, 418–28.

872       doi:10.1016/j.neuroimage.2015.01.040

873  Boersma, P., and Weenink, D. (**2017**). "Praat: doing phonetics by computer [Computer

874       Program].," Retrieved from http://www.praat.org

875  Breebaart, J. (**2017**). "No correlation between headphone frequency response and retail

876       price," J. Acoust. Soc. Am., , doi: 10.1121/1.4984044. doi:10.1121/1.4984044

877  Burnett, T. A., Freedland, M. B., Larson, C. R., and Hain, T. C. (**1998**). "Voice F0 responses to

878       manipulations in pitch feedback," J. Acoust. Soc. Am., **103**, 3153–3161.

879       doi:10.1121/1.423073

880  Cai, S., Ghosh, S. S., Guenther, F. H., and Perkell, J. S. (**2010**). "Adaptive auditory feedback

881       control of the production of formant trajectories in the Mandarin triphthong /iau/ and

882       its pattern of generalization," J. Acoust. Soc. Am., **128**, 2033–48. doi:10.1121/1.3479539

883  Casserly, E. D. (**2011**). "Speaker compensation for local perturbation of fricative acoustic

884       feedback," J. Acoust. Soc. Am., **129**, 2181–2190. doi:10.1121/1.3552883

885  Chen, S. H., Liu, H., Xu, Y., and Larson, C. R. (**2007**). "Voice F[sub 0] responses to pitch-shifted

886       voice feedback during English speech," J. Acoust. Soc. Am., **121**, 1157.

887       doi:10.1121/1.2404624

888  Darwin, C. J. (**1997**). "Auditory grouping," Trends Cogn. Sci., **1**, 327. Retrieved from

889       http://web.mit.edu/hst.723/www/ThemePapers/ASA/Darwin97.pdf

890    Darwin, C. J., Brungart, D. S., and Simpson, B. D. (**2003**). "Effects of fundamental frequency

891        and vocal-tract length changes on attention to one of two simultaneous talkers," J.

892        Acoust. Soc. Am., **114**, 2913. doi:10.1121/1.1616924

893    Dreschler, W. A., Verschuure, H., Ludvigsen, C., and Westermann, S. (**2001**). "ICRA Noises:

894        Artificial Noise Signals with Speech-like Spectral and Temporal Properties for Hearing

895        Instrument Assessment," Int. J. Audiol., **40**, 148–157. doi:10.3109/00206090109073110

896    Elman, J. L. (**1981**). "Effects of frequency-shifted feedback on the pitch of vocal productions,"

897        J. Acoust. Soc. Am., **70**, 45. doi:10.1121/1.386580

898    Epstein, A., Giolas, T. G., and Owens, E. (**1968**). "Familiarity and Intelligibility of Monosyllabic

899        Word Lists," J. Speech Hear. Res., **11**, 435–438. doi:10.1044/jshr.1102.435

900    Flagmeier, S. G., Ray, K. L., Parkinson, A. L., Li, K., Vargas, R., Price, L. R., Laird, A. R., et al.

901        (**2014**). "The neural changes in connectivity of the voice network during voice pitch

902        perturbation," Brain Lang., **132**, 7–13. doi:10.1016/j.bandl.2014.02.001

903    Fox, J., and Weisberg, S. (**2019**). *An {R} Companion to Applied Regression*, Sage, Thousand

904        Oaks, CA, Third edit. Retrieved from

905        https://socialsciences.mcmaster.ca/jfox/Books/Companion/

906    Franken, M. K., Acheson, D. J., McQueen, J. M., Hagoort, P., and Eisner, F. (**2018**). "Opposing

907        and following responses in sensorimotor speech control: Why responses go both ways,"

908        Psychon. Bull. Rev., **25**, 1458–1467. doi:10.3758/s13423-018-1494-x

909    Franken, M. K., Acheson, D. J., McQueen, J. M., Hagoort, P., and Eisner, F. (**2019**).

910        "Consistency influences altered auditory feedback processing," Q. J. Exp. Psychol., **72**,

911        2371–2379. doi:10.1177/1747021819838939

912    Franken, M. K., Eisner, F., Acheson, D. J., McQueen, J. M., Hagoort, P., and Schoffelen, J.

913        (**2018**). "Self-monitoring in the cerebral cortex: Neural responses to small pitch shifts in

914        auditory feedback during speech production," Neuroimage, **179**, 326–336.

915        doi:10.1016/j.neuroimage.2018.06.061

916    Guenther, F. H. (**2016**). *Neural Control of Speech*, The MIT Press, Cambridge, MA, 424 pages.

917    Hafke, H. Z. (**2008**). "Nonconscious control of fundamental voice frequency," J. Acoust. Soc.

918        Am., **123**, 273–8. doi:10.1121/1.2817357

919    Hain, T. C., Burnett, T. A., Kiran, S., Larson, C. R., Singh, S., and Kenney, M. K. (**2000**).

920        "Instructing subjects to make a voluntary response reveals the presence of two

921        components to the audio-vocal reflex," Exp. Brain Res., **130**, 133–141.

922        doi:10.1007/s002219900237

923    Hawco, C. S., Jones, J. A., Ferretti, T. R., and Keough, D. (**2009**). "ERP correlates of online

924        monitoring of auditory feedback during vocalization," Psychophysiology, **46**, 1216–

925        1225. doi:10.1111/j.1469-8986.2009.00875.x

926    Houde, J. F., and Jordan, M. I. (**1998**). "Sensorimotor adaptation in speech production,"

927        Science (80-. )., **279**, 1213–1216. Retrieved from

928        http://www.ncbi.nlm.nih.gov/pubmed/9469813

929    Houde, J. F., and Nagarajan, S. S. (**2011**). "Speech production as state feedback control,"

930        Front. Hum. Neurosci., , doi: 10.3389/fnhum.2011.00082.

931        doi:10.3389/fnhum.2011.00082

932    Johnston, J. D. (**1988**). "Transform Coding of Audio Signals Using Perceptual Noise Criteria,"

933        IEEE J. Sel. areas Commun., Retrieved from

934        https://pdfs.semanticscholar.org/c147/3395b8cc9714930507bde86c89d9d931f6ea.pdf.

935        Retrieved from

936        https://pdfs.semanticscholar.org/c147/3395b8cc9714930507bde86c89d9d931f6ea.pdf

937    Jones, J. A., and Munhall, K. G. (**2000**). "Perceptual calibration of F0 production: Evidence

938     from feedback perturbation," J. Acoust. Soc. Am., **108**, 1246. doi:10.1121/1.1288414

939     Keough, D., and Jones, J. A. (**2009**). "The sensitivity of auditory-motor representations to

940         subtle changes in auditory feedback while singing," J. Acoust. Soc. Am., **126**, 837–46.

941         doi:10.1121/1.3158600

942     Korzyukov, O., Bronder, A., Lee, Y., Patel, S., and Larson, C. R. (**2017**). "Bioelectrical brain

943         effects of one's own voice identification in pitch of voice auditory feedback,"

944         Neuropsychologia, **101**, 106–114. doi:10.1016/j.neuropsychologia.2017.04.035

945     Kuznetsova, A., Brockhoff, P. B., and Christensen, R. H. B. (**2016**). "lmerTest: Tests in Linear

946         Mixed Effects Models R Package version 20-33.," Retrieved from https://cran.r-

947         project.org/package=lmerTest

948     Lametti, D. R., Nasir, S. M., and Ostry, D. J. (**2012**). "Sensory Preference in Speech Production

949         Revealed by Simultaneous Alteration of Auditory and Somatosensory Feedback," J.

950         Neurosci., **32**, 9351–9358. doi:10.1523/JNEUROSCI.0404-12.2012

951     Lametti, D. R., Rochet-Capellan, A., Neufeld, E., Shiller, D. M., and Ostry, D. J. (**2014**).

952         "Plasticity in the Human Speech Motor System Drives Changes in Speech Perception," J.

953         Neurosci., **34**, 10339–10346. doi:10.1523/JNEUROSCI.0108-14.2014

954     Larson, C. R., Altman, K. W., Liu, H., and Hain, T. C. (**2008**). "Interactions between auditory

955         and somatosensory feedback for voice F0 control," Exp. brain Res., **187**, 613–21.

956         doi:10.1007/s00221-008-1330-z

957     Larson, C. R., Burnett, T. A., Bauer, J. J., Kiran, S., and Hain, T. C. (**2001**). "Comparison of voice

958         F[sub 0] responses to pitch-shift onset and offset conditions," J. Acoust. Soc. Am., **110**,

959         2845. doi:10.1121/1.1417527

960     Larson, C. R., Sun, J., and Hain, T. C. (**2007**). "Effects of simultaneous perturbations of voice

961         pitch and loudness feedback on voice F[sub 0] and amplitude control," J. Acoust. Soc.

962          Am., **121**, 2862. doi:10.1121/1.2715657

963    Lenth, R. (**2019**). "emmeans: Estimated Marginal Means, aka Least-Squares Means.,"

964          Retrieved from https://cran.r-project.org/package=emmeans

965    Li, W., Chen, Z., Liu, P., Zhang, B., Huang, D., and Liu, H. (**2013**). "Neurophysiological evidence

966          of differential mechanisms involved in producing opposing and following responses to

967          altered auditory feedback," Clin. Neurophysiol., **124**, 2161–71.

968          doi:10.1016/j.clinph.2013.04.340

969    Lind, A., Hall, L., Breidegard, B., Balkenius, C., and Johansson, P. (**2014**). "Auditory feedback

970          of one's own voice is used for high-level semantic monitoring: the 'self-comprehension'

971          hypothesis," Front. Hum. Neurosci., , doi: 10.3389/fnhum.2014.00166.

972          doi:10.3389/fnhum.2014.00166

973    Liu, H., and Larson, C. R. (**2007**). "Effects of perturbation magnitude and voice F0 level on the

974          pitch-shift reflex," J. Acoust. Soc. Am., **122**, 3671–7. doi:10.1121/1.2800254

975    Liu, H., Meshman, M., Behroozmand, R., and Larson, C. R. (**2011**). "Differential effects of

976          perturbation direction and magnitude on the neural processing of voice pitch

977          feedback," Clin. Neurophysiol., **122**, 951–7. doi:10.1016/j.clinph.2010.08.010

978    Liu, H., Wang, E. Q., Chen, Z., Liu, P., Larson, C. R., and Huang, D. (**2010**). "Effect of tonal

979          native language on voice fundamental frequency responses to pitch feedback

980          perturbations during sustained vocalizations," J. Acoust. Soc. Am., **128**, 3739–3746.

981          doi:10.1121/1.3500675

982    Liu, H., Wang, E. Q., Metman, L. V., and Larson, C. R. (**2012**). "Vocal responses to

983          perturbations in voice auditory feedback in individuals with Parkinson's disease," PLoS

984          One, **7**, e33629. doi:10.1371/journal.pone.0033629

985    Liu, P., Chen, Z., Larson, C. R., Huang, D., and Liu, H. (**2010**). "Auditory feedback control of

986        voice fundamental frequency in school children," J. Acoust. Soc. Am., **128**, 1306.

987        doi:10.1121/1.3467773

988    Maris, E., and Oostenveld, R. (**2007**). "Nonparametric statistical testing of EEG- and MEG-

989        data," J. Neurosci. Methods, **164**, 177–190. doi:DOI 10.1016/j.jneumeth.2007.03.024

990    Mitsuya, T., and Purcell, D. W. (**2016**). "Occlusion effect on compensatory formant

991        production and voice amplitude in response to real-time perturbation," Artic. J. Acoust.

992        Soc. Am., , doi: 10.1121/1.4968539. doi:10.1121/1.4968539

993    Parrell, B., Agnew, Z., Nagarajan, S., Houde, J., and Ivry, R. B. (**2017**). "Impaired Feedforward

994        Control and Enhanced Feedback Control of Speech in Patients with Cerebellar

995        Degeneration," J. Neurosci., **37**, 9249–9258. doi:10.1523/JNEUROSCI.3363-16.2017

996    Patel, S., Nishimura, C., Lodhavia, A., Korzyukov, O., Parkinson, A., Robin, D. a, and Larson, C.

997        R. (**2014**). "Understanding the mechanisms underlying voluntary responses to pitch-

998        shifted auditory feedback," J. Acoust. Soc. Am., **135**, 3036. doi:10.1121/1.4870490

999    Pörschmann, C. (**2000**). "Influences of Bone Conduction and Air Conduction on the Sound of

1000       One's Own Voice," Acust. Acta Acust., **86**, 1038–1045. Retrieved from

1001       www.ingentaconnect.com/content/dav/aaua/2000/00000086/00000006/art00018

1002    Puckette, M. (**1996**). "Pure data," Proc. Int. Comput. Music Conf., International Computer

1003       Music Association, San Francisco, CA, 269–272.

1004    Purcell, D. W., and Munhall, K. G. (**2006**). "Adaptive control of vowel formant frequency:

1005       Evidence from real-time formant manipulation," J. Acoust. Soc. Am., **120**, 966–977.

1006       doi:10.1121/1.2217714

1007    R Core Team (**2018**). "R: a language and environment for statistical computing.," Retrieved

1008       from http://www.r-project.org

1009    Scheerer, N. E., Behich, J., Liu, H., and Jones, J. A. (**2013**). "ERP correlates of the magnitude of

1010    pitch errors detected in the human voice," Neuroscience, **240**, 176–85.

1011    doi:10.1016/j.neuroscience.2013.02.054

1012    Schuerman, W. L., Nagarajan, S., McQueen, J. M., and Houde, J. (**2017**). "Sensorimotor

1013    adaptation affects perceptual compensation for coarticulation," J. Acoust. Soc. Am.,

1014    **141**, 2693–2704. doi:10.1121/1.4979791

1015    Shiller, D. M., Sato, M., Gracco, V. L., and Baum, S. R. (**2009**). "Perceptual recalibration of

1016    speech sounds following speech motor learning," J Acoust Soc Am, **125**, 1103–1113.

1017    doi:10.1121/1.3058638

1018    Sturgeon, B. A., Hubbard, R. J., Schmidt, S. A., and Loucks, T. M. (**2015**). "High F0 and

1019    musicianship make a difference: Pitch-shift responses across the vocal range," J. Phon.,

1020    **51**, 70–81. doi:10.1016/j.wocn.2014.12.001

1021    Subramaniam, K., Kothare, H., Mizuiri, D., Nagarajan, S. S., and Houde, J. F. (**2018**). "Reality

1022    Monitoring and Feedback Control of Speech Production Are Related Through Self-

1023    Agency," Front. Hum. Neurosci., **12**, 82. doi:10.3389/fnhum.2018.00082

1024    Warren, R. M., Riener, K. R., Bashford, J. A., and Brubaker, B. S. (**1995**). *Spectral redundancy:*

1025    *Intelligibility of sentences heard through narrow spectral slits* Percept. Psychophys.,Vol.

1026    57, 175–182 pages. Retrieved from

1027    https://link.springer.com/content/pdf/10.3758/BF03206503.pdf

1028

1029