

1 **Altered auditory feedback induces coupled changes in formant frequencies during speech**  
2 **production**

3 Running title: Formant changes induced by altered auditory feedback

4

5 Ding-lan Tang<sup>1</sup>, Daniel R. Lametti<sup>1,2</sup>, Kate E. Watkins<sup>1</sup>

6 <sup>1</sup>. *Wellcome Centre for Integrative Neuroimaging, Department of Experimental Psychology,*  
7 *University of Oxford*

8 <sup>2</sup>. *Department of Psychology, Acadia University, Wolfville, Nova Scotia, Canada*

9  
10

11

12 **Abstract**

13 Speaking is one of the most complicated motor behaviours, involving a large number of  
14 articulatory muscles which can move independently to command precise changes in speech  
15 acoustics. Here, we used real-time manipulations of speech feedback to test whether the  
16 acoustics of speech production (e.g. the formants) reflect independently controlled  
17 articulatory movements or combinations of movements. During repetitive productions of  
18 “head, bed, dead”, either the first (F1) or the second formant (F2) of vowels was shifted and  
19 fed back to participants. We then examined whether changes in production in response to  
20 these alterations occurred for only the perturbed formant or both formants. In Experiment  
21 1, our results showed that participants who received increased F1 feedback significantly  
22 decreased their F1 productions in compensation, but also significantly increased the  
23 frequency of their F2 productions. The combined F1-F2 change moved the utterances closer  
24 to a known pattern of speech production (i.e. the vowel category “hid, bid, did”). In  
25 Experiment 2, we further showed that a downshift in frequency of F2 feedback also induced  
26 significant compensatory changes in both the perturbed (F2) and the unperturbed formant  
27 (F1) that were in opposite directions. Taken together, the results demonstrate that a shift in  
28 auditory feedback of a single formant drives combined changes in related formants. The  
29 results suggest that, although formants can be controlled independently, the speech motor  
30 system may favour a strategy in which changes in formant production are coupled to  
31 maintain speech production within specific regions of the vowel space corresponding to  
32 existing speech-sound categories.

33

34 Key words: speech perception; adaptation; speech motor learning; feedback perturbation;  
35 auditory error

36

37 **New & Noteworthy**

38 Findings from previous studies examining responses to altered auditory feedback are  
39 inconsistent with respect to the changes speakers make to their production. Speakers can  
40 compensate by specifically altering their production to offset the acoustic error in feedback.  
41 Alternatively, they may compensate by changing their speech production more globally to  
42 produce a speech sound closer to an existing category in their repertoire. Our study shows  
43 support for the latter strategy.

44

45

46

## 47 **Introduction**

48 Speaking is one of the most complicated motor behaviours. Fluent speech production  
49 requires fine coordination of many muscles and integration with simultaneous auditory and  
50 somatosensory feedback. How this coordination of movement is achieved given the large  
51 number of muscles involved and degrees of freedom of movement and how these map on  
52 to sensory targets are largely unresolved problems in the field of speech motor control. On  
53 the sensory side, it is clear that auditory feedback is particularly important in speech  
54 development; pre-lingually deaf children have difficulties with speech articulation (Svirsky et  
55 al., 2004). It is also important for the maintenance of speech intelligibility, which can  
56 deteriorate in post-lingually deaf adults (Waldstein, 1990). Such individuals can rely on  
57 somatosensory feedback from orofacial muscles to achieve normal speech articulation  
58 (Nasir and Ostry 2008). Somatosensory feedback is also relied on by some individuals with  
59 normal hearing to monitor their speech production (Lametti et al., 2012; Tremblay et al.,  
60 2003).

61  
62 The relationship between the movements used to articulate speech sounds and the sensory  
63 consequences is learned during speech development and used to build internal models that  
64 encode these sensorimotor mappings and guide speech production. Internal models allow  
65 the brain to predict the sensory consequences of actions, perform online correction and  
66 rapidly detect errors. The exact nature of these representations is unknown. That is, are  
67 speech sound targets (e.g., phonemic categories) represented by a combination of  
68 movements used to achieve the sensory goal—or are the targets of speech the  
69 consequence of a series of independently controlled articulatory movements?  
70

71 Internal models need to be updated and trained to adapt for changes in the environment or  
72 in our bodies that affect speech perception and production—for example, when imitating  
73 the accent of an interlocutor or coping with the changes in shape of the vocal tract during  
74 adolescence. It is possible to explore these models experimentally by manipulating sensory  
75 feedback in real-time and measuring the ways in which speakers adapt their speech  
76 production (Houde & Jordan 1998; MacDonald et al., 2011; Lametti et al., 2018a).  
77

78 The basic speech adaptation paradigm involves altering frequencies that characterize vowel  
79 sounds, feeding these back to the speaker in near real time (e.g., with less than a 40 ms  
80 delay), and measuring the compensatory changes in the vowel sounds the speaker  
81 subsequently produces (Houde and Jordan 1998, 2002; Purcell and Munhall, 2006a, b;  
82 Mollaei et al., 2013). In speech, vowel sounds are defined by vocal tract resonances known  
83 as formants. Formants are distinctive frequency components evident in the acoustic signal  
84 reflecting the size and shape of different chambers in the vocal tract. The first two formants  
85 (F1 and F2) in a vowel sound contain most of the acoustical energy and are sufficient to  
86 disambiguate vowel sounds. By examining the relationship between the changes in the first  
87 two formants of vowels produced when speech feedback is altered, we can gain insight into  
88 how speech movements are controlled. For example, if the frequency of F1 for the vowel in  
89 “head” is increased, it sounds closer to the F1 frequency of the vowel in “had”. Participants  
90 can compensate by lowering the frequency of F1 in their subsequent productions, which  
91 would counteract the acoustic change applied and result in an F1 that with the frequency  
92 shift applied and fed back sounds closer to that of the intended vowel in “head”.  
93 Alternatively, participants could change the frequency of both F1 and F2 in their

94 productions; a combination of decreasing F1 frequency whilst simultaneously increasing F2  
95 could move their production closer to that of the vowel in “hid”, a familiar region of the  
96 vowel space (see Figure 1). Determining whether adaptive speech movements are  
97 characterized by specific changes along separate dimensions, (e.g., F1, which corresponds to  
98 tongue height, and F2, which corresponds to front-back position of the tongue body;  
99 Ladefoged, 2001) or are controlled more globally by the production of vocal tract shapes  
100 that involve both dimensions, would inform our understanding of the nature of  
101 representations in models of speech motor control.  
102

103 The current literature provides evidence for both specific changes in F1 during adaptation to  
104 F1 alterations and more global strategies involving changes in both F1 and F2. Previous work  
105 involving adaptation to gradual shifts in F1 showed compensation in F1 production but no  
106 significant change in F2 (Villacorta et al. 2007). In contrast, another study found changes in  
107 both F1 and F2 (MacDonald et al., 2011) when feedback in F1 was changed abruptly. In this  
108 case, however, the direction of change in F2 was the same for both upwards and  
109 downwards shifts of F1 and occurred more gradually than the change in F1. This suggests  
110 that F1 compensation, which is directional and compensatory for the applied F1 shift, is  
111 independent of the F2 change. On the other hand, there is also evidence to suggest that  
112 formants are not modified independently. For example, reliable F2 changes in response to  
113 shifts in F1 were observed in the opposite direction to the compensatory changes in F1 for  
114 utterances containing the vowel /ə/, such as “bed” and “head” (Rochet-Capellan and Ostry,  
115 2011; Lametti et al., 2018b).  
116

117 The inconsistency in results of previous studies could be due to the relatively small sample  
118 sizes and the large variability among participants. A sample size of 12 (or less) per group is  
119 typical for these kinds of studies yet a large proportion of people (10-40%) exhibit little or  
120 no compensation (Lametti et al., 2012, 2014; Rochet-Capellan et al. 2012). In addition,  
121 meta-analyses of such studies indicate that on average 22% of participants altered  
122 production in a pattern that follows (instead of opposes) the direction of the F1 shift  
123 (MacDonald et al., 2011). These individual differences could be explained by preferences  
124 for reliance on auditory or somatosensory feedback during speech perturbation studies  
125 (Lametti et al., 2012) and further complicates work aiming to explore the relationship  
126 between F1 and F2 production in response to alterations of a single formant (MacDonald et  
127 al., 2011).  
128

129 In the current study, we revisit the idea that speakers counteract errors in formant  
130 production by only changing the altered formant. In our first experiment, we increased the  
131 frequency of F1 and examined compensatory changes in both F1 and F2. Based on previous  
132 findings (see above), we predicted that participants would decrease their F1 to offset the  
133 perceived acoustical error. We also predicted based on studies where the shift was applied  
134 abruptly (e.g. Rochet-Capellan and Ostry 2011; Lametti et al., 2018b) that increases in F2  
135 would accompany the F1 decrease to move speech closer to a known region of the vowel  
136 space. In other words, when participants produced “head, bed, dead” and heard an upward  
137 shift in F1, an increase in F2 concurrent with the compensatory decrease in F1 would place  
138 the new production closer to “hid, bid, did” in vowel space (see Figure 1). In a second  
139 experiment, we decreased the frequency of F2 and examined compensatory changes in  
140 both F2 and F1. We predicted that increases in F2 would be observed to offset the induced  
141 acoustical error, and, as for responses to increased F1, decreases in the unaltered formant

142 (this time F1) would accompany these F2 increases to move speech closer to a known region  
143 of the vowel space (Figure 1). Such a result would support the idea that, although F1 and F2  
144 can be controlled independently (MacDonald et al., 2011), the speech motor system favours  
145 a strategy that moves formant production to specific regions of the vowel space.

146  
147

## 148 **Methods**

### 149 **Participants and Apparatus**

150 Forty native English speakers between the ages of 18 and 40 years, with no reported history  
151 of speech, language, or hearing disorders took part in the study. The Central University  
152 Research Ethics Committee (at the University of Oxford) approved the experimental  
153 protocol and participants gave informed consent. There were twenty participants in  
154 Experiment 1 ( $21.55 \pm 5.27$  years; 18 females) and 20 different participants in Experiment 2  
155 ( $20.15 \pm 3.59$  years, 14 females).

156

### 157 **Apparatus**

158 The experiment was conducted in a quiet room with participants seated in front of a  
159 computer screen, which was used to present words. During the experiment, participants  
160 wore over-ear headphones (Sennheiser, HD 280 Pro) and a head-mounted microphone  
161 (Shure, WH20) that was placed approximately 4-7 cm away from the right corner of the  
162 participant's mouth. A Matlab Mex-based program Audapter (Cai et al., 2008; Tourville et al.,  
163 2013) was used to alter speech and play it back to participants with an unnoticeable delay.

164

### 165 **Procedure**

166 In both experiments, consonant-vowel-consonant words were presented on the computer  
167 screen for 1500 ms, one at a time. The inter-trial-interval was 750 ms. Participants were told  
168 that they would have to read out the words on the screen, and that they would hear their  
169 own voice through headphones.

170

171 First, participants completed a vowel exploration session with normal feedback in which  
172 they produced nine different words ('dead' 'bed' 'head', 'dad' 'bad' 'had', 'did' 'bid' 'hid',  
173 containing three different vowels /æ/, /æ/ and /ɪ/, correspondingly) 15 times each. This  
174 session allowed participants to explore the vowel space. It also provided us with estimates  
175 of F1 and F2 frequencies for each vowel category in each individual. Vowel exploration was  
176 followed by three experimental phases: Baseline, Learning, and Unlearning (Figure 2).  
177 During the baseline phase, three different words "dead, bed, head", containing the same  
178 vowel /æ/, were produced 15 times each by participants with normal feedback (2 min). In  
179 the learning phase (10 min), participants produced "dead, bed, head" 75 times each with  
180 altered auditory feedback designed to induce sensorimotor learning (see Real-Time  
181 Feedback Alteration). In the final phase, Unlearning, "dead, bed, head" were produced 30  
182 times each with normal feedback (4 min). The order of words was randomized within each  
183 phase.

184

185 Three short noise-masked sessions preceded each of these experimental phases. During  
186 these noise-masked sessions, participants produced the words "dead, bed, head" five times  
187 each but, instead of hearing feedback of their speech, they heard masking noise that scaled  
188 with the amplitude of the produced speech signal (i.e. the signal-to-noise ratio was 0 dB).

189 White noise was also added to block participants' ability to hear their own speech as much  
190 as possible.

191

192 The three noise-masked sessions were carried out at the following time-points: 1) Pre-  
193 Baseline 2) Before-Learning, and 3) After-Learning. The aim of these sessions was to  
194 measure whether feedforward changes in speech production associated with sensorimotor  
195 learning would persist in the absence of auditory feedback. Figure 2 shows a schematic of  
196 the Experimental Procedure.

197

### 198 **Real-Time Feedback Alteration**

199 To induce sensorimotor learning in speech, the frequency of either F1 or F2 was altered and  
200 played back to participants with an imperceptible delay. The formant alteration was applied  
201 in mels — a perceptual scale characterizing changes in frequency that are judged by  
202 listeners to correspond to equal changes in pitch (Stevens et al., 1937). The transformation  
203 from Hz to mel used was:

204

$$205 \text{ mel} = 1127.01048 \times \log(1 + \text{Hz}/700) \quad (1)$$

206

207 Participants were randomly assigned to one of two groups. The F1+110 group (Experiment 1)  
208 experienced a 110 mel increase in first formant (F1) production, while the F2-110 group  
209 (Experiment 2) experienced a 110 mel decrease in second formant (F2) production. A  
210 feedback shift magnitude of 110 mel was selected for this study because it is comparable to  
211 the F1 feedback perturbation used previously (Lametti et al., 2018b).

212

213 Both the increase in F1 frequency (F1+110) and the decrease in F2 frequency (F2-110) of the  
214 vowel /ʌ/ in the produced words resulted in the percept of speech with formant frequencies  
215 that were closer to F1 or F2 frequencies respectively of the vowel /æ/ in the vowel space  
216 (Fig. 1). Compensatory responses were expected to be produced by speakers to counteract  
217 the perceived shift, either by reducing F1 (Experiment 1: F1+110 group) or increasing F2  
218 (Experiment 2: F2-110 group) (Fig. 1). We were also interested in the extent to which  
219 speakers altered their production of the unshifted formant in each experiment.

220

### 221 **Auditory analysis**

222 Speech was recorded at 16000 Hz in MATLAB and analysed using custom-written MATLAB  
223 scripts. First, the F1 and F2 frequencies of each vowel were extracted from the data output  
224 of Audapter, which estimates the first two formant frequencies using linear predictive  
225 coding. Then the average formant frequency was calculated from a 30-ms segment at the  
226 centre of each vowel. Utterances that contained production errors or contained formant  
227 frequencies that were greater than three standard deviations from a participant's mean F1  
228 and F2 values in each phase of the experiment were excluded from further analyses. This  
229 resulted in approximately 3% of the data being excluded.

230

### 231 **Statistical analysis**

232 In each participant, the F1 and F2 frequency data for utterances produced during the  
233 learning and unlearning phases of the experiments were normalised by subtracting the  
234 mean and dividing by the standard deviation of the formant frequencies produced during  
235 baseline (with normal feedback) to create a Z value for each datapoint.

236

237 Following Lametti et al. (2018b), we assessed the significance of compensatory changes in  
238 the shifted and unshifted formant frequencies in each experiment using a one-sample t-test  
239 on the average Z values of the last 30 trials in the learning phase (against  $z = 0$ , no change  
240 from baseline). As the directions of the compensatory change to the formant feedback shifts  
241 were predictable, we used directional t-test with a significance level of  $p < .05$ . In the  
242 current study, we also used noise-masked sessions before the baseline, before and after the  
243 learning phase. Data obtained in these sessions allowed us to compare the effect of the  
244 perturbation on speech production under the same conditions when feedback was blocked.  
245 The pre-baseline session was used to normalise the before- and after-learning sessions as  
246 described above (we first confirmed that there was no change in F1 and F2 in the sessions  
247 performed either side of the baseline;  $p > .05$ ). Then average change from before-learning  
248 session in produced F1 and F2 during after-learning session was calculated and compared  
249 with zero using directional t-tests to examine if there was a significant change from before  
250 to after-learning.

251

252 Following MacDonald et al. (2011), we assessed the significance of the return-to-baseline  
253 changes in the shifted and unshifted formant frequencies in the unlearning phase of each  
254 experiment. The average Z values of the last 15 trials in the unlearning phase were  
255 compared using a one-sample t-test (against  $Z = 0$ , i.e., the baseline average).

256

257 For each participant, we also calculated the Euclidean distance between the centre of the  
258 vowel clouds for “hid, bid, did” produced during the vowel exploration phase of the  
259 experiment and for “head, bed, dead” produced at baseline and at the end of (last 30  
260 utterances) the learning phase (Lametti et al., 2018b). Learning-related changes in this  
261 distance were then compared to zero (no change) with one-sample t-tests to examine  
262 whether there was a significant change in the distance induced by F1 (experiment 1) or F2  
263 (experiment 2) feedback perturbation. Distance changes were also compared between  
264 groups (experiment 1 and experiment 2) using a two-tailed t-test.

265

## 266 **Unplanned Analyses**

267 In an exploratory analysis, we compared the direction of compensation in Experiment 1 and  
268 Experiment 2. To do this, we calculated the angle between the vector that represents the  
269 direction of actual compensatory change from baseline and the vector representing the  
270 direction to move from “head, bed, dead” to “hid, bid, did” for each participant. The angles  
271 were then compared to zero (no difference from required angle) with one-sample t-tests to  
272 examine whether there was a significant difference between actual movement direction  
273 and the direction required in each Experiment separately. The angles were also compared  
274 between groups (experiment 1 and experiment 2) using a two-tailed t-test.

275

276

## 277 **Results**

### 278 **Experiment 1: F1 shift of +110 mel**

279

280 Figure 3A shows how F1 and F2 produced by participants, normalized as z-scores, changed  
281 during the sensorimotor learning and unlearning phases in response to a 110 mel increase  
282 in F1. By the end of the learning phase, participants had decreased their F1 on average by

283 30 mel to compensate for the increase in F1 and increased F2 on average by 15 mel (Figure  
284 3B). Directional t-tests showed that both F1 and F2 frequency at the end of the learning  
285 phase (last 30 utterances of learning phase) were significantly different from zero (baseline  
286 average) (F1:  $t(19) = 6.24$ ,  $p < .001$ ,  $d = 1.97$ ; F2:  $t(19) = 2.88$ ,  $p = .005$ ,  $d = 0.91$ ).

287  
288 Compensatory responses of F1 and F2 production were also evaluated under noise-masked  
289 conditions, which were used as a measure of changes from before to after learning in the  
290 same context (i.e. with no feedback). Figure 3C shows the change in normalized F1 and F2  
291 for heavily noise-masked speech from before and after learning. After learning, F1  
292 significantly decreased ( $t(18) = 4.12$ ,  $p < .001$ ,  $d = 0.98$ ) and F2 significantly increased ( $t(18) =$   
293  $3.00$ ,  $p = .004$ ,  $d = 0.69$ ).

294  
295 By the end of unlearning phase, F2 production returned to baseline (one-sample t-test  
296 against baseline average,  $Z = 0$ :  $t(19) = 1.42$ ,  $p = .172$ ), while F1 remained significantly  
297 reduced relative to baseline ( $t(19) = 2.33$ ,  $p = .031$ ) after these 90 trials with normal  
298 feedback (Figure 3D).

299  
300 In sum, in response to an F1 feedback perturbation, both F1 and F2 in production were  
301 observed to change though the magnitude of the change in the shifted formant (F1) was  
302 larger than in the unshifted one (F2). The results provide support for the idea that a single  
303 formant perturbation will induce a combination of changes in the shifted formant as well as  
304 the unshifted formant. Such combined changes could place the new production closer to a  
305 learnt pattern of speech production used to produce another vowel sound (“hid, bid, did” in  
306 current study; Figure 5A). This hypothesis was tested by comparing the magnitude  
307 (Euclidean distance) and direction of the compensatory shifts in F1 and F2 in vowel space  
308 between the pre- and post-learning productions of (“head, bed, dead”) and the production  
309 of “hid, bid, did” during the vowel exploration phase of the experiment before learning had  
310 occurred.

311  
312 As shown in Figure 5B, participants in Experiment 1, responding to an increase in F1  
313 feedback, moved their productions of “head, bed, dead” closer to the vowel space occupied  
314 by “hid, bid, did” by a magnitude of 32 mel on average (one-sample t-test against zero,  $t(19)$   
315  $= 5.44$ ,  $p < .001$ ,  $d = 1.72$ ). The direction of this change was not significantly different from  
316 that predicted (Figure 5C;  $t(19) = 1.45$ ,  $p = .164$ ). The same results were obtained by  
317 analysis of the data from the noise-masked session after learning.

318

319

## 320 **Experiment 2: F2 shift of -110 mels**

321

322 In Experiment 2, we gave participants feedback with a downshift in F2 and measured the  
323 changes made in production of F1 and F2. We predicted that participants would show a  
324 conceptually similar compensatory pattern to that observed in Experiment 1: F2 frequency  
325 would increase, to offset the perceived acoustical error, but F1 would also change (decrease)  
326 to move production changes closer to “hid, bid, did” (see Fig. 1).

327

328 Figure 4A shows how F1 and F2 produced by participants, normalized as z-scores, changed  
329 during the sensorimotor learning and unlearning phases in response to a 110 mel decrease

330 in F2. By the end of the learning phase, participants had increased their F2 on average by 50  
331 mel to compensate for the decrease in F2 and decreased F1 on average by 14 mel (Figure  
332 4B). Directional t-tests showed that both F1 and F2 frequency at the end of the learning  
333 phase (last 30 utterances of learning phase) were significantly different from zero (baseline  
334 average) (F2:  $t(19) = 7.79$ ,  $p < .001$ ,  $d = 2.46$ ; F1:  $t(19) = 2.30$ ,  $p = .017$ ,  $d = 0.73$ ).

335

336 Figure 4C shows the change in normalized F1 and F2 frequency for heavily noise-masked  
337 speech from before and after learning in response to a decrease in F2 feedback. There was a  
338 significant increase in F2 production ( $t(19) = 7.12$ ,  $p < .001$ ,  $d = 1.59$ ), and a significant  
339 reduction in F1 ( $t(19) = 1.99$ ,  $p = .030$ ,  $d = 0.45$ ) from before to after learning.

340

341 By the end of the unlearning phase, F2 frequency returned to baseline (one-sample t-test  
342 against baseline average,  $Z = 0$ :  $t(19) = 1.95$ ,  $p = .066$ ), while F1 frequency remained  
343 significantly reduced relative to baseline ( $t(19) = 2.31$ ,  $p = .032$ ) after these 90 trials with  
344 normal feedback (Fig. 4D).

345

346 The results of Experiment 2 provide further support for the idea that feedback  
347 perturbations induce changes not only in the shifted formant, but also in the unshifted  
348 formant. As in Experiment 1, we evaluated whether these combined changes served to  
349 move the produced speech token closer to another vowel sound. The vector analysis  
350 confirmed a learning-related reduction in the Euclidean distance between “head, bed, dead”  
351 and “hid, bid, did” in F1-F2 space Experiment 2 (Figure 5B). By the end of learning, the  
352 productions of “head, bed, dead” were 36 mel closer to the productions of “hid, bid, did” in  
353 response to the F2 shift-down (one-sample t-test against zero,  $t(19) = 5.96$ ,  $p < .001$ ,  $d =$   
354  $1.88$ ). However, the direction of the change was significantly different from the angle  
355 predicted (Figure 5C;  $t(19) = 6.86$ ,  $p < .001$ ,  $d = 2.17$ ). Thus, although the results of  
356 Experiment 2 met our predictions—a decrease in F2 drove an increase in F2 and a decrease  
357 in F1—changes in production were not as well directed towards the “hid, bid, did” part of  
358 vowel space as they were in Experiment 1 (Figure 5A). The same results were obtained by  
359 analysis of the data from the noise-masked session after learning.

360

### 361 **Comparison of Experiment 1 and 2**

362 We compared the changes in shifted and unshifted formants between the two experiments.  
363 The data were rectified as the directions differed between the two experiments (all data  
364 were made positive) and analysed using a two-tailed t-test. For the shifted formants, the  
365 magnitude of the learning-related change (last 30 utterances of the learning phase) was  
366 significantly greater in the F2-shift experiment (Exp 2) compared with the F1-shift  
367 experiment (Exp 1) ( $t(38) = 2.48$ ,  $p = .018$ ;  $d = 0.79$ , Fig. 3B and 4B). However, for the  
368 unshifted formants, the magnitudes of the shift at the end of the learning phases of  
369 Experiments 1 and 2 did not differ ( $p = .329$ ; Figs. 3B and 4B). In the vector analysis, there  
370 was no significant difference in the magnitude of the learning-related decrease in distance  
371 between Experiments 1 and 2 ( $t(38) = 0.52$ ,  $p = .610$ ; Fig. 5B). However, the differences  
372 between the actual and predicted direction of the changes in each experiment were  
373 significantly different ( $t(38) = 4.385$ ,  $p < .001$ ;  $d = 1.39$ ), with a larger difference from the  
374 predicted direction for the actual direction of changes in Experiment 2 (Fig. 5C).

375



376 Taken together, these findings suggest that the compensatory responses seen in response  
377 to shifts applied to F1 and F2 formants in separate experiments each served to move the  
378 vowel production towards another vowel in the vowel space and an existing set of  
379 articulatory movements associated with that learned production. The significantly larger  
380 amount of compensation in the shifted formant (F2) observed in Experiment 2 (in response  
381 to a shift down in F2) may explain the discrepancy between the actual and predicted  
382 direction required to move from the “head, bed, dead” to the “hid, bid, did” vowel space for  
383 that experiment. Nevertheless, the magnitude of the difference between the two vowel  
384 productions was significantly reduced following learning in both experiments. It is worth  
385 noting that when participants were debriefed after the study, 75% of participants who  
386 experienced the F2 shift indicated they were aware of the changes, while only 20% of the  
387 participants who experienced the shift upward in F1 noticed them.

388  
389

## 390 Discussion

391 The current study investigated how speakers altered their production of the first two  
392 formants of vowels by manipulating the auditory feedback for F1 and F2 in two separate  
393 experiments. Our hypothesis was that, although the movements required to produce  
394 individual formants can be controlled independently to compensate for specific acoustic  
395 changes in feedback, speech movements are controlled in a more global way. Since changes  
396 in both formants are required to move from one vowel category to another, we predicted  
397 that altered auditory feedback of a single formant during speech production would induce  
398 changes not only in the shifted formant, but also in the unshifted formant. Furthermore, we  
399 predicted that these combined changes would result in a pattern of speech production  
400 closer to the frequencies of an existing vowel category.

401

402 Consistent with our first prediction, the results of Experiments 1 and 2 showed that by the  
403 end of the adaptation phase, participants had significantly changed production of both the  
404 shifted and the unshifted formants. Specifically, those who received feedback with an  
405 increase in F1, produced speech with a significant decrease in F1 as well as a significant  
406 increase in F2. Those who received feedback with a decrease in F2, produced speech with a  
407 significant increase in F2 and a significant decrease in F1. The changes relative to speech  
408 prior to the application of the shift were evident both in the last 30 trials of the learning  
409 phase and in the noise-masked blocks that followed the learning phase (see Figs. 3 & 4).

410

411 The results of Experiment 1, in which F1 feedback was increased and participants  
412 compensated by decreasing F1 and increasing F2 in their productions, are consistent with  
413 previous findings using the same task with a different altered feedback system (Lametti et  
414 al., 2018b; Rochet-Capellan and Ostry, 2011). However, the findings of both experiments in  
415 the current study contrast with previous work in which participants showed independent  
416 changes to the perturbed and unperturbed formants when feedback was gradually shifted  
417 in either F1 or F2 (Villacorta et al. 2007). In the latter study, compensatory F1 and F2  
418 changes were observed in the opposite direction to the applied shift in F1 and F2  
419 respectively, as seen in our study. However, F1 did not change during F2 feedback  
420 perturbation, which contrasts with the findings for Experiment 2 in our study. Small F2  
421 changes induced by F1 feedback shift were observed by MacDonald et al (2011) but these  
422 changes were decreases in F2 when F1 feedback was shifted either up or down. In contrast,

423 in our study, F2 production increased when F1 was shifted up. Furthermore, the small  
424 reductions in F2 induced by perturbed F1 in the previous study were found to persist even  
425 after the feedback was returned to normal, whereas in our experiments, F2 returned to  
426 baseline values but the F1 changes while reduced were not quite returned to baseline.

427  
428 The difference in the findings from the current study and the previous one by MacDonald  
429 and colleagues (2011) may be explained by some experimental differences between the two  
430 studies. For example, in the previous study (MacDonald et al., 2011), only 50 utterances  
431 were produced during the learning phase, which was assumed to be sufficient for speech  
432 adaptation to reach a plateau. In the current study, we witnessed changes in the perturbed  
433 and unperturbed formants that continued to occur after 50 utterances (see Figs. 3A and 4A).  
434 In our experiments, utterances were of three different words with the same vowel “head,  
435 bed, dead” rather than just one word (“head”) used in MacDonald et al., (2011). Intuitively,  
436 the varied context should reduce the rate of learning but it could also allow for transfer or  
437 generalisation of the adapted vowel to different contexts and possibly a more stable effect  
438 by the end of learning. Finally, the magnitude of the shift in F1 feedback was somewhat  
439 larger in the MacDonald study (200Hz) than the equivalent of 110 mel, which we applied in  
440 Experiment 1 (around 150Hz, but tuned to each participant’s formant to account for the  
441 logarithmic relationship between frequency and perception), whereas the F2 shift in the  
442 previous study (250Hz) was closer to our Experiment 2 (~265Hz). The degree to which  
443 participants were aware of the altered feedback may have affected their compensatory  
444 responses in the perturbed formant in particular (though see Munhall et al., 2009). In our  
445 Experiment 2, the perturbed formant F2 was noticeably changed (by 75% of participants)  
446 and the compensatory response in F2 production was of a much larger magnitude than the  
447 change to the perturbed formant F1 in Experiment 1 (only 20% of participants noticed this  
448 change). The magnitude of the change in production of the unperturbed formant did not  
449 differ between the two experiments, however.

450  
451 Our findings are compatible with the idea that speech movement parameters can be  
452 represented in internal models as coordinated and combined sets to achieve specific  
453 sensory goals rather than as precise one-to-one mappings between individual motor  
454 dimensions and acoustic outcomes. However, as we discussed above, our results differ from  
455 those found previously using a similar paradigm (MacDonald et al., 2011). On the basis of  
456 their findings, these authors concluded that the nervous system is capable of independent  
457 control of F1 and F2. Even so, the authors suggested that such precise control of single  
458 formant may be observed only during feedback error correction and this may not extend to  
459 general speech motor control. In our view, the representations in the internal models  
460 tested by feedback perturbation paradigms are likely to be the same as those used in  
461 general speech motor control. We speculate that the inverse and forward internal models  
462 may differ in the extent to which they exert independent and precise articulatory  
463 adjustments to each acoustic parameter during feedback perturbation. A previous study  
464 showed that stimulation to the cerebellum affected only compensatory responses in F1  
465 production in response to an F1 perturbation, whereas stimulation to the motor cortex  
466 affected the amount of compensation in both F1 and F2 in response to the same F1  
467 perturbation (Lametti et al., 2018b). This would suggest that the forward model, which is  
468 used to make rapid ongoing adjustments to speech production can independently adjust  
469 movements in single dimensions corresponding to separate acoustic changes. On the other

470 hand, the inverse model may encode sensorimotor mappings that represent the  
471 combination or combinations of articulatory movements required to achieve a sensory  
472 target that exists in the speech repertoire.

473

474 Using altered auditory feedback, we showed that adaptive responses to feedback shifts  
475 applied to a single formant resulted in changes in production that extended from the shifted  
476 to the unshifted formant. These combined changes moved productions closer to the  
477 frequencies of nearby existing categories in the vowel space. We conclude that although  
478 formants can be controlled independently by adjustments to separate articulatory  
479 dimensions, the motor system underlying speech production may prefer to make combined  
480 changes in formant production. This latter strategy would serve to maintain speech  
481 production within specific regions of the vowel space corresponding to existing speech-  
482 sound categories.

483

484 **Acknowledgements:**

485 This research was supported by a studentship to D-L.T. from the China Scholarship  
486 Council. The Wellcome Centre for Integrative Neuroimaging is supported by core funding  
487 from the Wellcome Trust (203139/Z/16/Z).

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510  
511  
512  
513  
514  
515

## 516 References

517

518 Cai, S., Boucek, M., Ghosh, S. S., Guenther, F. H., & Perkell, J. S. (2008). A system for online  
519 dynamic perturbation of formant trajectories and results from perturbations of the Mandarin  
520 triphthong/iau. *Proceedings of the 8th ISSP*, 65–68.

521 Houde, J. F., & Jordan, M. I. (1998). Sensorimotor adaptation in speech production. *Science*,  
522 279(5354), 1213–1216.

523 Houde, J. F., & Jordan, M. I. (2002). Sensorimotor adaptation of speech I: Compensation and  
524 adaptation. *Journal of Speech, Language, and Hearing Research*, 45(2), 295–310.

525 Lametti, D. R., Nasir, S. M., & Ostry, D. J. (2012). Sensory preference in speech production  
526 revealed by simultaneous alteration of auditory and somatosensory feedback. *J Neurosci*, 32(27),  
527 9351-9358.

528 Lametti, D. R., Rochet-Capellan, A., Neufeld, E., Shiller, D. M., & Ostry, D. J. (2014). Plasticity in  
529 the human speech motor system drives changes in speech perception. *J Neurosci*, 34(31), 10339-  
530 10346.

531 Lametti, D. R., Smith, H. J., Watkins, K. E., & Shiller, D. M. (2018a). Robust Sensorimotor  
532 Learning during Variable Sentence-Level Speech. *Curr Biol*, 28(19), 3106-3113.

533 Lametti, D. R., Smith, H. J., Freidin, P. F., & Watkins, K. E. (2018b). Cortico-cerebellar Networks  
534 Drive Sensorimotor Learning in Speech. *Journal of cognitive neuroscience*, 30(4), 540–551.

535 Ladefoged P (2001) Vowels and consonants: an introduction to the sounds of languages.  
536 Oxford: Blackwell Publishers.

537 MacDonald, E. N., Purcell, D. W., & Munhall, K. G. (2011). Probing the independence of  
538 formant control using altered auditory feedback. *The Journal of the Acoustical Society of America*,  
539 129(2), 955–965.

540 Maeda, S. (1990). Compensatory articulation during speech; Evidence from the analysis and  
541 synthesis of vocal-tract shapes using an articulatory model. In Hardcastle, William J. & Alain Marcha.

542 Matthies, M., Perrier, P., Perkell, J. S., & Zandipour, M. (2001). Variation in anticipatory  
543 coarticulation with changes in clarity and rate. *J Speech Lang Hear Res*, 44(2), 340-353.

544 Mollaei, F., Shiller, D. M., & Gracco, V. L. (2013). Sensorimotor adaptation of speech in  
545 Parkinson's disease. *Movement Disorders*, 28(12), 1668–1674.

546 Nasir SM, Ostry DJ (2008). Speech motor learning in profoundly deaf adults. *Nat Neurosci*  
547 11:1217–1222.

548 Purcell, D. W., & Munhall, K. G. (2006a). Adaptive control of vowel formant frequency:  
549 Evidence from real-time formant manipulation. *The Journal of the Acoustical Society of America*,  
550 120(2), 966–977.

551 Purcell, D. W., & Munhall, K. G. (2006b). Compensation following real-time manipulation of  
552 formants in isolated vowels. *The Journal of the Acoustical Society of America*, 119(4), 2288–2297.

553 Rochet-Capellan, A., & Ostry, D. J. (2011). Simultaneous acquisition of multiple auditory-  
554 motor transformations in speech. *Journal of Neuroscience*, 31(7), 2657–2662.

555 Rochet-Capellan, A., Richer, L., & Ostry, D. J. (2012). Nonhomogeneous transfer reveals  
556 specificity in speech motor learning. *Journal of neurophysiology*, 107(6), 1711–1717.

557 Stevens, S; Volkman; John & Newman, E. B. (1937). "A scale for the measurement of the  
558 psychological magnitude pitch". *Acoust Soc Am*, 8 (3): 185–190.

559           Svirsky, M. A., Teoh, S. W., & Neuburger, H. (2004). Development of language and speech  
560 perception in congenitally, profoundly deaf children as a function of age at cochlear implantation.  
561 *Audiol Neurootol*, 9(4), 224-233.  
562           Tourville, J. A., Cai, S., & Guenther, F. (2013). Exploring auditory-motor interactions in normal  
563 and disordered speech. In *Proceedings of Meetings on Acoustics ICA2013* (Vol. 19, p. 060180). ASA.  
564           Tremblay, S., Shiller, D. M., & Ostry, D. J. (2003). Somatosensory basis of speech production.  
565 *Nature*, 423(6942), 866-869.  
566           Villacorta, V. M., Perkell, J. S., & Guenther, F. H. (2007). Sensorimotor adaptation to feedback  
567 perturbations of vowel acoustics and its relation to perception. *The Journal of the Acoustical Society*  
568 *of America*, 122(4), 2306–2319.  
569           Waldstein, R. S. (1990). Effects of postlingual deafness on speech production: Implications for  
570 the role of auditory feedback. *The Journal of the Acoustical Society of America*.  
571  
572  
573  
574  
575  
576

577 **Figure Captions**

578 **Figure 1. Vowel production.** Schematic showing the relationship between tongue height  
579 and position in the vocal tract (left) and the corresponding frequencies of the first and  
580 second formants of vowels (/I/, /ε/, and /æ/) in vowel space from a single participant in the  
581 current study (right). Data points correspond to F1 and F2 values of single utterances during  
582 the vowel exploration phase in a representative participant. The ellipses represent a 90%  
583 confidence interval around the data points of same colour.

584  
585 **Figure 2. Experimental procedure.** The study involved Vowel Exploration, Baseline, Learning,  
586 and Unlearning, interleaved with three short noise-masked sessions.

587  
588 **Figure 3. Experiment 1: F1 shifted upwards by 110 mel. Change in normalized formant**  
589 **frequencies during speech production.** Values for F1 are shown in blue and F2 in red. The  
590 dashed line reflects no change from baseline. The graph in the top panel (A) shows the  
591 group average changes across learning and unlearning phases of the experiment relative to  
592 baseline; solid circles indicate the group average for each trial; the shaded area represents  
593 the trials where the altered feedback was applied. The graphs in the bottom panel (B, C and  
594 D) show average changes in normalised formant frequencies for individual participants for  
595 the Learning (B), Noise-masked (C) and Unlearning (D) phases of the experiment relative to  
596 baseline; the solid lines represent the group mean; \* indicates significant change from  
597 baseline.

598  
599 **Figure 4. Experiment 2: F2 shifted downwards by 110 mel.** See legend to Fig. 3 for details.

600  
601 **Figure 5. Comparison of changes in vowel production in Experiment 1 and Experiment 2.**  
602 Top panel: (A) Vector figures of experiment 1 (yellow) and experiment 2 (green). The vectors  
603 represent the direction and magnitude of the compensatory change from baseline in vowel  
604 space. The thin vectors represent data from individual participants and the thick vectors  
605 represent the group average for each experiment. The black arrow represents the average  
606 of the direction required to move from “head, bed, dead” to “hid, bid, did”, while the  
607 dotted arrow represents the formant frequency shift we applied (F1+110 mel in Experiment  
608 1 and F2-110 mel in Experiment 2). The light grey circle indicates 110 mel. Bottom panel: (B)  
609 Learning-related reduction in Euclidean distance between “head, bed, dead” and “hid, bid,  
610 did” in F1-F2 space. Negative values indicate that production of “head, bed, dead” by the  
611 end of learning has moved closer to “hid, bid, did”. Bottom panel: (C) The difference  
612 between the angles of the direction of actual compensatory movement and predicted  
613 movement (to “hid, bid, did”) for each participant. Open circles represent individual  
614 participants, the solid lines represent the group means. Distance changes (B) and Actual-  
615 Predicted angles (C) were compared to zero in each experiment. \*  $p < 0.05$  and \*\*  $p < 0.01$

616  
617  
618  
619  
620











