

1 **Title**

2 Fluidity in the perception of auditory speech: Cross-modal recalibration of voice gender
3 and vowel identity by a talking face

4

5 **Authors**

6 Merel A. Burgering¹, Thijs van Laarhoven¹, Martijn Baart^{1,2} & Jean Vroomen¹

7

8 ¹ Department of Cognitive Neuropsychology, Tilburg University, Warandelaan 2, P.O.

9 Box 90153, 5000 LE Tilburg, the Netherlands

10 ² BCBL. Basque Center on Cognition Brain and Language, Donostia - San Sebastián,

11 Spain

12

13 **Corresponding author**

14 Jean Vroomen

15 Email address: j.vroomen@uvt.nl

16 Work phone number: +31 134662394

17

18

19

20

21

22

23

24

25

26 **Abstract**

27 Humans quickly adapt to variations in the speech signal. Adaptation may surface as
28 *recalibration*, a learning effect driven by error-minimization between a visual face and an
29 ambiguous auditory speech signal, or as *selective adaptation*, a contrastive aftereffect
30 driven by the acoustic clarity of the sound. Here, we examined whether these
31 aftereffects occur for vowel identity and voice gender. Participants were exposed to
32 male, female, or androgynous tokens of speakers pronouncing /e/, /ø/, (embedded in
33 words with a consonant-vowel-consonant structure), or an ambiguous vowel halfway
34 between /e/ and /ø/ dubbed onto the video of a male or female speaker pronouncing /e/
35 or /ø/. For both voice gender and vowel identity, we found assimilative aftereffects after
36 exposure to auditory ambiguous adapter sounds, and contrastive aftereffects after
37 exposure to auditory clear adapter sounds. This demonstrates that similar principles for
38 adaptation in these dimensions are at play.

39

40 **Keywords:** audiovisual integration, gender, vowel, recalibration, selective adaptation

41

42

43 **Introduction**

44 Humans constantly integrate different types of sensory input to form coherent
45 representations of the world. This is particularly relevant in social interactions, in which
46 we quickly combine the voice we hear with the face we see when watching our
47 interlocutor. In less than half a second, audiovisual integration processes are initiated
48 that, for example, support perception of the speaker's biological sex – here referred to
49 as gender – (Latinus, VanRullen, & Taylor, 2010), emotion (de Gelder & Vroomen,
50 2000), and phonetic detail of the spoken input (Baart, Lindborg, & Andersen, 2017;
51 Klucharev, Möttönen, & Sams, 2003; Pilling, 2009; Saint-Amour, De Sanctis, Molholm,
52 Ritter, & Foxe, 2007; Stekelenburg & Vroomen, 2007; Sumbly & Pollack, 1954; van
53 Wassenhove, Grant, & Poeppel, 2005).

54 Visual information is helpful to classify voice gender because there is substantial
55 variability in the acoustic parameters that contribute to voice gender (i.e., fundamental
56 frequency, (F0), corresponding to the perceived pitch, (Fenn et al., 2011; Pernet &
57 Belin, 2012; Titze, 1989). Seeing the speaker's face while hearing their voice facilitates
58 categorization of both voice and face gender in terms of response times (Joassin,
59 Maurage, & Campanella, 2011). Also, when facial gender is incongruent with the voice,
60 effects are detrimental rather than facilitatory (Huestegge & Raettig, 2018). The effect of
61 seeing a face on voice gender categorization are also stronger than the effect of hearing
62 a voice on face categorization, suggesting that visual information is more dominant in
63 face-voice gender integration than auditory information (Latinus et al., 2010).

64 Although audiovisual incongruent stimulus materials can contribute to our
65 understanding of multi-sensory integration, it is not clear whether these effects are

66 caused by a genuine perceptual change or by a response bias. For example, an
67 incorrect voice gender response – such as identifying a female voice as ‘male’ when it is
68 presented in combination with a male face – may be caused by visual ‘capture’
69 (participants really perceived a male voice), but it is also possible that participants
70 simply based their response on the visual information only.

71 Under natural circumstances, large incongruencies between a face and voice
72 (such as hearing a male voice and seeing a female face) are rare, but what is much
73 more common is that there is a small discrepancy between what is heard and seen,
74 typically because one of the two signals is unclear, degraded, or ambiguous. This
75 distinction is important, because when the auditory signal is ambiguous rather than fully
76 incongruent with the visual input, listeners may use visual facial cues to perceptually
77 adjust/recalibrate their voice gender categories, as they do for phonetic boundaries
78 (Bertelson, Vroomen, & de Gelder, 2003; Sumbly & Pollack, 1954). This perceptual shift
79 in the auditory modality minimizes the error between the two signals and induces a
80 learning effect that can be measured as an aftereffect in audio-only trials.

81 In the phonetic domain, this effect was first demonstrated by Bertelson et al.
82 (2003) who exposed listeners to a moderate phonetic audiovisual conflict. Participants
83 saw a speaker who pronounced /aba/ (or /ada/) while an ambiguous speech sound
84 halfway between /aba/ and /ada/ – A? for auditory ambiguous – was delivered
85 simultaneously. Immediately after exposure, listeners indicated whether ambiguous
86 audio-only test sounds were either /aba/ or /ada/. Identification of the ambiguous
87 sounds was shifted towards the previously seen lip-read information, so the same test
88 sound was perceived more likely as /aba/ when the previous exposure contained lip-

89 read /aba/ videos, and more likely as /ada/ when exposure contained lip-read /ada/
90 videos. The rationale behind this effect was that during exposure, the perceptual system
91 minimizes the inter-sensory discrepancy by shifting the auditory phonetic boundary,
92 which leads to longer-term assimilative auditory aftereffects. Bertelson et al. (2003)
93 termed the effect phonetic recalibration, which has proven to be a robust phenomenon
94 (Baart, de Boer-Schellekens, & Vroomen, 2012; Baart & Vroomen, 2010; Franken et al.,
95 2017; Keetels, Pecoraro, & Vroomen, 2015; Keetels, Stekelenburg, & Vroomen, 2016;
96 Kilian-Hütten, Vroomen, & Formisano, 2011; van Linden & Vroomen, 2007; Vroomen &
97 Baart, 2009, 2012; Vroomen, Keetels, De Gelder, & Bertelson, 2004; Vroomen, van
98 Linden, Keetels, de Gelder, & Bertelson, 2004).

99 Typically, in the paradigm described above, a control condition is included in
100 which participants are exposed to visual information that is paired with canonical/clear
101 and congruent speech sounds that lead to selective adaptation (Eimas & Corbit, 1973).
102 Selective adaptation differs from recalibration in two important ways. Although the same
103 visual information is presented during exposure, selective adaptation is in the opposite
104 direction of recalibration (a contrastive aftereffect, so after exposure to audiovisual
105 /aba/, listeners show *less* /aba/-responses during the auditory test). This effect is not
106 driven by an inter-sensory conflict, but by the repeated presentation of the unambiguous
107 speech sound itself, and is thus independent of the visual information (Roberts &
108 Summerfield, 1981; Saldaña & Rosenblum, 1994). Contrastive aftereffects may reflect
109 neural fatigue of hypothetical 'linguistic feature detectors' (Eimas & Corbit, 1973), but it
110 has also been proposed that they reflect a criterion shift (see Vroomen & Baart (2012)
111 for an overview) or neural sharpening (Kleinschmidt & Jaeger, 2011).

112 Audiovisual recalibration is quite ubiquitous, as it has also been found to occur
113 for the perception of space (Wozny & Shams, 2011), time (Bermant & Welch, 1976;
114 Bertelson & Aschersleben, 1998; Fujisaki, Shimojo, Kashino, & Nishida, 2004; Keetels
115 & Vroomen, 2007; Radeau & Bertelson, 1974; Vroomen, Keetels, et al., 2004), and for
116 the perception of emotional affect (Baart & Vroomen, 2018). Audiovisual recalibration
117 thus may be a domain-general learning mechanism through which the perceptual
118 system makes necessary adjustments whenever confronted with relatively mild inter-
119 sensory conflicts. Here, the critical question was whether audiovisual recalibration also
120 occurs for the perception of voice gender, which has never been demonstrated before,
121 and vowel identity.

122 Previous studies on phonetic recalibration mostly focused on consonants
123 because consonants have sharper category boundaries than vowels, see for example
124 (Kuhl, 1991). However, there is some evidence that recalibration also occurs for vowels
125 (Franken et al., 2017; Keetels, Bonte, & Vroomen, 2018). Given that identification of
126 voice gender is mainly driven by fundamental frequency of the sound (Gelfer & Mikos,
127 2005), and fundamental frequency is more discernible in vowels than in consonants, we
128 envisaged that vowels would provide an ideal platform to simultaneously assess
129 aftereffects of gender and vowel identity. We therefore used audiovisual recordings of a
130 canonical low-pitched male speaker and a high-pitched female speaker pronouncing the
131 vowels /e/ and /ø/. These vowels were chosen because they are close in F1/F2 acoustic
132 space, and easy to discriminate when lip-reading because the rounding of /ø/ is clearly
133 visible. The vowels were embedded in the context of two Dutch words with a similar
134 frequency of occurrence ('*beek*' [*stream*] and '*beuk*' [*beech*]). These stimuli then allowed

135 us to investigate recalibration and selective adaptation of vowels and voice gender in a
136 within-participant and within-stimulus design.

137 We expected to obtain contrastive aftereffects (indicative of selective adaptation)
138 of voice gender if the auditory tokens were clearly from a male or female speaker
139 (Schweinberger et al., 2008; Zäske, Perlich, & Schweinberger, 2016). Assimilative
140 aftereffects of voice gender (indicative of recalibration) have never been demonstrated
141 before, but as in the phonetic domain, we expected assimilation of voice gender to
142 occur if an androgynous voice was combined with a male or female face. Finding an
143 assimilative effect of voice gender is of interest because it would speak to the generality
144 of the phenomenon since perception of voice gender is quite different from perception of
145 phonemes. For example, voice gender is a more or less stable property over time in the
146 speech signal, which is quite different from phonetic information that is very short-lived
147 and variable between, but also *within* speakers. Furthermore, while vowel categorization
148 occurs in a dense multidimensional acoustic space (largely depending of first and
149 second formant, F1 and F2) that is fine-tuned by language-specific rules, voice gender
150 categorization is, arguably, less complex (a binary male/female distinction, mainly
151 based on fundamental frequency) that is largely shaped by the anatomical differences
152 between the male and female vocal apparatus.

153

154

155

156

157

158 **Methods**

159 *Participants*

160 Thirty students (11 males, 26 right-handed, mean age of 20.6 years, SD = 2.1)
161 from Tilburg University participated in return for course credits or 8 euro/hour¹. All
162 participants reported normal hearing, had (corrected to) normal vision and were naïve to
163 the stimuli and research question. Participants provided written informed consent, and
164 the study was conducted in accordance with the Declaration of Helsinki. The Ethics
165 Review Board of the School of Social and Behavioral Sciences of Tilburg University
166 approved the experimental procedures (EC-2016.48).

167

168 *Stimulus material*

169 *Auditory material.* We selected four artefact-free audiovisual recordings of a male
170 and female native Dutch speaker pronouncing *beek* and *beuk*. The original speech
171 sound *beek* was pronounced as /e/ (the close-mid front unrounded vowel in IPA with F1
172 = 471 Hz and F2 = 2013 Hz for the male speaker and F1 = 498 and F2 = 2261 for
173 female speaker) and the original speech sound *beuk* was pronounced as /ø/ (the close-
174 mid front rounded vowel in IPA with F1 = 455 Hz and F2 = 1539 Hz for the male
175 speaker and F1 = 485 Hz and F2 = 1734 Hz for the female speaker). Tokens were
176 chosen to have matching duration of their vowels (duration of male /beek/ = 702 ms,
177 duration of /e/ = 192 ms; duration of male /beuk/ = 631 ms, duration of /ø/ = 205 ms;
178 duration of female /beek/ = 580 ms, duration of /e/ = 191 ms; duration of female /beuk/ =
179 539 ms, duration of /ø/ = 210 ms). In order to minimize other accidental acoustic

¹ The sample size was larger than in previous work from our lab (see e.g. Bertelson et al., 2003), and was chosen without conducting a formal power analysis.

180 differences between tokens that might serve as a cue for gender or vowel
181 discrimination, we deleted the release of the final consonant /k/ from *beek* and *beuk* (the
182 unvoiced portions) and replaced them by an identical release from /k/ taken from a
183 /beek/ or /beuk/ recording spoken by a different male. These sounds then served as
184 anchors for two male-female gender continua (one for *beuk* and the other for *beek*).
185 They were created using Tandem-STRAIGHT with a step-size of 2% between adjacent
186 tokens (Kawahara et al., 2008). Tandem-STRAIGHT decomposes a speech sound into
187 five sound parameters, namely spectrum, frequency, aperiodicity, fundamental
188 frequency, and time. Each parameter can be adjusted independently. For each speech
189 sound, we manually identified time landmarks (corresponding with the transitions in the
190 spectrogram, such as on- and offsets of the phonation) and frequency landmarks
191 (corresponding with the first three formants in the spectrogram). Morphed stimuli were
192 then generated by re-synthesization based on interpolation (linear for time; logarithmic
193 for F0, frequency and amplitude) (Schweinberger, Kawahara, Simpson, Skuk, & Zäske,
194 2014).

195 We also created two *beuk-beek* vowel continua, one for the male speaker and
196 the other for the female speaker in the same way as described before. We used tokens
197 from the morphing continuum from 5-95% with a step size of 5% from the endpoints
198 towards 40 and 60% and step size of 2% to have higher sampling between 40-60%. We
199 ran a pilot study on seven participants to determine the male-female boundaries ($40.6 \pm$
200 3.3 for the word *beek* [$A_{\text{gender?}}$] and 40.8 ± 4.1 for the word *beuk* [$A_{\text{gender?}}$]), and the
201 *beuk-beek* vowel boundaries (55.8 ± 3.2 for the male speaker [$A_{\text{vowel?male}}$] and $57.1 \pm$
202 2.1 for the female speaker [$A_{\text{vowel?female}}$]). The sounds closest to these boundaries

203 were designated as the ambiguous exposure stimulus and test sound (40 for $A_{\text{gender?}}$;
204 40 for $A_{\text{gender?}}$; 56 for $A_{\text{vowel?male}}$ and 58 for $A_{\text{vowel?female}}$). In order to have variation
205 in the test sounds, we also used stimuli of +8% and -8% (denoted as A_{+1} and A_{-1}).
206 The ambiguous boundary tokens and their ambiguous neighbors were used across all
207 participants.

208 *Visual material.* During exposure, participants saw the video of a male or female
209 speaker pronouncing *beek* or *beuk*. Recordings were framed as frontal headshots. The
210 entire face of the speaker was visible against a neutral black background and measured
211 17° horizontally (ear to ear) and 20° vertically (hairline to chin). The videos were edited
212 in Adobe Premiere. A single exposure phase contained four repetitions of either the
213 male or female speaker saying *beek* or *beuk*. It contained a fade-in and fade-out of two
214 frames at the start and the end of the video resulting in a total duration ~5.48 sec. The
215 audio (clear or ambiguous) was dubbed onto the videos without any noticeable
216 synchronization error.

217

218 *Procedure*

219 *General.* The experiment took place in a dimly lit sound-attenuated room.
220 Instructions and the face of the speaker were presented on a 25-in monitor (BenQ
221 Zowie XL 2540, 240 Hz refresh rate) positioned at eye-level, ~70 cm from the
222 participant's head. The sound was presented through headphones (Sennheiser HD-
223 203) with a peak intensity of 60 dB SPL. The participant responded by pressing one of
224 two buttons on a response box placed in front of the monitor. Participants were
225 instructed to pay attention to the videos displayed on the monitor, which was checked
226 by the experimenter via a live-feed from a camera in the testing booth. These

227 instructions were repeated during the breaks between tasks, and after 24 consecutive
228 exposure-test blocks within each task.

229 *Voice gender identification after audiovisual exposure.*

230 In order to induce voice gender recalibration, participants were exposed to four
231 repetitions (ISI=425 ms) of one of the four audiovisual exposure stimuli containing an
232 *androgynous* voice saying *beek/beuk* dubbed onto a male/female face: $A_{\text{gender?}}V_{\text{male}}$,
233 $A_{\text{gender?}}V_{\text{female}}$, $A_{\emptyset\text{gender?}}V_{\emptyset\text{male}}$ and $A_{\emptyset\text{gender?}}V_{\emptyset\text{female}}$. The exposure phase was
234 immediately followed by a test phase wherein three test sounds were randomly
235 presented, namely the ambiguous voice gender stimulus with the same vowel that was
236 delivered during exposure (henceforth, $/A_{\text{gender?}}/$), and the two close speech morphs on
237 the same continuum $/A_{-1}/$ and $/A_{+1}/$ (Fig. 1A). After each test sound, participants
238 decided whether the test token was 'male' or 'female' in a 2AFC task with two buttons
239 on a response box. The next test sound was played 250 ms after a button press.

240 In order to induce voice gender selective speech adaptation, the exact same
241 procedure was used as for recalibration except that the audiovisual exposure stimuli
242 now contained the *clear* and *gender congruent* audio: (instead of androgynous):
243 $A_{\text{male}}V_{\text{male}}$, $A_{\text{female}}V_{\text{female}}$, $A_{\emptyset\text{male}}V_{\emptyset\text{male}}$, $A_{\emptyset\text{female}}V_{\emptyset\text{female}}$ (Fig. 1B). There were twelve
244 repetitions for each unique exposure-test mini-block, all delivered in pseudo-random
245 order, so in total there were 48 exposure-test mini-blocks for gender recalibration, and
246 48 mini-blocks for gender selective adaptation.

247

248

249

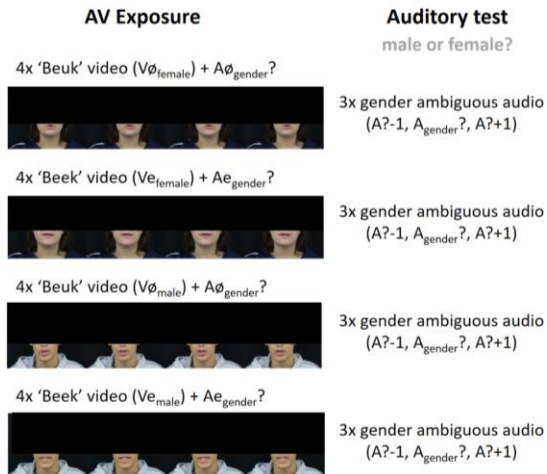
250 *Vowel identification after audiovisual exposure.*

251 To induce vowel recalibration, the same procedures were used as for gender
252 recalibration, except that the four exposure stimuli to assess recalibration were
253 ambiguous with respect to vowel identity: $A_{\text{vowel?male}}V_{\text{emale}}$, $A_{\text{vowel?male}}V_{\text{ømale}}$,
254 $A_{\text{vowel?female}}V_{\text{øfemale}}$ and $A_{\text{vowel?female}}V_{\text{efemale}}$ (henceforth $A_{\text{vowel?}}$). The test sounds
255 were $A_{\text{vowel?}}$ and two neighboring sounds on the *beuk-beek* continua. The exposure
256 stimuli to assess selective adaptation of vowels were, as in voice gender selective
257 adaptation, the gender- and vowel-congruent audiovisual stimuli containing clear audio:
258 $A_{\text{emale}}V_{\text{emale}}$, $A_{\text{efemale}}V_{\text{efemale}}$, $A_{\text{ømale}}V_{\text{ømale}}$, $A_{\text{øfemale}}V_{\text{øfemale}}$.

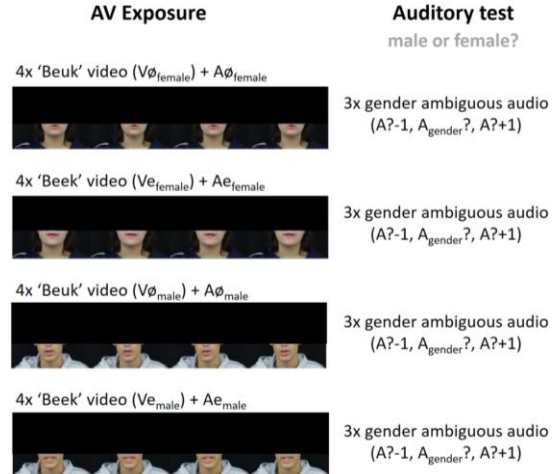
259 Aftereffects of gender and vowel were assessed sequentially with block order
260 counterbalanced across participants. Preliminary analyses showed that block order did
261 not have significant effects on voice gender recalibration and selective adaptation
262 effects, $F_s \leq 1.453$, $p_s \geq .245$, or on vowel recalibration and selective adaptation, $F_s <$
263 $.111$, $p_s > .065$. There was also no significant effect of participant gender on voice
264 gender recalibration and selective adaptation, $F_s \leq .737$, $p_s \geq .401$, or on vowel
265 recalibration and selective adaptation, $F_s \leq 3.358$, $p_s \geq .082$, so block order and gender
266 of the participant were not further analyzed.

267

(A) Recalibration



(B) Selective adaptation



268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

Fig. 1 Overview of the audiovisual exposure-auditory test design. Recalibration (A): four repetitions of a dynamic video of a speaker pronouncing 'beuk' or 'beek' combined with audio of ambiguous voice gender were followed by an auditory-only test in which the participant had to categorize the stimulus into the male or female category. Selective adaptation (B): four repetitions of a dynamic video of a speaker pronouncing 'beuk' or 'beek' combined with audio of either a male or a female speaker were followed by an auditory-only test in which the participant had to categorize the stimulus into the male or female category. The black bars across the upper half of the faces in the figure were included to anonymize the speakers, but were not presented during the experiment.

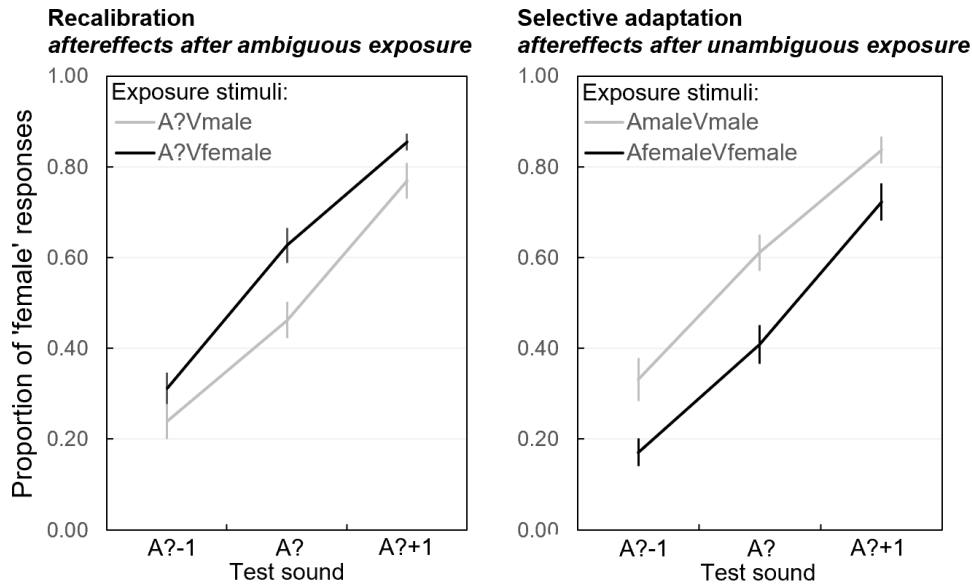
Results

Gender recalibration and adaptation

Individual proportions of 'female' responses on the auditory-only test trials were calculated for each combination of Visual exposure gender (female or male), Auditory exposure type (ambiguous or unambiguous), Vowel (/e/ or /ø/), and Test sound ($A_{\text{gender?}}-1, A_{\text{gender?}}, A_{\text{gender?}}+1$). Data from 9 participants were excluded from the analyses due to unambiguous floor or ceiling effects (see supplementary materials for individual data plots), indicating that they did not adhere to the task instructions or were unable to perform the task correctly. For the remaining 21 participants, grand average proportions of 'female' responses as a function of Visual exposure gender, Vowel, and

287 Test sound are shown for ambiguous and unambiguous auditory exposure types
 288 separately in Figure 2.

289



290

291 **Figure 2.** Averaged proportion of 'female' responses on the auditory test that followed AV exposure
 292 (N=21) in the Gender identification task, averaged across /e/ and /ø/ vowels. Error bars represent one
 293 standard error of the mean.

294

295 A generalized linear mixed-effects model with a logistic linking function to
 296 account for the dichotomous dependent variable was fitted to the single-trial data (lme4
 297 package in R version 3.5.3). The fitted model included Response (male or female
 298 response) as the dependent variable. The model included fixed effects for Visual
 299 exposure gender (male or female), Auditory exposure type (ambiguous or
 300 unambiguous), Vowel (/e/ or /ø/), and Test sound ($A_{\text{gender}}?-1$, $A_{\text{gender}}?$, $A_{\text{gender}}?+1$), with
 301 uncorrelated random intercepts and slopes by participants for the within-participant
 302 variables Visual exposure gender and Auditory exposure type, and their interaction. All
 303 categorical factors were recoded such that their values were centered around 0. Hence,

304 the fitted coefficients could be interpreted as the difference in ‘female’ responses (in log-
305 odds) between two factor levels (e.g. Visual exposure gender male vs female, Auditory
306 exposure type ambiguous vs unambiguous). The fitted model was: $\text{Response} \sim 1 +$
307 $\text{VisualExposureGender} * \text{AuditoryExposureType} * \text{Vowel} * \text{TestSound} + (1 +$
308 $\text{VisualExposureGender} * \text{AuditoryExposureType} | \text{Participant})$. Fixed effect coefficient
309 estimates are shown in Table 1.

310 The analysis revealed a main effect of Test sound ($b = 1.36$, $SE = 0.04$, $p <$
311 0.001), indicative of more ‘female’ responses to the more female-like test sounds, and a
312 main effect of Auditory exposure type ($b = 0.08$, $SE = 0.03$, $p = 0.01$). Importantly, a
313 significant interaction between Visual exposure gender and Auditory exposure type was
314 found ($b = -0.37$, $SE = 0.09$, $p < 0.001$), indicating that the aftereffects of gender were
315 different for ambiguous and unambiguous auditory exposure stimuli. This interaction
316 effect was further examined with post hoc pairwise contrasts (Bonferroni corrected),
317 testing the effect of visual exposure gender at each auditory exposure type. These
318 contrasts showed a higher proportion of ‘female’ responses to the test sounds after
319 exposure to ambiguous sounds paired with a visual female speaker, compared to
320 ambiguous sounds paired with a visual male speaker, thereby demonstrating gender
321 recalibration ($b = 0.58$, $SE = 0.18$, $p = 0.001$). In addition, a higher proportion of male
322 responses was reported after exposure to unambiguous sounds paired with a visual
323 female speaker compared to unambiguous sounds paired with a visual male speaker -
324 indicating gender adaptation, $b = -0.91$, $SE = 0.25$, $p < 0.001$).

325

326

Vowel and voice gender recalibration

327

328

329

Table 1. Fixed effect coefficients and standard errors for the fitted mixed effects regression model:

Response ~ 1 + VisualExposureGender * AuditoryExposureType * Vowel * TestSound + (1 + VisualExposureGender * AuditoryExposureType | Participant)

Fixed factor	Estimate	Standard error	z-value	<i>p</i>
(Intercept)	0.16	0.13	1.242	0.21
VisualExposureGender	0.08	0.06	1.44	0.15
AuditoryExposureType	0.08	0.03	2.56	0.01*
Vowel	-0.02	0.03	-0.66	0.51
TestSound	1.36	0.04	32.74	< 0.001***
VisualExposureGender * AuditoryExposureType	-0.37	0.09	-4.06	< 0.001***
VisualExposureGender * TestSound	-0.03	0.04	-0.76	0.45
VisualExposureGender * Vowel	0.06	0.03	1.78	0.07
AuditoryExposureType * Vowel	0.04	0.03	1.18	0.24
AuditoryExposureType * TestSound	-0.01	0.04	-0.28	0.78
Vowel * Testsound	0.08	0.04	1.99	0.05
VisualExposureGender * AuditoryExposureType * Vowel	-0.04	0.03	-1.21	0.23
VisualExposureGender * AuditoryExposureType * Testsound	0.01	0.04	0.32	0.75
VisualExposureGender * Vowel * Testsound	-0.00	0.04	-0.08	0.94
AuditoryExposureType * Vowel * Testsound	0.01	0.04	0.21	0.83
VisualExposureGender * AuditoryExposureType * Vowel * Testsound	0.05	0.04	1.36	0.17

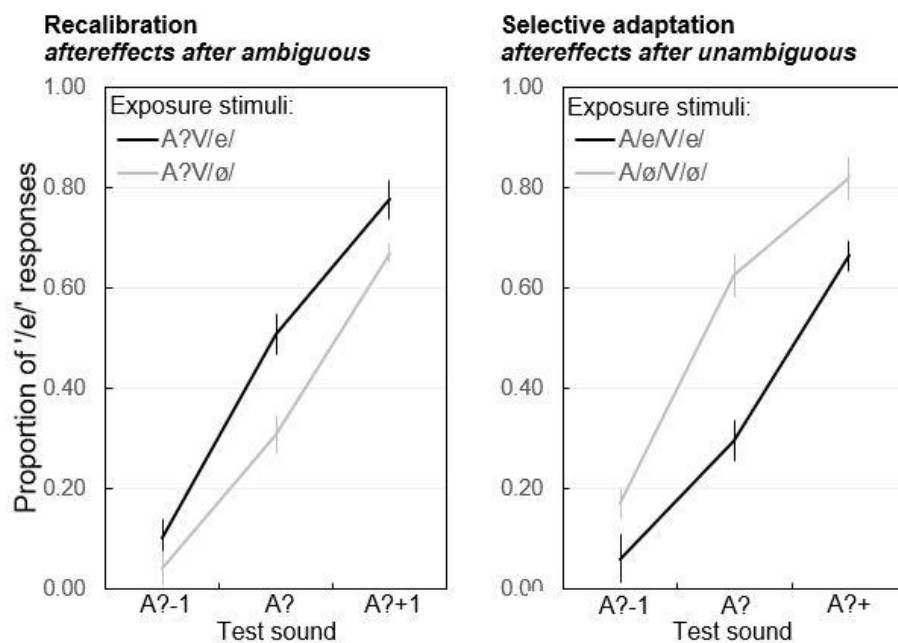
330 **p* < .05; ***p* < .01; ****p* < .001

331

332

333 *Vowel recalibration and adaptation*

334 Individual proportions of /e/ responses on the auditory-only test trials were
 335 calculated for each combination of Visual exposure vowel (/e/ or /ø/), Auditory exposure
 336 type (ambiguous or unambiguous), Gender (female or male), and Test sound ($A_{\text{vowel?}-1}$,
 337 $A_{\text{vowel?}}$, $A_{\text{vowel?}+1}$). Data from 3 participants were excluded from the analyses due to
 338 unambiguous floor or ceiling effects (see supplementary materials for individual data
 339 plots), indicating that they did not adhere to the task instructions or were unable to
 340 perform the task correctly. For the remaining 27 participants, grand average proportions
 341 of /e/ responses as a function of Vowel, Visual exposure gender, and Test sound are
 342 shown for ambiguous and unambiguous auditory exposure types separately in Figure 3.



343
 344 **Figure 3.** Averaged proportion of /e/ responses on the auditory test that followed AV exposure (N=27) in
 345 the Vowel identification task, averaged across male and female sounds. Error bars represent one
 346 standard error of the mean.

347
 348 A generalized linear mixed-effects model with a logistic linking function to

349 account for the dichotomous dependent variable was fitted to the single-trial data (lme4
 350 package in R version 3.5.3). The fitted model included Response (/e/ or /ø/ response)
 351 as the dependent variable, and fixed effects for Visual exposure vowel (/e/ or /ø/),
 352 Auditory exposure type (ambiguous or unambiguous), Gender (female or male), and
 353 Test sound ($A_{\text{vowel?}-1}$, $A_{\text{vowel?}}$, $A_{\text{vowel?}+1}$), with uncorrelated random intercepts and
 354 slopes by participant for the within-participant variables Visual exposure vowel and
 355 Auditory exposure type, and their interaction. All categorical factors were recoded such
 356 that their values were centered around 0. Hence, the fitted coefficients could be
 357 interpreted as the difference in /e/ responses (in log-odds) between two factor levels
 358 (e.g. Visual exposure vowel /e/ vs /ø/, Auditory exposure type ambiguous vs
 359 unambiguous). The fitted model was: $\text{Response} \sim 1 + \text{VisualExposureVowel} * \text{AuditoryExposureType} * \text{Gender} * \text{TestSound} + (1 + \text{VisualExposureVowel} * \text{AuditoryExposureType} | \text{Participant})$. Fixed effect coefficient estimates are shown in
 360
 361
 362 Table 2.
 363

Table 2. Fixed effect coefficients and standard errors for the fitted mixed effects regression model:
 $\text{Response} \sim 1 + \text{VisualExposureVowel} * \text{AuditoryExposureType} * \text{Gender} * \text{TestSound} + (1 + \text{VisualExposureVowel} * \text{AuditoryExposureType} | \text{Participant})$.

Fixed factor	Estimate	Standard error	z-value	P
(Intercept)	-0.52	0.10	-5.38	< 0.001***
VisualExposureVowel	0.11	0.04	2.67	< 0.01**
AuditoryExposureType	-0.12	0.03	-3.62	< 0.001***
Gender	0.25	0.03	8.21	< 0.001***
TestSound	1.79	0.04	42.06	< 0.001***
VisualExposureVowel * AuditoryExposureType	-0.52	0.04	-13.07	< 0.001***

Vowel and voice gender recalibration

VisualExposureVowel * TestSound	0.00	0.04	0.09	0.93
VisualExposureVowel * Gender	-0.07	0.03	-2.23	0.03*
AuditoryExposureType * Gender	-0.01	0.03	-0.42	0.67
AuditoryExposureType * TestSound	0.03	0.04	0.81	0.42
Gender * Testsound	-0.10	0.04	-2.31	0.02*
VisualExposureVowel * AuditoryExposureType * Gender	0.08	0.03	2.70	< 0.01**
VisualExposureVowel * AuditoryExposureType * Testsound	0.06	0.04	1.49	0.14
VisualExposureVowel * Gender * Testsound	0.04	0.04	0.92	0.36
AuditoryExposureType * Gender * Testsound	-0.02	0.04	-0.60	0.55
VisualExposureVowel * AuditoryExposureType * Gender * Testsound	0.01	0.04	0.36	0.72

364 * $p < .05$; ** $p < .01$; *** $p < .001$

365

366 The analysis revealed a negative effect for the intercept ($b = -0.52$, $SE =$
367 0.10 , $p < 0.001$), which indicates a slight overall bias towards /ø/ responses. There was
368 a positive main effect of Test sound ($b = 1.79$, $SE = 0.04$, $p < 0.001$), indicative of more
369 /e/ responses to the more /e/-like test sounds. In addition, there were main effects of
370 Visual exposure vowel ($b = 0.11$, $SE = 0.04$, $p < 0.01$), Auditory exposure type ($b =$
371 -0.12 , $SE = 0.03$, $p < 0.001$), and Gender ($b = 0.25$, $SE = 0.03$, $p < 0.001$), and
372 significant interactions between Visual exposure vowel and Gender ($b = -0.07$, $SE =$
373 0.03 , $p = 0.03$), and between Gender and Test sound ($b = -0.10$, $SE = 0.04$, $p = 0.02$).
374 Importantly, a significant interaction between Visual exposure vowel and Auditory
375 exposure type was found ($b = -0.52$, $SE = 0.04$, $p < 0.001$), indicating that the
376 aftereffects of vowel were different for ambiguous and unambiguous Auditory exposure
377 types. Finally, there was a significant interaction between Visual exposure vowel,
378 Auditory exposure type, and Gender ($b = 0.08$, $SE = 0.03$, $p < 0.01$), indicating that the

379 difference in aftereffects of vowel between the ambiguous and unambiguous Auditory
380 exposure types depended on speaker Gender.

381

382 The three-way interaction effect between Visual exposure vowel, Auditory
383 exposure type, and Gender was further examined with post hoc pairwise contrasts
384 (Bonferroni corrected), testing the Visual exposure vowel \times Auditory exposure
385 interaction at each level of Gender. These contrasts showed a significant Visual
386 exposure vowel \times Auditory exposure interaction for both the male and female speaker
387 (male speaker: $b = -1.73$, $SE = 0.19$, $p < 0.001$, female speaker: $b = -2.40$, $SE = 0.21$, p
388 < 0.001). These interaction effects were further explored with post hoc pairwise
389 contrasts (Bonferroni corrected), which showed significant recalibration and adaptation
390 effects for both the male and female speaker. Specifically, a higher proportion of /e/
391 responses to the auditory-only test trials was reported after exposure to ambiguous
392 sounds paired with visual /e/, compared to ambiguous sounds paired with visual /ø/ (i.e.
393 recalibration), male speaker: $b = 0.78$, $SE = 0.13$, $p < 0.001$, female speaker: $b = 0.84$,
394 $SE = 0.14$, $p < 0.001$). In addition, a higher proportion of /e/ responses was reported
395 after exposure to *unambiguous* sounds paired with visual /ø/ compared to unambiguous
396 sounds paired with visual /e/ (i.e. selective adaptation), male speaker: $b = -0.96$, $SE =$
397 0.15 , $p < 0.001$, female speaker: $b = -1.57$, $SE = 0.16$, $p < 0.001$).

398 As can be seen in Table 3, vowel recalibration was alike across gender of the
399 exposure stimuli, whereas selective adaptation was larger after female than male
400 exposure stimuli, $t(26) = 2.44$, $p = .022$.

401

402

403

404

Table 3. Vowel recalibration and selective adaptation per exposure gender, averaged across test-tokens. Aftereffects were quantified as the difference between proportion of /e/-responses after Visual /e/ and Visual /ø/, resulting in *positive* values for recalibration, and *negative* values for selective adaptation. The ambiguous exposure sound A? was ambiguous in terms of vowel identity (not in terms of gender).

Aftereffect type	Exposure gender (Exposure stimulus)	Aftereffect
Recalibration	Male (A?Vmale)	+.12***
	Female (A?Vfemale)	+.12***
Selective adaptation	Male (A?Vmale)	-.16***
	Female (A?Vfemale)	-.24***

405 * $p < .05$; ** $p < .01$; *** $p < .001$ when tested against 0.

406

407

408 Discussion

409 We found, for the first time, compelling evidence that listeners use the gender of
 410 a male or female face to perceptually adjust (recalibrate) their voice gender category
 411 boundary, which is presumably based on pitch differences between a male/female
 412 voice. When an androgynous voice was dubbed onto the video of a female (instead of
 413 male) face during an audiovisual exposure phase, listeners were more likely to
 414 categorize an androgynous voice as female in auditory-only posttest trials.

415 A similar assimilative effect was found for vowels: an ambiguous vowel halfway
 416 between /e/ and /ø/ dubbed onto the video of a speaker saying /e/ (instead of /ø/) led to
 417 more /e/ responses in auditory-only posttest trials. Gender of the stimulus materials can
 418 modulate vowel identification (Johnson, Strand, & D'Imperio, 1999), and we indeed

419 observed a main effect of Gender on the auditory vowel identification task that followed
420 audiovisual exposure (overall, more /e/ responses were given after exposure to a male
421 rather than female face). Most importantly however, we did not observe a difference in
422 recalibration effect size for vowels induced by male and female exposure materials. We
423 did, however, observe that selective adaptation for vowels was larger after exposure to
424 female adapters rather than male adapters. Johnson et al. (1999) reported that rating
425 female talkers – but not male talkers – as ‘stereotypical’ is correlated with voice
426 breathiness (in addition to fundamental frequency). Perhaps then, breathiness in the
427 female adapter sound constituted an additional acoustic cue that increased the size of
428 the selective adaptation effect, consistent with the notion that the contrastive adaptation
429 effect is mainly driven by the (unambiguous) exposure sound, and not by the video.

430 In order to exclude the possibility that assimilative aftereffects were generated
431 by other mechanisms than recalibration (e.g., priming or a simple response strategy to
432 repeat the exposure stimulus), we included a condition in which the exposure stimuli
433 were audio-visually congruent and thus without inter-sensory conflict. With these stimuli,
434 we found in line with previous studies contrastive aftereffects indicative of selective
435 adaptation (Diehl, 1975; Eimas & Corbit, 1973; Schweinberger et al., 2008; Zäske et al.,
436 2016). Selective adaptation of phonetic information is most likely driven by the
437 unambiguous nature of the auditory component of the audiovisual exposure stimulus
438 and appears to be independent of the visual information (Roberts & Summerfield, 1981;
439 Saldaña & Rosenblum, 1994) The same applies for selective adaptation of voice
440 gender, where the visual information also does not seem to be very relevant. For

441 example, silent articulating faces did not induce adaptation of perceived auditory gender
442 (Schweinberger et al., 2008).

443 It remains to be examined in future studies *what* representation listeners adjusted
444 in the case of the gender recalibration task: listeners might have shifted their
445 male/female voice category in general, or only for these two talkers that they heard
446 during the exposure phase. Previous studies on *phonetic* calibration have demonstrated
447 that recalibration is extremely token-specific, and that it even can be ear- and location-
448 specific so that the same ambiguous sound can be simultaneously adapted to two
449 opposing phonetic interpretations if presented in the left and right ear (Keetels et al.,
450 2015). Generalization of recalibration of voice gender, though, might be different. In an
451 informal pilot study (Burgering, Baart, & Vroomen, 2018), we had switched talkers - but
452 not gender - between exposure and test and observed comparable aftereffects. This
453 result, at least tentatively, suggests that voice gender recalibration is not speaker-, or
454 token-specific, but rather generalizes across speakers and tokens.

455 Another intriguing question for future research is to examine to which extent
456 adaptation in voice gender and voice identity rely on common or separate neural
457 mechanisms. It seems likely that some mechanisms will be shared, while others will be
458 separate. As an example, a study by Green and colleagues (Green, Kuhl, Meltzoff, &
459 Stevens, 1991) provided behavioral evidence that perception of gender and phonetic
460 information rely on dimension-specific mechanisms. The authors showed that the
461 McGurk illusion – such as hearing /da/ when auditory /ba/ is delivered in combination
462 with a face articulating /ga/ – was not modulated by gender incongruency in the
463 audiovisual stimulus, despite the fact that the face-voice gender mismatch was perfectly

464 clear. Audiovisual integration of phonetic information thus seems to be, at least partially,
465 independent of audiovisual integration of gender information. A reason for this might be
466 that indexical information such as emotional affect or gender is quite holistic in nature
467 and can be acquired from an image or a simple vocalization. In contrast, phonetic
468 processing of speech relies on the fine-grained temporal coherence between what is
469 seen and heard (Cellerino, Borghetti, & Sartucci, 2004; Curby, Johnson, & Tyson, 2012;
470 Lewin & Herlitz, 2002; Sun, Gao, & Han, 2010; Tottenham et al., 2009).

471 The timing of when gender and phonetic information becomes available, though,
472 might be similar. In an EEG (Electroencephalography) study, Latinus et al. (2010)
473 observed that congruency between facial and vocal gender modulated brain processes
474 within 180 ms and 230 ms after stimulus onset, which aligns with the time-frame during
475 which auditory-only gender differences are processed (Latinus & Taylor, 2012; Zäske,
476 Schweinberger, Kaufmann, & Kawahara, 2009). Interestingly, processing of phonetic
477 congruency is also (partially) realized during this time-window (Arnal, Morillon, Kell, &
478 Giraud, 2009; Baart et al., 2017; Baart, Stekelenburg, & Vroomen, 2014; Stekelenburg
479 & Vroomen, 2007) and audiovisual congruency processing of gender and phonetic
480 information thus overlap in time.

481 It also remains for future studies to examine whether there is a common neural
482 mechanism for recalibration of voice gender and vowel identity,, especially since there
483 seems to be a good candidate brain region that should be involved in this process: the
484 superior temporal sulcus (STS). Specifically, the STS is involved in lip-read-induced
485 phonetic recalibration (Kilian-Hütten, Valente, Vroomen, & Formisano, 2011), as well as
486 text-induced phonetic recalibration (especially in the right hemisphere, see (Bonte,

487 Correia, Keetels, Vroomen, & Formisano, 2017), and is also part of a right hemisphere
488 dominated network related to processing vocal gender (Belin et al., 2000; Imaizumi et
489 al., 1997; Von Kriegstein, Eger, Kleinschmidt, & Giraud, 2003; von Kriegstein, Smith,
490 Patterson, Kiebel, & Griffiths, 2010), and cross modal integration of face and voice
491 (Blank, Anwender, & von Kriegstein, 2011; Campanella & Belin, 2007; Von Kriegstein,
492 Kleinschmidt, Sterzer, & Giraud, 2005).

493 To conclude, humans can flexibly adjust their perceived voice gender categories
494 based on previous exposure. The results are in line with previous studies on voice-face
495 interaction, and the underlying mechanisms seem to operate like those that underlie
496 phonetic selective adaptation and recalibration. The current study inspires future work
497 on the domain general versus domain specific aspects of recalibration.

498

499 **Acknowledgement**

500 This research was supported by Gravitation Grant 024.001.006 of the Language in
501 Interaction Consortium from Netherlands Organization for Scientific Research. The third
502 author was supported by The Netherlands Organization for Scientific Research (NWO:
503 VENI Grant 275-89-027). We thank Alicia Driessen for help with the data collection.

504

505

506

507

508

509

510

511

512

513 **References**

- 514 Arnal, L. H., Morillon, B., Kell, C. A., & Giraud, A. L. (2009). Dual neural routing of visual facilitation in
 515 speech processing. *Journal of Neuroscience*, *29*(43), 13445-13453.
 516 doi:10.1523/JNEUROSCI.3194-09.2009
- 517 Baart, M., de Boer-Schellekens, L., & Vroomen, J. (2012). Lipread-induced phonetic recalibration in
 518 dyslexia. *Acta Psychologica*, *140*(1), 91-95. doi:10.1016/j.actpsy.2012.03.003
- 519 Baart, M., Lindborg, A., & Andersen, T. S. (2017). Electrophysiological evidence for differences between
 520 fusion and combination illusions in audiovisual speech perception. *Eur J Neurosci*, *46*(10), 2578-
 521 2583. doi:10.1111/ejn.13734
- 522 Baart, M., Stekelenburg, J. J., & Vroomen, J. (2014). Electrophysiological evidence for speech-specific
 523 audiovisual integration. *Neuropsychologia*, *53*, 115-121.
 524 doi:10.1016/j.neuropsychologia.2013.11.011
- 525 Baart, M., & Vroomen, J. (2010). Phonetic recalibration does not depend on working memory.
 526 *Experimental brain research*, *203*(3), 575-582. doi:10.1007/s00221-010-2264-9
- 527 Baart, M., & Vroomen, J. (2018). Recalibration of vocal affects by a dynamic face. *Experimental brain*
 528 *research*, 1-8.
- 529 Belin, P., Fecteau, S., & Bedard, C. (2004). Thinking the voice: neural correlates of voice perception.
 530 *Trends in cognitive sciences*, *8*(3), 129-135.
- 531 Belin, P., Zatorre, R. J., Lafaille, P., Ahad, P., & Pike, B. (2000). Voice-selective areas in human auditory
 532 cortex. *Nature*, *403*(6767), 309.
- 533 Bermant, R. I., & Welch, R. B. (1976). Effect of degree of separation of visual-auditory stimulus and eye
 534 position upon spatial interaction of vision and audition. *Perceptual and Motor Skills*, *43*(2), 487-
 535 493. doi:10.2466/pms.1976.43.2.487
- 536 Bertelson, P., & Aschersleben, G. (1998). Automatic visual bias of perceived auditory location.
 537 *Psychonomic bulletin & review*, *5*(3), 482-489.
- 538 Bertelson, P., Vroomen, J., & de Gelder, B. (2003). Visual recalibration of auditory speech identification:
 539 A McGurk aftereffect. *Psychological Science*, *14*(6), 592-597. doi:10.1046/j.0956-
 540 7976.2003.psci_1470.x
- 541 Bestelmeyer, P. E., Belin, P., & Grosbras, M. H. (2011). Right temporal TMS impairs voice detection.
 542 *Current Biology*, *21*(20), R838-R839. doi:10.1016/j.cub.2011.08.046
- 543 Blank, H., Anwander, A., & von Kriegstein, K. (2011). Direct structural connections between voice-and
 544 face-recognition areas. *Journal of Neuroscience*, *31*(96), 12906-12915.
 545 doi:10.1523/JNEUROSCI.2091-11.2011
- 546 Bonte, M., Correia, J. M., Keetels, M., Vroomen, J., & Formisano, E. (2017). Reading-induced shifts of
 547 perceptual speech representations in auditory cortex. *Scientific reports*, *7*. doi:10.1038/s41598-
 548 017-05356-3
- 549 Bosker, H. R., Reinisch, E., & Sjerps, M. J. (2017). Cognitive load makes speech sound fast, but does not
 550 modulate acoustic context effects. *Journal of Memory and Language*, *94*, 166-176.
- 551 Burgering, M. A., Baart, M., & Vroomen, J. (2018, June 14-17). *Audiovisual recalibration and selective*
 552 *adaptation for vowels and speaker sex*. Paper presented at the 19th International Multisensory
 553 Research Forum (IMRF), Toronto, Canada.
- 554 Campanella, S., & Belin, P. (2007). Integrating face and voice in person perception. *Trends in cognitive*
 555 *sciences*, *11*(12), 535-543. doi:10.1016/j.tics.2007.10.001
- 556 Cellerino, A., Borghetti, D., & Sartucci, F. (2004). Sex differences in face gender recognition in humans.
 557 *Brain research bulletin*, *63*(6), 443-449. doi:10.1016/j.brainresbull.2004.03.010
- 558 Charest, I., Pernet, C., Latinus, M., Crabbe, F., & Belin, P. (2012). Cerebral processing of voice gender
 559 studied using a continuous carryover fMRI design. *Cerebral Cortex*, *23*(4), 958-966.

- 560 Curby, K. M., Johnson, K. J., & Tyson, A. (2012). Face to face with emotion: Holistic face processing is
 561 modulated by emotional state. *Cognition & Emotion*, *26*(1), 93-102.
 562 doi:10.1080/02699931.2011.555752
- 563 de Gelder, B., & Vroomen, J. (2000). The perception of emotions by ear and by eye. *Cognition &*
 564 *Emotion*, *14*(3), 289-311.
- 565 Diehl, R. L. (1975). The effect of selective adaptation on the identification of speech sounds. *Perception*
 566 *& psychophysics*, *17*(1), 48-52.
- 567 Eimas, P. D., & Corbit, J. D. (1973). Selective adaptation of linguistic feature detectors. *Cognitive*
 568 *psychology*, *4*(1), 99-109.
- 569 Feng, G., Yi, H. G., & Chandrasekaran, B. (2018). The Role of the Human Auditory Corticostriatal Network
 570 in Speech Learning. *Cerebral Cortex*.
- 571 Fenn, K. M., Shintel, H., Atkins, A. S., Skipper, J. I., Bond, V. C., & Nusbaum, H. C. (2011). When less is
 572 heard than meets the ear: Change deafness in a telephone conversation. *The Quarterly Journal*
 573 *of Experimental Psychology*, *64*(7), 1442-1456. doi:10.1080/17470218.2011.570353
- 574 Franken, M., Eisner, F., Schoffelen, J., Acheson, D. J., Hagoort, P., & McQueen, J. M. (2017). *Audiovisual*
 575 *recalibration of vowel categories*. Paper presented at the Proceedings of Interspeech 2017.
- 576 Fujisaki, W., Shimojo, S., Kashino, M., & Nishida, S. Y. (2004). Recalibration on audiovisual simultaneity.
 577 *Nature Neuroscience*, *7*(7), 773.
- 578 Gelfer, M. P., & Mikos, V. A. (2005). The relative contributions of speaking fundamental frequency and
 579 formant frequencies to gender identification based on isolated vowels. *Journal of Voice*, *19*(4),
 580 544-554. doi:10.1016/j.jvoice.2004.10.006
- 581 Green, K. P., Kuhl, P. K., Meltzoff, A. N., & Stevens, E. B. (1991). Integrating speech information across
 582 talkers, gender, and sensory modality: Female faces and male voices in the McGurk effect.
 583 *Perception & psychophysics*, *50*(6), 524-536.
- 584 Huestegge, S. M., & Raettig, T. (2018). Crossing gender borders: bidirectional dynamic interaction
 585 between face-based and voice-based gender categorization. *Journal of Voice*.
- 586 Imaizumi, S., Mori, K., Kiritani, S., Kawashima, R., Sugiura, M., Fukuda, H., . . . Hatano, K. (1997). Vocal
 587 identification of speaker and emotion activates different brain regions. *Neuroreport*, *8*(12),
 588 2809-2812.
- 589 Jäncke, L., Wüstenberg, T., Scheich, H., & Heinze, H. J. (2002). Phonetic perception and the temporal
 590 cortex. *NeuroImage*, *15*(4), 733-746.
- 591 Joassin, F., Maurage, P., & Campanella, S. (2011). The neural network sustaining the crossmodal
 592 processing of human gender from faces and voices: An fMRI study. *NeuroImage*, *54*(2), 1654-
 593 1661.
- 594 Johnson, K., Strand, E. A., & D'Imperio, M. (1999). Auditory–visual integration of talker gender in vowel
 595 perception. *Journal of phonetics*, *27*(4), 359-384.
- 596 Kawahara, H., Morise, M., Takahashi, T., Nisimura, R., Irino, T., & Banno, H. (2008). Tandem-STRAIGHT: A
 597 temporally stable power spectral representation for periodic signals and applications to
 598 interference-free spectrum, F0, and aperiodicity estimation. *Acoustics, Speech and Signal*
 599 *Processing, ICASSP 2008. IEEE International Conference*, 3933-3936.
- 600 Keetels, M., Bonte, M., & Vroomen, J. (2018). A Selective Deficit in Phonetic Recalibration by Text in
 601 Developmental Dyslexia. *Frontiers in psychology*, *9*.
- 602 Keetels, M., Pecoraro, M., & Vroomen, J. (2015). Recalibration of auditory phonemes by lipread speech
 603 is ear-specific. *Cognition*, *141*, 121-126. doi:10.1016/j.cognition.2015.04.019
- 604 Keetels, M., Stekelenburg, J. J., & Vroomen, J. (2016). A spatial gradient in phonetic recalibration by
 605 lipread speech. *Journal of phonetics*, *56*, 124-130. doi:10.1016/j.wocn.2016.02.005
- 606 Keetels, M., & Vroomen, J. (2007). No effect of auditory-visual spatial disparity on temporal
 607 recalibration. *Experimental brain research*, *182*(4), 559-565.

- 608 Kilian-Hütten, N., Valente, G., Vroomen, J., & Formisano, E. (2011). Auditory cortex encodes the
 609 perceptual interpretation of ambiguous sound. *Journal of Neuroscience*, *31*(5), 1715-1720.
- 610 Kilian-Hütten, N., Vroomen, J., & Formisano, E. (2011). Brain activation during audiovisual exposure
 611 anticipates future perception of ambiguous speech. *NeuroImage*, *57*(4), 1601-1607.
 612 doi:10.1016/j.neuroimage.2011.05.043
- 613 Kleinschmidt, D., & Jaeger, T. F. (2011). A Bayesian belief updating model of phonetic recalibration and
 614 selective adaptation. *Proceedings of the 2nd Workshop on Cognitive Modeling and*
 615 *Computational Linguistics*, 10-19.
- 616 Klucharev, V., Möttönen, R., & Sams, M. (2003). Electrophysiological indicators of phonetic and non-
 617 phonetic multisensory interactions during audiovisual speech perception. *Cognitive Brain*
 618 *Research*, *18*(1), 65-75. doi:10.1016/j.cogbrainres.2003.09.004
- 619 Kuhl, P. K. (1991). Human adults and human infants show a “perceptual magnet effect” for the
 620 prototypes of speech categories, monkeys do not. *Perception & psychophysics*, *50*(2), 93-107.
- 621 Latinus, M., & Taylor, M. J. (2012). Discriminating male and female voices: differentiating pitch and
 622 gender. *Brain topography*, *25*(2), 194-204.
- 623 Latinus, M., VanRullen, R., & Taylor, M. J. (2010). Top-down and bottom-up modulation in processing
 624 bimodal face/voice stimuli. *BMC neuroscience*, *11*(1), 36.
- 625 Lewin, C., & Herlitz, A. (2002). Sex differences in face recognition - Women's faces make the difference.
 626 *Brain and cognition*, *50*(1), 121-128.
- 627 Liebenthal, E., Binder, J. R., Spitzer, S. M., Possing, E. T., & Medler, D. A. (2005). Neural substrates of
 628 phonemic perception. *Cerebral Cortex*, *15*(10), 1621-1631.
- 629 Liebenthal, E., Sabri, M., Beardsley, S. A., Mangalathu-Arumana, J., & Desai, A. (2013). Neural dynamics
 630 of phonological processing in the dorsal auditory stream. *Journal of Neuroscience*, *33*(39),
 631 15414-15424.
- 632 Modelska, M., Pourquié, M., & Baart, M. (2019). No “Self” Advantage for Audiovisual Speech
 633 Aftereffects. *Frontiers in psychology*, *10*(658).
- 634 Pernet, C. R., & Belin, P. (2012). The role of pitch and timbre in voice gender categorization. *Frontiers in*
 635 *psychology*, *3*, 23.
- 636 Pilling, M. (2009). Auditory event-related potentials (ERPs) in audiovisual speech perception. *Journal of*
 637 *Speech, Language, and Hearing Research*, *52*(4), 1073-1081.
- 638 Radeau, M., & Bertelson, P. (1974). The after-effects of ventriloquism. *The Quarterly Journal of*
 639 *Experimental Psychology*, *26*(1), 63-71.
- 640 Reinisch, E., & Sjerps, M. J. (2013). The uptake of spectral and temporal cues in vowel perception is
 641 rapidly influenced by context. *Journal of phonetics*, *41*(2), 101-116.
- 642 Roberts, M., & Summerfield, Q. (1981). Audiovisual presentation demonstrates that selective adaptation
 643 in speech perception is purely auditory. *Perception & psychophysics*, *30*(4), 309-314.
- 644 Saint-Amour, D., De Sanctis, P., Molholm, S., Ritter, W., & Foxe, J. J. (2007). Seeing voices: High-density
 645 electrical mapping and source-analysis of the multisensory mismatch negativity evoked during
 646 the McGurk illusion. *Neuropsychologia*, *45*(3), 587-597.
 647 doi:10.1016/j.neuropsychologia.2006.03.036
- 648 Saldaña, H. M., & Rosenblum, L. D. (1994). Selective adaptation in speech perception using a compelling
 649 audiovisual adaptor. *The Journal of the Acoustical Society of America*, *95*(6), 3658-3661.
- 650 Schweinberger, S. R., Casper, C., Hauthal, N., Kaufmann, J. M., Kawahara, H., Kloth, N., . . . Zäske, R.
 651 (2008). Auditory Adaptation in Voice Perception. *Current Biology*, *18*, 684-688.
 652 doi:10.1016/j.cub.2008.04.015
- 653 Schweinberger, S. R., Kawahara, H., Simpson, A. P., Skuk, V. G., & Zäske, R. (2014). Speaker perception.
 654 *Wiley Interdisciplinary Reviews: Cognitive Science*, *5*(1), 15-25.

- 655 Stekelenburg, J. J., & Vroomen, J. (2007). Neural correlates of multisensory integration of ecologically
656 void audiovisual events. *Journal of cognitive neuroscience*, *19*(12), 1964-1973.
- 657 Sugano, Y., Keetels, M., & Vroomen, J. (2016). Auditory dominance in motor-sensory temporal
658 recalibration. *Experimental brain research*, *234*(5), 1249-1262. doi:10.1007/s00221-015-4497-0
- 659 Sumbly, W. H., & Pollack, I. (1954). Visual contribution to speech intelligibility in noise. *The Journal of the*
660 *Acoustical Society of America*, *26*(2), 212-215.
- 661 Sun, Y., Gao, X., & Han, S. (2010). Sex differences in face gender recognition: an event-related potential
662 study. *Brain research*, *1327*(69-76).
- 663 Titze, I. R. (1989). Physiologic and acoustic differences between male and female voices. *The Journal of*
664 *the Acoustical Society of America*, *85*(4), 1699-1707.
- 665 Tottenham, N., Tanaka, J. W., Leon, A. C., McCarry, T., Nurse, M., Hare, T. A., . . . Nelson, C. (2009). The
666 NimStim set of facial expressions: Judgements from untrained research participants. *Psychiatry*
667 *Research*, *168*(3), 242-249.
- 668 van Linden, S., & Vroomen, J. (2007). Recalibration of phonetic categories by lipread speech versus
669 lexical information. *Journal of Experimental Psychology: Human Perception & Performance*,
670 *33*(6), 1483-1494. doi:10.1037/0096-1523.33.6.1483
- 671 van Wassenhove, V., Grant, K. W., & Poeppel, D. (2005). Visual speech speeds up the neural processing
672 of auditory speech. *Proceedings of the National Academy of Sciences of the United States of*
673 *America*, *102*(4), 1181-1186. doi:10.1073/pnas.0408949102
- 674 Von Kriegstein, K., Eger, E., Kleinschmidt, A., & Giraud, A. L. (2003). Modulation of neural responses to
675 speech by directing attention to voices or verbal content. *Cognitive Brain Research*, *17*(1), 48-55.
- 676 Von Kriegstein, K., Kleinschmidt, A., Sterzer, P., & Giraud, A. L. (2005). Interaction of face and voice areas
677 during speaker recognition. *Journal of cognitive neuroscience*, *17*(3), 367-376.
- 678 von Kriegstein, K., Smith, D. R. R., Patterson, R. D., Kiebel, S. J., & Griffiths, T. D. (2010). How the human
679 brain recognizes speech in the context of changing speakers. *Journal of Neuroscience*, *30*(2),
680 629-638.
- 681 Vroomen, J., & Baart, M. (2009). Phonetic recalibration only occurs in speech mode. *Cognition*, *110*(2),
682 254-259. doi:10.1016/j.cognition.2008.10.015
- 683 Vroomen, J., & Baart, M. (2012). Phonetic Recalibration in Audiovisual Speech. In M. M. Murray,
684 Wallace, M.T. (Ed.), *The Neural Bases of Multisensory Processes*. Frontiers in Neuroscience: CRC
685 Press/Taylor & Francis.
- 686 Vroomen, J., Keetels, M., De Gelder, B., & Bertelson, P. (2004). Recalibration of temporal order
687 perception by exposure to audio-visual asynchrony. *Cognitive Brain Research*, *22*(1), 32-35.
- 688 Vroomen, J., van Linden, S., Keetels, M., de Gelder, B., & Bertelson, P. (2004). Selective adaptation and
689 recalibration of auditory speech by lipread information: Dissipation. *Speech Communication*, *44*,
690 55-61.
- 691 Wozny, D. R., & Shams, L. (2011). Recalibration of auditory space following milliseconds of cross-modal
692 discrepancy. *Journal of Neuroscience*, *31*(12), 4607-4612.
- 693 Zäske, R., Perlich, M. C., & Schweinberger, S. R. (2016). To hear or not to hear: Voice processing under
694 visual load. *Attention Perception & Psychophysics*, *78*(5), 1488-1495. doi:10.3758/s13414-016-
695 1119-2
- 696 Zäske, R., Schweinberger, S. R., Kaufmann, J. M., & Kawahara, H. (2009). In the ear of the beholder:
697 neural correlates of adaptation to voice gender. *European Journal of Neuroscience*, *30*, 527-534.
698 doi:10.1111/j.1460-9568.2009/06839.x
- 699 Zäske, R., Schweinberger, S. R., & Kawahara, H. (2010). Voice aftereffects of adaptation to speaker
700 identity. *Hearing Research*, *268*, 38-45. doi: 10.1016/j.heares.2010.04.011.

701

Supplementary materials

702

703

704

705

706

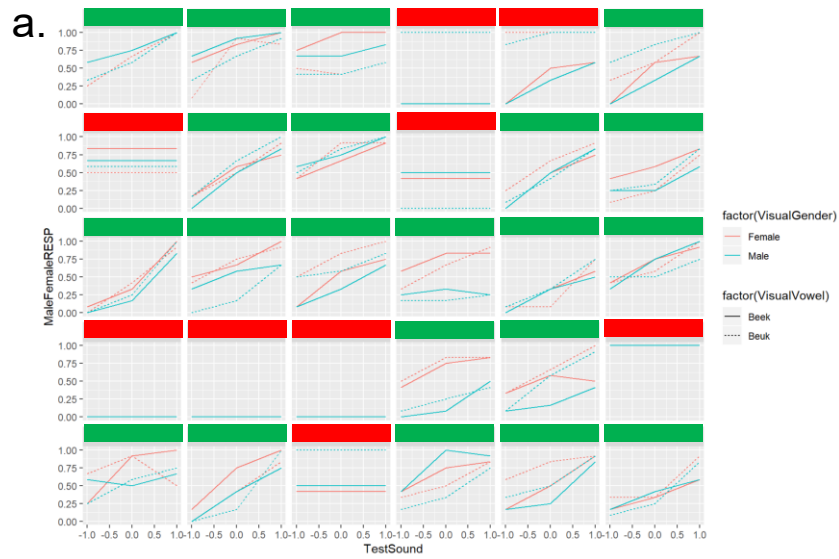
707

708

709

710

711



712

713

714

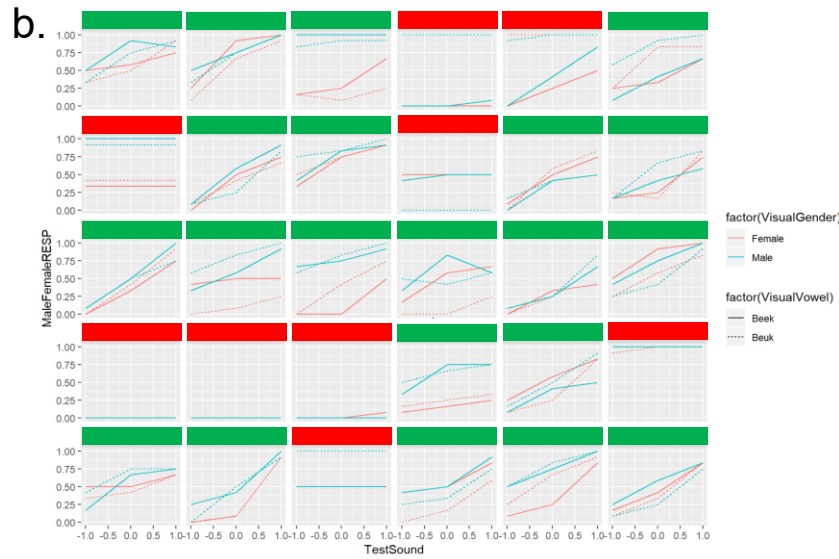
715

716

717

718

719



720

Figure S1. Proportion of female responses in the auditory Gender identification task after AV exposure

721

for all individual participants (N = 30). Participants highlighted by red bars were excluded (N = 9) from the

722

analyses due to ceiling effects (indicating that the test tokens did not represent their perceptual

723

boundaries, and/or participants simply pressed only one key during the test for unknown reasons), or

724

otherwise questionable data patterns. Panel a. represents the data after exposure to ambiguous

725

adapters, panel b. represents the data after exposure to unambiguous adapters.

726

Vowel and voice gender recalibration

727

728

729

730

731

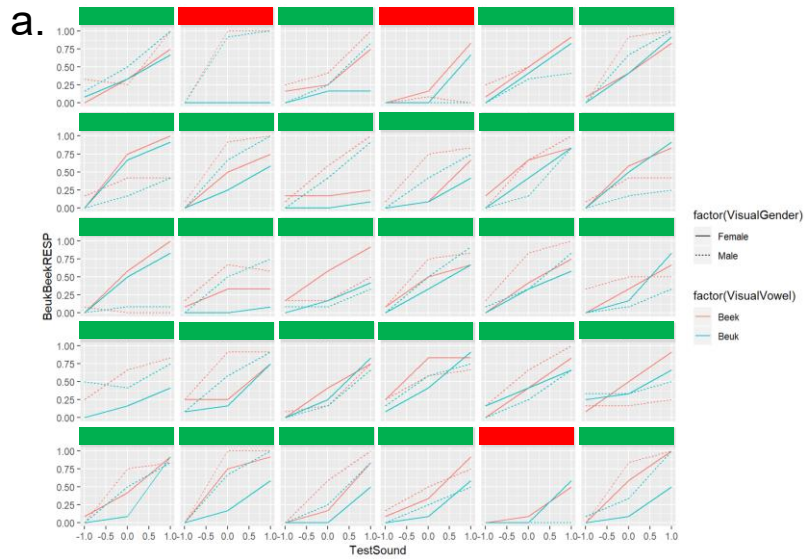
732

733

734

735

736



737

738

739

740

741

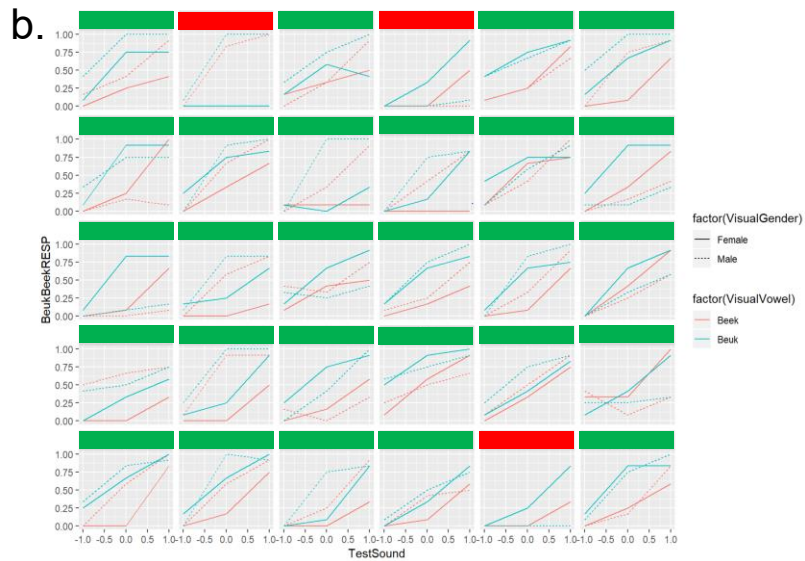
742

743

744

745

746



747 **Figure S2.** Proportion of /e/ responses in the auditory Gender identification task after AV exposure for all

748 individual participants (N = 30). Participants highlighted by red bars (N = 3) were excluded from the

749 analyses due to ceiling effects (indicating that the test tokens did not represent their perceptual

750 boundaries, and/or participants simply pressed only one key during the test for unknown reasons), or

751 otherwise questionable data patterns. Panel a. represents the data after exposure to ambiguous

752 adapters, panel b. represents the data after exposure to unambiguous adapters.

753