

1 **Familiar voices are more intelligible, even if they are not recognized as familiar**

2 Holmes, E.¹, Domingo, Y.¹, & Johnsrude, I. S.^{1,2}

3 ¹Brain and Mind Institute, University of Western Ontario, Canada

4 ²School of Communication Sciences and Disorders, University of Western Ontario, Canada

5

6

7 Corresponding author: Emma Holmes; E-mail: emma.holmes@ucl.ac.uk; Phone: +44 7597

8 967397; Mailing address: Wellcome Centre for Human Neuroimaging, Institute of Neurology,

9 University College London, 12 Queen Square, London WC1N 3BG, U.K.

Abstract

10

11

12

13

14

15

16

17

18

19

20

21

22

We can recognize familiar people by their voices, and familiar talkers are more intelligible than unfamiliar talkers when competing talkers are present. However, whether the acoustic voice characteristics that permit recognition and those that benefit intelligibility are the same or different is unknown. Here, we recruited pairs of participants who had known each other for 6 months or longer, and manipulated the acoustic correlates of two voice characteristics (vocal tract length and glottal pulse rate). These had different effects on explicit recognition of, and the speech-intelligibility benefit realized from, familiar voices. Furthermore, even when explicit recognition of familiar voices was eliminated, they were still more intelligible than unfamiliar voices—demonstrating that familiar voices do not need to be explicitly recognized to benefit intelligibility. Processing familiar-voice information appears therefore to depend on multiple, at least partially independent, systems that are recruited depending on the perceptual goal of the listener.

23

Introduction

24 When we converse with other people, we become familiar with their voices, and this
25 enables us to subsequently recognize those people by voice. Historically, the components of
26 speech that convey talker-identity information ('the carrier') were considered separately from
27 those that convey the spoken message ('the content'; Halle, 1985; Joos, 1948). Indeed, brain
28 activity differs when participants attend to speech content or the speaker's identity (von
29 Kriegstein, Kleinschmidt, Sterzer, & Giraud, 2005), showing that information about the carrier is
30 encoded at least partially separately from the content. Intriguingly, however, familiar-voice
31 information can aid intelligibility of degraded speech content. In the presence of a competing
32 talker, listeners find speech more intelligible if it is spoken by a familiar than unfamiliar talker
33 (Domingo, Holmes, & Johnsrude, submitted; Johnsrude et al., 2013; Kreitewolf, Mathias, & von
34 Kriegstein, 2017; Levi, Winters, & Pisoni, 2011; Nygaard & Pisoni, 1998; Nygaard, Sommers, &
35 Pisoni, 1994; Yonan & Sommers, 2000). Thus, experience with a carrier aids in identification of
36 content. However, the acoustic characteristics that underlie the benefit to speech intelligibility
37 from a familiar voice—and whether they are the same as those that are critical for recognizing a
38 voice as familiar—are currently unknown.

39 Speech spoken by different talkers varies on several dimensions. The source-filter
40 model of speech production (Fant, 1960; Chiba & Kajiyama, 1941) assumes that the acoustics
41 of speech result from the action of the articulatory filter upon the vocal source, which is created
42 through vocal-fold vibration. The rate of vocal-fold vibration (which is also known as the glottal
43 pulse rate) is related to the mass of the vocal folds. The rate of vibration determines the
44 fundamental frequency (f_0) of the speech signal. This source is dynamically filtered by the vocal
45 tract, which differs in length and shape between different talkers. These properties of the vocal
46 tract determine the resonances, or formants, of speech, which are frequency-specific
47 concentrations of sound energy. Both f_0 and formant spacing are somewhat variable within

48 talkers. Although vocal-tract characteristics are relatively fixed within a talker, the shape of the
49 vocal cavity changes when talkers alter the positions of the articulators (e.g., lips and tongue) to
50 create different sounds (e.g., Hillenbrand, Getty, Clark, & Wheeler, 1995). The length of the
51 vocal tract also changes the location (spacing) of the formants in lawful ways (Turner et al.,
52 2009). The length and tension of the vocal folds can be controlled by the talker; for example, f_0
53 contour differs between statements and questions (Eady & Cooper, 1986) and instantaneous f_0
54 fluctuates throughout a sentence when a talker speaks emotively (Bänziger & Scherer, 2005).
55 Nevertheless, average f_0 and formant spacing both differ reliably between different people, due
56 to physical constraints, and are informative about the gender (Titze, 1989) and size (Smith et
57 al., 2005) of a talker.

58 These two cues (f_0 and formant spacing) also contribute to listeners' judgements of
59 talker identity. They both influence the perceived similarity of unfamiliar talkers (f_0 : Baumann &
60 Belin, 2009; Gaudrain, Li, Ban, & Patterson, 2009; Matsumoto, Hiki, Sone, & Nimura, 1973;
61 Murry & Singh, 1980; Walden, Montgomery, Gibeily, Prosek, & Schwartz, 1978; formant
62 spacing: Baumann & Belin, 2009; Gaudrain et al., 2009; Matsumoto et al., 1973; Murry & Singh,
63 1980). In addition, they allow listeners to recognize familiar people from their voices (f_0 :
64 Abberton & Fourcin, 1978; LaRiviere, 1975; Lavner, Gath, & Rosenhouse, 2000; Lavner,
65 Rosenhouse, & Gath, 2001; van Dommelen, 1987, 1990; formant spacing: LaRiviere, 1975;
66 Lavner et al., 2000, 2001). Lavner et al. (2000) found that changing formant positions or f_0
67 reduced familiar-talker recognition, but recognition was more greatly affected by changes to
68 formant positions than by changes to f_0 —thus suggesting that vocal tract features contribute
69 more than glottal source features to familiar-talker recognition. This previous work is specific to
70 the acoustic cues that allow listeners to recognize talkers as familiar; the acoustic cues that
71 allow listeners to find familiar voices more intelligible have not been explored. Given that brain
72 activity differs when participants attend to speech content or the speaker's identity (von

73 Kriegstein et al., 2005), it seems plausible that the acoustic cues that underlie the speech-
74 intelligibility benefit for familiar voices may be different to those underlying recognition.

75 We recruited pairs of participants who had known each other for 6 months or longer. We
76 used a closed-set (rather than open-set) task to assess speech intelligibility, so that differences
77 between familiar and unfamiliar voice conditions could not be attributed to a difference in the
78 tendency to guess when uncertain. Each participant recorded sentences from the “BUG” speech
79 corpus (Kidd, Best, & Mason, 2008), where every sentence is of the form ““<Name> <verb>
80 <number> <adjective> <noun>” (e.g., “Bob bought five green bags”). We investigated whether
81 manipulating the acoustic correlates of glottal pulse rate (i.e., f_0) or of vocal tract length (VTL;
82 i.e. formant spacing) reduced the ability to recognise the voice as familiar and/or the speech-
83 intelligibility benefit gained from a familiar compared to unfamiliar target talker in the presence of
84 a competing talker.

85 **Methods**

86 ***Participants***

87 We recruited 11 pairs of participants (7 male, 15 female) who had known each other for
88 0.5–9.0 years (median = 2.0 years, interquartile range = 1.5) and who spoke regularly (> 5
89 hours per week). Pairs of participants were friends or couples. Seven were opposite-sex pairs
90 and three were same-sex (female-female) pairs. Twenty-one participants completed the entire
91 experiment. This sample size is sufficient to detect within-subjects effects of size $f = 0.41$ with
92 0.95 power (Faul et al., 2007); Johnsrude et al. (2013) reported a familiar-talker benefit to
93 speech intelligibility of size $f = 0.72$, which should be detectable with the current sample. The 21
94 participants were aged 19–24 years (median = 22.5 years, interquartile range = 2.6) and were
95 native Canadian English speakers who reported no history of hearing difficulty. Participants had
96 average pure-tone hearing levels of 15 dB HL or better in each ear (at four octave frequencies

97 between 0.5 and 4 kHz). The experiment was cleared by Western University's Health Sciences
98 Research Ethics Board. Informed consent was obtained from all participants.

99 ***Apparatus***

100 The experiment was conducted in a single-walled sound-attenuating booth (Eckel
101 Industries of Canada, Ltd.; Model CL-13 LP MR). Participants sat in a comfortable chair facing a
102 24-inch LCD visual display unit (either ViewSonic VG2433SMH or Dell G2410t).

103 Acoustic stimuli were recorded using a Sennheiser e845-S microphone connected to a
104 Steinberg UR22 sound card (Steinberg Media Technologies). During the listening tasks,
105 acoustic stimuli were presented through the Steinberg UR22 sound card (Steinberg Media
106 Technologies) and were delivered binaurally through Grado Labs SR225 headphones.

107 ***Stimuli***

108 Each participant recorded 480 sentences from the Boston University Gerald (BUG)
109 corpus (Kidd et al., 2008), which follow the structure: "<Name> <verb> <number> <adjective>
110 <noun>". In the sub-set used in the experiment, there were two names ('Bob' and 'Pat'), eight
111 verbs ('bought', 'found', 'gave', 'held', 'lost', 'saw', 'sold', 'took'), eight numbers ('two', 'three',
112 four', 'five', 'six', 'eight', 'nine', 'ten'), eight adjectives ('big', 'blue', 'cold', 'hot', 'new', 'old', 'red',
113 'small'), and eight nouns ('bags', 'cards', 'gloves', 'hats', 'pens', 'shoes', 'socks', 'toys'). An
114 example is "Bob bought three blue bags". To ensure that all sentences were spoken at similar
115 rates—and thus the five words from two different sentences would overlap when used in the
116 speech intelligibility task—we played videos indicating the desired pace for each sentence
117 (Holmes, 2018) while participants completed the recordings. The sentences had an average
118 duration of 2.5 seconds ($s = 0.3$). The levels of the digital recordings of the sentences were
119 normalised to the same root mean square (RMS) power.

120 Sentences were processed using the 'Change Gender' function in Praat (Boersma &
121 Weenink, 2013). Fundamental frequency (f_0) was changed by shifting the 'median pitch' of the

122 sentence upwards. Changes in vocal tract length (VTL) were simulated by shifting the
123 frequencies of the formants upwards by a percentage, which also increases their spacing. We
124 created 'unshifted' versions by shifting the median pitch and formants upwards, then downwards
125 again by the same amount, to restore the median pitch and formant positions of the original
126 sentence. The reason for creating 'unshifted' versions was to preserve any distortions
127 introduced by the signal processing, but maintain the original f_0 and formant values.

128 We aimed to manipulate f_0 and VTL by approximately the same perceptual amount, so
129 that any differences in the extent to which the two attributes influenced task performance was
130 not due to differences in perceptual discriminability of the two cues. To this aim, we estimated
131 listeners' thresholds for discriminating f_0 and VTL and used a multiple of this just-noticeable-
132 difference threshold in the main experiment. We wanted to make the manipulations large, so we
133 multiplied the median threshold (across participants) by 5, which was the largest manipulation
134 possible before the sentences became distorted by the signal processing algorithm. We
135 estimated the thresholds for discriminating changes to f_0 and VTL in a group of 5 participants
136 who did not take part in the main experiment. These participants performed a two-alternative
137 forced-choice (2AFC) task with a weighted (9:1) up-down adaptive procedure (Kaernbach,
138 1991) that estimated the 90% threshold for discriminating f_0 and VTL manipulations of the
139 familiar voice (i.e., the participant's partner's voice). On each trial, participants heard three
140 different sentences spoken by their partner's voice, presented sequentially. The first sentence
141 was presented with the original f_0 and VTL (unshifted version). Either the second or third
142 sentence was the manipulated version and the remaining sentence was unshifted, like the first
143 sentence. Participants indicated whether the second or third sentence was manipulated. We
144 used separate, but interleaved, runs for f_0 and VTL, each with a starting manipulation value of
145 1.15% above the original recording. The procedure stopped after 8 reversals and threshold
146 values were calculated as the median of the last 5 reversals (f_0 : 8.05%; VTL: 5.35%). We set the
147 manipulation magnitude at five times the median threshold from the group of 5 participants,

148 which produced stimuli with median pitches (corresponding to f_0) that were 40.25% higher than
149 that of the original sentences and sentences with formant frequencies (corresponding to VTL)
150 that were 26.75% higher than those of the original sentences. We refer to these stimuli as f_0 -
151 manipulated and VTL-manipulated stimuli, respectively. We created ‘both-manipulated’
152 sentences by shifting median pitch by 40.25% and formants by 26.75%.

153 During the experiment, each participant heard sentences spoken by their familiar partner
154 and sentences spoken by two unfamiliar talkers, who were the partners of other participants in
155 the experiment, sex matched to the familiar talker. The advantage of this aspect of the design
156 was that acoustic stimuli were counterbalanced across the familiar and unfamiliar voice
157 conditions; so that, across the group, these two types of condition were acoustically as similar
158 as possible. Each voice was presented to one participant (i.e. their partner) as a familiar talker
159 and to two other participants as an unfamiliar talker. The only exception was the participant
160 whose partner did not complete the experiment. This voice was presented as unfamiliar twice,
161 but never as familiar. For the same reason, two other voices were presented once as familiar
162 and only once as unfamiliar.

163 ***Procedure***

164 Participants completed two tasks: a speech intelligibility task and an explicit recognition
165 task. Half completed the speech intelligibility task first and the other half completed the explicit
166 recognition task first. Each task included three voice-manipulation conditions: (1) the original f_0
167 and VTL were preserved (unshifted condition), (2) f_0 was manipulated (f_0 -manipulated
168 condition), (3) VTL was manipulated (VTL-manipulated condition), and (4) f_0 and VTL were both
169 manipulated in combination (both-manipulated condition).

170 In the speech intelligibility task, participants heard two sentences spoken simultaneously
171 by different talkers. They identified the four remaining words of a sentence that began with a
172 particular target name (“Bob” or “Pat”), by clicking buttons on a screen. On each trial, either the
173 target sentence was spoken by the participant’s partner and the masker sentence was spoken

174 by an unfamiliar talker (“Familiar Target” condition), or both sentences were spoken by
175 unfamiliar talkers (“Both Unfamiliar” condition). The target and masker sentences were always
176 spoken by different talkers but were both manipulated in the same way (i.e. VTL-manipulated,
177 f_0 -manipulated, both-manipulated, or unshifted). Target and masker sentences were presented
178 at two different target-to-masker ratios (TMRs): -6 and +3 dB. For all participants, acoustic
179 stimuli were presented at a comfortable listening level (approximately 67 dB(A) SPL), which was
180 roved over a range of 3 dB. All trial types (2 familiarity conditions x 4 manipulation conditions x 2
181 TMRs) were randomly interleaved. Participants completed 768 trials (i.e., 32 trials in each
182 condition), with a short break every 64 trials and a longer break after 384 trials, after which the
183 target name word (i.e. “Bob” or “Pat”) was switched.

184 In the explicit recognition task, listeners heard one sentence on each trial. The sentence
185 could be spoken by the participant’s partner or by one of the two unfamiliar voices. We used the
186 same four voice manipulations as in the speech intelligibility task (VTL-manipulated, f_0 -
187 manipulated, both-manipulated, or unshifted). Participants were told that some of the sentences
188 had been manipulated and were instructed to report whether they thought each sentence was
189 spoken by their partner or not, regardless of any manipulation. Participants completed 84 trials
190 (21 for each manipulation condition).

191 At the end of the experiment, we checked that participants could accurately discriminate
192 between sentences that had been manipulated in f_0 and/or correlates of VTL and sentences in
193 which the original f_0 and correlates of VTL had been preserved. On each trial, participants heard
194 three different sentences spoken by their partner, presented sequentially. On each trial, all three
195 sentences were spoken by either the familiar talker or one of the two unfamiliar talkers. The first
196 sentence was always presented in its ‘unshifted’ version, as a reference. Of the two remaining
197 sentences, one was the manipulated version and the other was the ‘unshifted’ version. In a
198 2AFC task, participants had to indicate whether the second or third sentence had been

199 manipulated. Participants completed 48 trials, with 16 in each of the three manipulation
200 conditions (VTL-manipulated, f_0 -manipulated, or both-manipulated).

201 **Analyses**

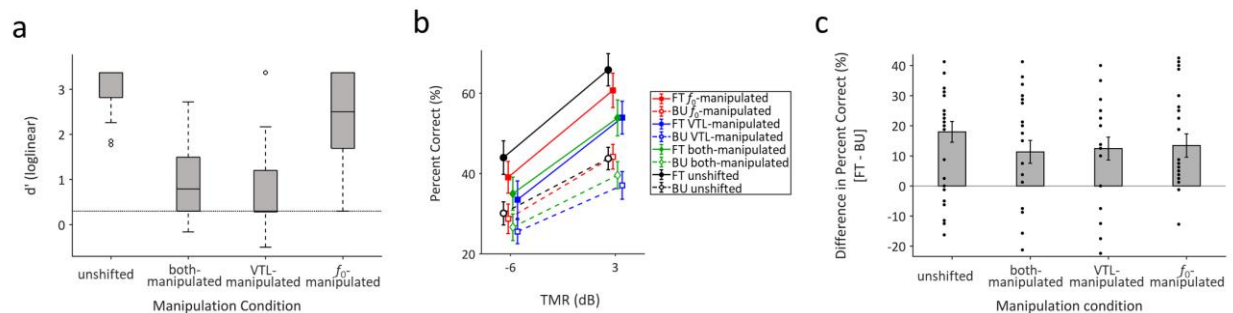
202 We calculated sensitivity (d') for the explicit-recognition data using loglinear correction
203 (Hautus, 1995), so chance d' is 0.3. For the speech intelligibility task, we calculated the
204 percentage of sentences in which participants reported all four words (after the name) correctly.

205 To assess the familiar-talker benefit to speech intelligibility, we compared percent correct
206 between the Familiar Target and Both Unfamiliar conditions. In both conditions, participants had
207 to report words from a target sentence in the presence of a masker sentence that was spoken
208 by a different (unfamiliar) talker. The masker voices were identical in the two conditions—the
209 only difference between these two conditions was whether the target sentence was spoken by a
210 familiar talker or by one of the unfamiliar talkers. We also analysed whether performance on the
211 speech intelligibility and explicit recognition tasks were affected by the manipulation condition
212 (VTL-manipulated, f_0 -manipulated, both-manipulated, or unshifted).

213 To assess whether there was a relationship between recognition performance and
214 speech-intelligibility benefit (e.g. to assess whether there is a greater intelligibility benefit for
215 voices that are better recognized), we calculated Spearman's rank correlation coefficients
216 between performance in the explicit recognition task and the magnitude of the speech-
217 intelligibility benefit for the familiar voice (i.e., the difference in percent correct between the
218 Familiar Target and Both Unfamiliar conditions). We did this separately for each manipulation
219 condition.

220 **Results**

221 Results from the manipulation discrimination task showed that participants could
222 discriminate changes in f_0 (mean [\bar{x}] = 91.6%, standard deviation [s] = 18.5), VTL (\bar{x} = 95.9%, s
223 = 18.2), and both cues



224

225 **Fig 1.** Explicit recognition and speech intelligibility (N=21). (a) Sensitivity (d') in the Explicit
 226 Recognition task. Open circles illustrate data from participants who were outliers. (b) Percent of
 227 trials in which participants reported the words from the target sentence correctly in the Speech
 228 Intelligibility task (c) Familiar-voice benefit (i.e. difference in percent correct between Familiar
 229 Target and Both Unfamiliar conditions), collapsed across target-to-masker ratios, in the Speech
 230 Intelligibility task. Error bars represent ± 1 standard error of the mean. Filled circles display
 231 results from individual participants. See the Results section for a description of significant
 232 differences between conditions. FT = Familiar Target; BU = Both Unfamiliar.

233

234

235 combined ($\bar{x} = 94.7$, $s = 22.3$) with high accuracy. One participant achieved below-chance
 236 performance (12.5%) on the discrimination task, but performed similarly to the other participants
 237 in the explicit recognition and speech intelligibility tasks, so we included this participant in the
 238 analyses (excluding this participant did not affect the pattern of results).

239 **Explicit recognition**

240 As shown in Figure 1a, sensitivity (d') in the explicit recognition task depended strongly
 241 on condition. Sensitivity was much lower in VTL-manipulated and both-manipulated conditions
 242 than in the unshifted and f_0 -manipulated conditions. The d' data violated the assumption of
 243 normality (skewed distributions and $p < .05$ in Shapiro-Wilk test), so non-parametric tests are
 244 reported.

245 We compared d' across the four manipulation conditions using Wilcoxon signed-rank
246 tests. Participants were significantly better at recognizing their partner's voice in the unshifted
247 condition compared to all others ($Z \geq 2.67, p \leq .008$). They were also better in the f_0 -manipulated
248 condition than in both conditions in which VTL was manipulated ($Z \geq 3.62, p < .001$). Sensitivity
249 (d') did not differ between the two conditions in which VTL was manipulated (VTL-manipulated
250 and both-manipulated; $Z = .71, p = .48$).

251 Sign tests, evaluating d' scores against chance level (0.3), showed that participants were
252 unable to recognize their partner's voice (i.e., chance sensitivity) in the two VTL-manipulated
253 conditions (VTL-manipulated: $S = 8, p = .38$; both-manipulated: $S = 13, p = .38$) but were
254 significantly better than chance in the unshifted ($S = 21, p < .001$) and f_0 -manipulated ($S = 18, p$
255 $= .001$) conditions.

256 To investigate whether the manipulations affected recognition differently for male and
257 female voices we conducted a 2x4 Mixed ANOVA (Sex x Manipulation). We found no main
258 effect of voice sex [$F(1, 19) = 1.13, p = .30, \omega = .01$] and no significant interaction between Sex
259 and Manipulation condition [$F(1, 19) = .26, p = .62, \omega = -.04$].

260 ***Speech intelligibility***

261 Baseline performance in the Both Unfamiliar condition was similar across the four
262 manipulation conditions (Figure 1b). Therefore, for each manipulation, we calculated the
263 familiar-voice speech-intelligibility benefit by subtracting percent correct in the Both Unfamiliar
264 condition from percent correct in the Familiar Target condition.

265 The data met the assumptions of normality, as assessed by the Shapiro-Wilk test and by
266 observing box-plots and Q-Q plots. We analyzed the data using a two-way within-subjects
267 ANOVA with the factors Manipulation (unshifted, f_0 -manipulated, VTL-manipulated, both-
268 manipulated) and TMR (-6, +3). The main effect of Manipulation was significant [$F(3, 60) = 3.69,$
269 $p = .017, \omega = .11$]. Planned comparisons showed that the familiar-voice benefit in the unshifted
270 condition was significantly larger than in all other conditions ($p \leq .036$). The familiar-voice benefit

271 did not differ significantly between any of the other conditions ($p \geq .31$). Participants received a
272 significantly greater familiar-voice benefit at +3 dB TMR ($\bar{x} = 10.1$, $s = 13.7$) than at -6 dB TMR
273 ($\bar{x} = 17.4$, $s = 19.6$) [$F(1, 20) = 9.17$, $p = .007$, $\omega = .27$]. The interaction between Manipulation
274 and TMR was not significant [$F(3, 60) = .24$, $p = .87$, $\omega = -.04$].

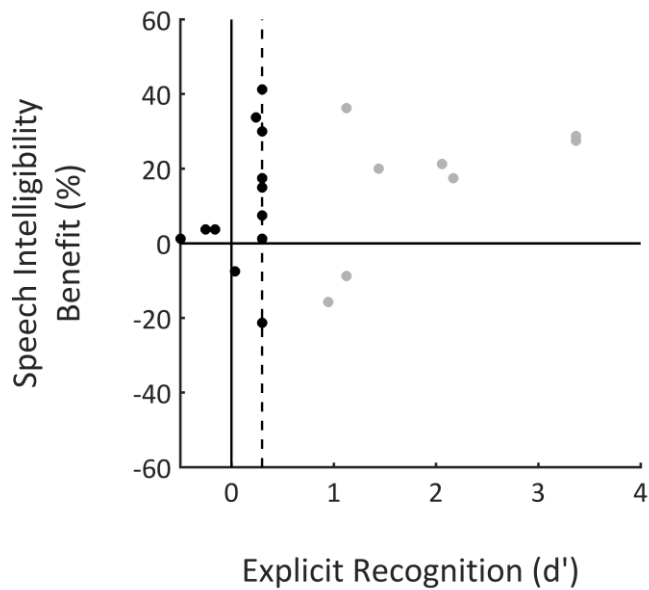
275 Figure 1c illustrates the familiar-voice benefit to speech intelligibility across the four
276 manipulations, collapsed across TMRs. One-sample t-tests for each manipulation showed that
277 the familiar-voice benefit was significantly greater than zero in all four conditions ($p \leq .007$).

278 We split the data by whether the voices were male or female and conducted a 2x4 (Sex
279 x Manipulation) Mixed ANOVA on the magnitude of the speech-intelligibility benefit for the
280 familiar voice. There was no main effect of voice sex [$F(1, 19) = 1.65$, $p = .21$, $\omega = .03$] and no
281 significant interaction between Sex and Manipulation [$F(1, 19) = 1.92$, $p = .18$, $\omega = .04$].

282 ***Voice manipulations affected recognition and intelligibility differently***

283 There was no significant relationship between recognition performance and the speech-
284 intelligibility benefit for any of the four manipulations ($r \leq .34$, $p \geq .13$). Thus, speech-intelligibility
285 benefit for a familiar voice does not appear to relate to the ability to explicitly recognize that
286 person from their voice.

287 To examine whether the pattern of results across manipulations differed significantly
288 between the speech-intelligibility and explicit-recognition tasks, we converted d' from the explicit
289 recognition task and percent improvement in speech intelligibility from the familiar talker into z-
290 scores and entered the data into a 2-way within-subjects ANOVA. We tested the two-way
291 interaction between Task (speech intelligibility and explicit recognition) and Manipulation
292 (unshifted, f_0 -manipulated, VTL-manipulated, and both-manipulated). The interaction was
293 significant [$F(3, 60) = 35.35$, $p < .001$, $\omega = .62$], confirming that the pattern across manipulations
294 indeed differed between the two tasks.



295

296 **Fig 2.** VTL-manipulated condition: Relationship between explicit recognition d' and the
 297 magnitude of the speech-intelligibility benefit for the familiar voice (i.e., Familiar Target – Both
 298 Unfamiliar). The vertical dashed line indicates chance performance ($d' = 0.3$) in the explicit
 299 recognition task. Each point illustrates one participant. Points that are coloured in black
 300 represent participants who scored at or below chance level in the explicit recognition task for the
 301 VTL-manipulated condition.

302

303

304 To further examine whether participants were able to gain a speech-intelligibility benefit
 305 from distorted voices that they were not able to explicitly recognize, we selected a sub-set of
 306 participants ($N = 13$) whose sensitivity was at or below chance ($d' \leq 0.3$) in the VTL-manipulated
 307 condition of the explicit recognition task (Figure 2). We performed a sign test for these 13
 308 participants to determine whether the speech-intelligibility benefit for the VTL-manipulated
 309 familiar voice differed from zero. Indeed, these participants gained a speech-intelligibility benefit
 310 for the VTL-manipulated familiar voice that was significantly greater than zero (median = 7.50%,

311 $S = 11, p = .022$). This result demonstrates that participants are able to gain a speech-
312 intelligibility benefit from a distorted familiar voice, even when they are not able to explicitly
313 recognize that voice as familiar.

314 **Discussion**

315 When the acoustic correlates of VTL were manipulated (27% shift in formant
316 frequencies), participants could no longer recognize a familiar voice, but still found it more
317 intelligible than sex-matched unfamiliar voices. In contrast, when f_0 was manipulated (shifted by
318 40%) participants could still recognize the familiar voice as well as finding it more intelligible.
319 Importantly, the patterns of results for these two manipulations differed significantly from each
320 other, to the point that participants who were unable to recognize the VTL-modified familiar
321 voice still found it more intelligible than unfamiliar voices. Thus, the two abilities rely on (at least
322 partially) distinct cognitive (and possibly neural) substrates. If you are using voice acoustics to
323 recognize someone you know, VTL information seems to be much more important than pitch
324 information. If, however, you are using voice acoustics to understand a familiar talker better,
325 pitch and VTL information play a partial role, but neither are critical.

326 In the face-recognition literature, a distinction has been drawn between identity and
327 expression processing (for a review, see Calder & Young, 2005). Patients with prosopagnosia
328 are able to identify emotional expressions in faces, despite impaired recognition of facial identity
329 (Humphreys et al., 1993). Similarly, patient studies have revealed a double dissociation
330 between voice-identity processing and speech processing (e.g., Van Lancker & Canter, 1982).

331 The 'auditory face' model (Belin et al., 2004), which is based on an influential model of
332 face perception (Bruce & Young, 1986), has been used to describe voice perception. This
333 model suggests that voice perception is multi-dimensional, with different systems specialised for
334 identity, speech recognition and emotional expression identification. The dissociation between
335 explicit recognition and the speech-intelligibility benefit in the current study is intriguing, because

336 it predicts that patients who are impaired in their ability to recognize voices might still find
337 familiar voices more intelligible when they are masked by a competing talker. Our results are
338 consistent with the idea that familiar-voice information may feed into (at least partially) separate
339 voice recognition and speech analysis systems.

340 The acoustic correlates of VTL appear to be critical for explicit recognition, whereas f_0
341 contributes to a lesser extent. This finding is consistent with the results of other studies that
342 compared the contributions of f_0 and VTL to explicit recognition (Lavner et al., 2000; Gaudrain et
343 al., 2009). The current results extend those previous findings by showing that the greater
344 influence of acoustic correlates of VTL on voice recognition cannot be explained by differences
345 in perceptual discriminability of the two sets of acoustic features. We approximately equated the
346 discriminability of the manipulations by selecting manipulation magnitudes from discrimination
347 (just-noticeable difference) thresholds in a separate group of participants. Thus, we conclude
348 that recognition of a voice as familiar is more robust to perceived differences in f_0 than to
349 perceived differences in correlates of VTL. Gaudrain et al. (2009) speculate that greater within-
350 talker variation in f_0 than VTL could explain the smaller contribution of f_0 to talker recognition.
351 Here, the average within-talker variability was 39.30% ($s = 21.19$) for f_0 and 0.39% ($s = 0.06$) for
352 formant spacing. The majority ($N = 12$) of the talkers had f_0 ranges less than our f_0 manipulation
353 of 40.25%, whereas all had formant spacing ranges substantially less than our formant
354 manipulation of 26.75%. Thus, based on our recorded sentences, it seems plausible that
355 differences in within-talker variability explains the greater effect of the VTL than the f_0
356 manipulation on recognition.

357 Although the VTL manipulation eliminated the ability to recognize a voice as familiar, it
358 did not eliminate the ability to gain a speech-intelligibility benefit from the familiar voice.
359 Manipulating f_0 and acoustic correlates of VTL decreased speech intelligibility (compared to the
360 unshifted condition) similarly. There was no additional decrement when both cues were
361 manipulated together compared to when f_0 or VTL were manipulated alone. It is important for

362 the interpretation of our results that speech intelligibility in the Both Unfamiliar condition was
363 similar across the manipulations (see Figure 1b), meaning that the baselines used to calculate
364 the familiar-voice benefit were at a similar place on the psychometric function for all
365 manipulation conditions. Thus, the difference in the familiar-target benefit to intelligibility is real,
366 rather than an artifact of differences in baseline performance.

367 The manipulations we used were as large as we could impose without distorting the
368 recordings, and were almost as large as the average difference between male and female
369 voices (Titze, 1989). Given that even these manipulations failed to eradicate the intelligibility
370 difference, listeners must rely on acoustic information other than average f_0 and the formant
371 ratio to better understand speech spoken by a familiar talker when a competing talker is
372 present. For example, f_0 contour, formant patterns, harmonic-to-noise ratio, intonation, and
373 rhythm might be important for the familiar-talker benefit to intelligibility. However, the same cues
374 were present in the VTL-manipulated stimuli in the explicit recognition task, and participants
375 performed at chance. Therefore, these cues are not sufficient for recognizing a voice as familiar.

376 In a separate group of participants ($N = 18$), we repeated the experiment using smaller
377 manipulations of f_0 and acoustic correlates of VTL. For each listener, we manipulated f_0 and
378 acoustic correlates of VTL at the listener's 90% threshold for discriminating manipulations to
379 those cues (i.e., manipulations were shifts of one just-noticeable difference unit, not five; the
380 range of thresholds were 1.7–6.3% for VTL and 3.9–9.9% for f_0). Although these manipulations
381 were perceptually discriminable (by definition), we found no effect of the manipulations on the
382 ability to recognize the voice as familiar or on the magnitude of the speech-intelligibility benefit
383 for the familiar voice. This result demonstrates that larger deviations to a familiar voice are
384 required to reduce explicit recognition and the speech-intelligibility benefit for familiar voices.

385 Across both experiments, we replicated the familiar-voice benefit to speech intelligibility
386 (Domingo et al., submitted; Johnsrude et al., 2013; Kreitewolf et al., 2017; Levi et al., 2011;
387 Nygaard & Pisoni, 1998; Nygaard et al., 1994; Yonan & Sommers, 2000) when the original f_0

388 and information about the original VTL of the familiar voice was preserved. The familiar-voice
389 intelligibility benefit is similar in magnitude in the current experiments (10–25%) as Johnsrude et
390 al. (2013) found for spouses' voices (10–20%), which is consistent with recent data indicating
391 that even 6 months of experience with a friend or partner's voice is sufficient to yield a large
392 intelligibility benefit (Domingo et al., submitted).

393 Overall, our results demonstrate a large improvement in speech intelligibility when
394 participants listened to a friend's voice in the presence of a competing talker than when they
395 listened to a stranger's voice. This benefit was relatively robust to large manipulations of f_0 and
396 acoustic correlates of VTL. Indeed, participants gained an intelligibility benefit from a
397 manipulated familiar voice even when they were no longer able to explicitly recognize that voice
398 as familiar. The findings demonstrate a dissociation between explicit recognition of a familiar
399 voice and the speech-intelligibility benefit gained from a familiar voice in the presence of a
400 competing talker. The findings imply that different mechanisms may be involved in processing
401 familiar-voice information, depending on the context in which the information is used.

402

403

Acknowledgements

404 This work was supported by funding from the Canadian Institutes of Health Research
405 (CIHR; Operating Grant: MOP 133450) and the Natural Sciences and Engineering Research
406 Council of Canada (NSERC; Discovery Grant: 327429-2012). We would like to thank Grace To
407 and Shivaani Shanawaz for assisting with data collection.

408

409

Author Contributions

410 E.H. and I.S.J. designed the research. E.H. and Y.D. collected the data. E.H. analysed
411 the data. E.H., Y.D., and I.S.J. wrote the paper.

412

Declaration of Conflicting Interests

413 The authors declare no conflicts of interest with respect to the authorship or the
414 publication of this article.

415

416

References

417 Abberton, E., & Fourcin, A. J. (1978). Intonation and speaker identification. *Language and*
418 *Speech*, 21(4), 305–318.

419 Anderson, A. K., & Phelps, E. A. (2000). Expression without recognition: Contributions of the
420 human amygdala to emotional communication. *Psychological Science*, 11(2), 106–111.

421 <https://doi.org/10.1111/1467-9280.00224>

422 Bänziger, T., & Scherer, K. R. (2005). The role of intonation in emotional expressions. *Speech*
423 *Communication*, 46(3–4), 252–267. <https://doi.org/10.1016/j.specom.2005.02.016>

424 Baumann, O., & Belin, P. (2009). Perceptual scaling of voice identity: Common dimensions for
425 different vowels and speakers. *Psychological Research*, 74(1), 110–120.

426 <https://doi.org/10.1007/s00426-008-0185-z>

427 Belin, P., Fecteau, S., & Bédard, C. (2004). Thinking the voice: Neural correlates of voice
428 perception. *Trends in Cognitive Sciences*, 8(3), 129–135.

429 <http://doi.org/10.1016/j.tics.2004.01.008>

430 Boersma, P., & Weenink, D. (2013) Praat: doing phonetics by computer [Computer program].
431 Version 5.4.04, retrieved 26 January 2015 from <http://www.praat.org/>

432 Bruce, V., & Young, A. (1986). Understanding face recognition. *British Journal of Psychology*,
433 77, 305–327.

434 Calder, A. J., & Young, A. W. (2005). Understanding the recognition of facial identity and facial
435 expression. *Nature Reviews Neuroscience*, 6(8), 641–651. <http://doi.org/10.1038/nrn1724>

436 Chiba, T., & Kajiyama, M. (1941). The vowel, its nature and structure, Tokyo-Kaiseikan Pub Co.,

437 Tokyo.

438 Domingo, Y., Holmes, E., & Johnsrude, I. S. (submitted). The benefit to speech intelligibility of
439 hearing a familiar voice.

440 Eady, S. J., & Cooper, W. E. (1986). Speech intonation and focus location in matched
441 statements and questions. *The Journal of the Acoustical Society of America*, 80(2), 402–
442 15.

443 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G*Power 3: A flexible statistical power
444 analysis program for the social, behavioral, and biomedical sciences. *Behavior Research*
445 *Methods*, 39, 175-191.

446 Fant, G. (1960). *Acoustic Theory of Speech Production*. Netherlands: The Hague.

447 Gaudrain, E., Li, S., Ban, V. S., & Patterson, R. D. (2009). The role of glottal pulse rate and
448 vocal tract length in the perception of speaker identity. *Proceedings of the Annual*
449 *Conference of the International Speech Communication Association, INTERSPEECH*,
450 148–151.

451 Halle, M. (1985). Speculations about the representation of words in memory. *Phonetic*
452 *Linguistics*.

453 Hautus, M. J. (1995). Corrections for extreme proportions and their biasing effects on estimated
454 values of d'. *Behavior Research Methods, Instruments, & Computers*, 27(1), 46–51.
455 <https://doi.org/10.3758/BF03203619>

456 Hillenbrand, J. M., Getty, L. A., Clark, M. J., & Wheeler, K. (1995). Acoustic characteristics of
457 American English vowels. *Journal of the Acoustical Society of America*, 97(5), 3099–3111.
458 <https://doi.org/10.1121/1.411872>

459 Holmes, E. (2018). Speech Recording Videos (Version v1.0.0). Zenodo.
460 <http://doi.org/10.5281/zenodo.1165402>

461 Humphreys, G. W., Donnelly, N., & Riddoch, M. J. (1993). Expression is computed separately
462 from facial identity, and it is computed separately for moving and static faces:

463 Neuropsychological evidence. *Neuropsychologia*, 31(2), 173–181.
464 [https://doi.org/10.1016/0028-3932\(93\)90045-2](https://doi.org/10.1016/0028-3932(93)90045-2)

465 Johnsruide, I. S., Mackey, A., Hakyemez, H., Alexander, E., Trang, H. P., & Carlyon, R. P.
466 (2013). Swinging at a cocktail party: voice familiarity aids speech perception in the
467 presence of a competing voice. *Psychological Science*, 24(10), 1995–2004.
468 <https://doi.org/10.1177/0956797613482467>

469 Joos, M. (1948). Acoustic Phonetics. *Language*, 24(Suppl. 2), 5–136.
470 <https://doi.org/10.2307/522229>

471 Kaernbach, C. (1991). Simple adaptive testing with the weighted up-down method. *Perception &*
472 *Psychophysics*, 49(3), 227–229. <https://doi.org/10.3758/BF03214307>

473 Kidd, G., Best, V., & Mason, C. R. (2008). Listening to every other word: examining the strength
474 of linkage variables in forming streams of speech. *The Journal of the Acoustical Society of*
475 *America*, 124(6), 3793–3802. <https://doi.org/10.1121/1.2998980>

476 Kreitewolf, J., Mathias, S. R., & von Kriegstein, K. (2017). Implicit talker training improves
477 comprehension of auditory speech in noise. *Frontiers in Psychology*, 8, 1584.
478 <https://doi.org/10.3389/fpsyg.2017.01584>

479 LaRiviere, C. (1975). Contributions of fundamental frequency and formant frequencies to
480 speaker identification. *Phonetica*, 31, 185–197. <https://doi.org/10.1159/000259668>

481 Lavner, Y., Gath, I., & Rosenhouse, J. (2000). Effects of acoustic modifications on the
482 identification of familiar voices speaking isolated vowels. *Speech Communication*, 30(1), 9–
483 26. [https://doi.org/10.1016/S0167-6393\(99\)00028-X](https://doi.org/10.1016/S0167-6393(99)00028-X)

484 Lavner, Y., Rosenhouse, J., & Gath, I. (2001). The prototype model in speaker identification by
485 human listeners. *International Journal of Speech Technology*, 4(1), 63–74.
486 <https://doi.org/10.1023/A:1009656816383>

487 Levi, S. V, Winters, S. J., & Pisoni, D. B. (2011). Effects of cross-language voice training on
488 speech perception: whose familiar voices are more intelligible? *The Journal of the*

489 *Acoustical Society of America*, 130(6), 4053–62. <https://doi.org/10.1121/1.3651816>

490 Matsumoto, H., Hiki, S., Sone, T., & Nimura, T. (1973). Multidimensional representation of
491 personal quality of vowels and its acoustical correlates. *IEEE Transactions on Audio and*
492 *Electroacoustics*, 21(5), 428–436. <https://doi.org/10.1109/TAU.1973.1162507>

493 Murry, T., & Singh, S. (1980). Multidimensional analysis of male and female voices. *The Journal*
494 *of the Acoustical Society of America*, 68(5), 1294–1300. <https://doi.org/7440851>

495 Nygaard, L. C., & Pisoni, D. B. (1998). Talker-specific learning in speech perception. *Perception*
496 *& Psychophysics*, 60(3), 355–376. <https://doi.org/10.3758/BF03206860>

497 Nygaard, L. C., Sommers, M. S., & Pisoni, D. B. (1994). Speech perception as a talker-
498 contingent process. *Psychological Science*, 5(1), 42–46.

499 Titze, I. R. (1989). Physiologic and acoustic differences between male and female voices.
500 *Journal of the Acoustical Society of America*, 85(4), 1699–1707.
501 <https://doi.org/10.1121/1.397959>

502 Turner, R. E., Walters, T. C., Monaghan, J. J. M., & Patterson, R. D. (2009). A statistical,
503 formant-pattern model for segregating vowel type and vocal-tract length in developmental
504 formant data. *The Journal of the Acoustical Society of America*, 125(4), 2374–2386.
505 <http://doi.org/10.1121/1.3079772>

506 van Dommelen, W. A. (1987). The contribution of speech rhythm and pitch to speaker
507 recognition. *Language and Speech*, 30(4), 325–338.
508 <https://doi.org/10.1177/002383098703000403>

509 van Dommelen, W. A. (1990). Acoustic parameters in human speaker recognition. *Language*
510 *and Speech*, 33(3), 259–272.

511 Van Lancker, D. R., & Canter, G. J. (1982). Impairment of voice and face recognition in patients
512 with hemispheric damage. *Brain and Cognition*, 1(2), 185–195.
513 [https://doi.org/10.1016/0278-2626\(82\)90016-1](https://doi.org/10.1016/0278-2626(82)90016-1)

514 von Kriegstein, K., Kleinschmidt, A., Sterzer, P., & Giraud, A.-L. (2005). Interaction of face and

515 voice areas during speaker recognition. *Journal of Cognitive Neuroscience*, 17(3), 367–
516 376. <https://doi.org/10.1162/0898929053279577>

517 Walden, B. E., Montgomery, A. A., Gibeily, G. J., Prosek, R. A., & Schwartz, D. M. (1978).
518 Correlates of psychological dimensions in talker similarity. *Journal of Speech, Language,*
519 *and Hearing Research*, 21, 265–275.

520 Yonan, C. A., & Sommers, M. S. (2000). The effects of talker familiarity on spoken word
521 identification in younger and older listeners. *Psychology and Aging*, 15(1), 88–99.
522 <https://doi.org/10.1037/0882-7974.15.1.88>