



PRIFYSGOL
BANGOR
UNIVERSITY

Assessing susceptibility to distraction along the vocal processing hierarchy

Muhl, Constanze; Bestelmeyer, Patricia

Quarterly Journal of Experimental Psychology

DOI:

[10.1177/1747021818807183](https://doi.org/10.1177/1747021818807183)

Published: 01/07/2019

Peer reviewed version

[Cyswllt i'r cyhoeddiad / Link to publication](#)

Dyfyniad o'r fersiwn a gyhoeddwyd / Citation for published version (APA):

Muhl, C., & Bestelmeyer, P. (2019). Assessing susceptibility to distraction along the vocal processing hierarchy. *Quarterly Journal of Experimental Psychology*, 72(7), 1657-1666. <https://doi.org/10.1177/1747021818807183>

Hawliau Cyffredinol / General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Assessing susceptibility to distraction along the vocal processing hierarchy

Short title: *Voice perception ability tests*

Constanze Mühl^{a*} & Patricia E.G. Bestelmeyer^a

^aSchool of Psychology, Bangor University, LL57 2AS, UK

*Correspondence:

Constanze Mühl

School of Psychology

Bangor University

Brigantia Building

Penrallt Road

Bangor, Gwynedd

LL57 2AS

UK

E-mail: constanze.muhl@bangor.ac.uk

Phone: +44 (0)1248 383254

Abstract

Recent models of voice perception propose a hierarchy of steps leading from a more general, “low-level” acoustic analysis of the voice signal to a voice-specific, “higher-level” analysis. We aimed to engage two of these stages: First, a more general detection task in which voices had to be identified amidst environmental sounds, and, second, a more voice-specific task requiring a same/different decision about unfamiliar speaker pairs (Bangor Voice Matching Test, BVMT). We explored how vulnerable voice recognition is to interfering distractor voices, and whether performance on the aforementioned tasks could predict resistance against such interference. Additionally, we manipulated the similarity of distractor voices to explore the impact of distractor similarity on recognition accuracy. We found moderate correlations between voice detection ability and resistance to distraction ($r = .44$), and BVMT and resistance to distraction ($r = .57$). A hierarchical regression revealed both tasks as significant predictors of the ability to tolerate distractors ($R^2 = .36$). The first stage of the regression (BVMT as sole predictor) already explained 32% of the variance. Descriptively, the “higher-level” BVMT was a better predictor ($\beta = .47$) than the more general detection task ($\beta = .25$), although further analysis revealed no significant difference between both beta weights. Furthermore, distractor similarity did not affect performance on the distractor task. Overall, our findings suggest the possibility to target specific stages of the voice perception process. This could help explore different stages of voice perception and their contributions to specific auditory abilities, possibly also in forensic and clinical settings.

Keywords: voice perception, voice detection, voice recognition

23 Successful social interaction relies on our capacity to extract relevant information
24 from our surroundings and the people with whom we are interacting. While there is an
25 extensive amount of research into the perception of such cues from faces, the perception of
26 these cues from voices has been neglected until recently (Blank, Wieland, & von Kriegstein,
27 2014; Gainotti, 2014). Theoretical models of voice perception closely follow those already
28 established for face perception but have received little empirical evaluation. Belin and
29 colleagues suggest a voice perception model adapted from Bruce and Young’s (1986) model
30 of familiar face perception (Belin, Fecteau, & Bédard, 2004). This voice perception model
31 proposes that after an initial low-level analysis of the voice signal, a number of different
32 independent modules are responsible for the analysis of vocal speech, vocal affect, and
33 speaker identity information, before additional semantic knowledge about a person is
34 accessed through the activation of Person Identity Nodes (Belin et al., 2004; Campanella &
35 Belin, 2007). This proposal suggests that the independent levels and modules can be
36 investigated separately.

37 An alternative model by Kreiman and Sidtis (2013) suggests that the recognition
38 process for voices relies simultaneously on the Gestalt perception of the whole (pattern
39 recognition) and the analysis of specific auditory cues within the voice (feature analysis). The
40 degree to which both are engaged depends on the familiarity of the voices. Recognition of
41 unfamiliar voices calls for the extraction of features more than for an overall pattern
42 recognition, possibly also involving comparison to a known “average” voice, and is more
43 stimulus-driven. Familiar voice recognition is more top-down in that it relies heavily on the
44 overall voice pattern, with only voice-identity specific features becoming salient throughout
45 recognition. As such, recognising an unfamiliar voice is a question of discriminating and
46 matching two voice signals, and is therefore often described as the ability of voice
47 discrimination. Recognising a familiar voice, in contrast, is the recognition of an overall vocal
48 pattern specific to a single person. The term “voice recognition” therefore often applies to the
49 recognition of voice identity for familiar speakers in particular (see also van Lancker &

50 Kreiman, 1987). Furthermore, a recent neuroimaging study with lesion patients
51 (Roswadowitz, Kappes, Obrig, & von Kriegstein, 2018) has also found that different brain
52 structures are involved in the perception of newly-learned unfamiliar vs. familiar voices, which
53 supports this distinction.

54 Although Kreiman and Sidtis' model does not indicate independent feature-specific
55 modules (e.g. for vocal affect perception) like Belin and colleagues' model does, it
56 nevertheless posits the involvement of several distinct brain regions. Tasks related to voice
57 perception therefore recruit the distributed areas that are relevant for solving a specific task.
58 Findings of distributed time scales, for example in vocal affect perception (Iredale, Rushby,
59 McDonald, Dimoska-Di Marco, Swift, 2013; see also model for vocal affect processing by
60 Schirmer & Kotz, 2006, and Bestelmeyer et al., 2014), suggests that voice perception
61 involves hierarchical stages. According to these, earlier stages represent more general
62 analyses, and in the case of unfamiliar voices possibly also more stimulus-driven analyses,
63 before voices are processed in a more abstract, integrative manner (e.g. Warren, Jennings,
64 & Griffiths, 2005; Schirmer & Kotz, 2006).

65 The need for research on this topic, and indeed support for the existence of different
66 independent voice perception modules, becomes more apparent when surveying the
67 diversity of clinical symptoms reported for individuals with phonagnosia, or an impairment in
68 voice perception. For example, an extensive study of patients with brain lesions revealed
69 that while most patients with voice recognition deficits (in this case the recognition of famous
70 familiar voices) were still able to discriminate between two different unfamiliar voices, one of
71 the patients showed an impairment in both (Neuner & Schweinberger, 2000). However, in
72 this sample no further tests were reported to see whether other domains of voice perception
73 like the perception of gender or affect were selectively impaired as well. In recent years,
74 cases of individuals with developmental phonagnosia have emerged. To assess the extent
75 of their voice recognition deficits, these individuals often complete a number of voice
76 perception tests that target specific voice perception abilities. Usually, only certain functions

77 of voice perception are impaired (e.g. identity perception), while others like gender
78 perception remain intact (see also the first reported case of developmental phonagnosia in
79 Garrido et al., 2009). Both acquired and developmental voice perception deficits underline
80 the need for a more in-depth assessment of possible singular processing stages in order to
81 establish the range of functions that can be selectively impaired.

82 Apart from clinical contexts and the focus on general perception mechanisms, voice
83 identity perception has also received attention in non-clinical contexts, particularly in the field
84 of forensic psychology. As Kreiman and Sidtis (2013) point out, recognising an unfamiliar
85 person by voice alone is not a task we often encounter in natural settings, yet witnesses to a
86 crime might only be exposed to a perpetrator's voice. The reliability of witness testimony
87 therefore depends on a witness's ability to extract identity information from a typically
88 unfamiliar voice (i.e. process and compare the features of that voice to a stored
89 representation of average voices) and store this information for the newly heard voice. Then,
90 at a later point, the witness needs to distinguish the initial target voice from other unfamiliar
91 voices (all of which require the same processing steps), and match it to its correct target at a
92 later voice line-up. In terms of Belin and colleagues' more general model of possible distinct
93 modules, this forensic line-up task requires structural encoding of the perpetrator's voice
94 beyond just low-level auditory processing. Ideally, identity-specific features of the target
95 voices also have to be accessible at a later time point to allow for correct identification of the
96 perpetrator. This process is, of course, prone to error (Legge, Grosman, & Pieper, 1984;
97 Yarmey, 1995), and studies on it are often tailored to match specific criminal cases, making
98 connections to existing, more general voice perception literature difficult (Kreiman & Sidtis,
99 2013).

100 Despite the ecological validity of such voice line-up tasks, more controlled, lab-based
101 experiments are necessary. A recent study by Stevenage and colleagues (2013) explored
102 the detrimental impact of interference on speaker perception. Listeners heard an unfamiliar
103 speaker articulating a single sentence. In a fixed 16 s interval, participants then heard either

104 nothing, or two or four distractor voices. This was followed by a test voice. Participants had
105 to decide whether this test voice was identical to the initial target voice or not. Accuracy on
106 this task was reduced as soon as any distractor voice was introduced. The detrimental effect
107 distractors had on overall task performance occurred both when the distracting voices were
108 similar (as defined by same speaker sex as target voice) or different (opposite speaker sex).

109 Our aim for the current study was, [on the one hand](#), to test two potentially separate
110 abilities that occur at different stages of voice perception. [On the other hand, we also](#)
111 [wanted](#) to explore their impact on a third, complex auditory task [that has been used](#)
112 [previously and in more ecologically valid contexts](#). The aforementioned potentially separate
113 abilities are first, the ability to detect voices as a discrete class of sound objects (voice
114 detection ability), and, second, the ability to determine whether two utterances were spoken
115 by the same speaker or not (voice matching ability). To investigate whether both are suitable
116 to determine the accuracy on a more complex auditory task, we chose a distractor task
117 [examining](#) how vulnerable or susceptible someone is to the interference of a distracting
118 voice. [This third task follows the example of voice perception tasks common in forensic](#)
119 [contexts \(same/different decisions about a voice that one had previously been exposed to,](#)
120 [following interfering information\)](#). However, for the current study this takes place within a lab-
121 [based environment, allowing for stricter control of voice variables. For this reason, we also](#)
122 [wanted to revisit the issue of distractor similarity, i.e. whether distractors that are either](#)
123 [similar or different from the initial target voice affect the accuracy of one's same/different](#)
124 [decision](#).

125 Voices are arguably the most salient sound in our environment. Although there is
126 some debate about the timescale of this development, several studies have reported that
127 infants already show preferential brain activation patterns for vocal sounds within the first
128 twelve months after birth (e.g. Blasi et al., 2011; Grossman, 2011; Cheng, Lee, Chen, Wang,
129 & Decety, 2012). Additionally, lesions studies have shown that voices are processed
130 independently of other object sounds (Peretz et al., 1994; Neuner & Schweinberger, 2000).

131 As such, the detection of voices should be part of the earlier processing stream of vocal
132 sounds (as described in Belin and colleagues' model). In our study we aimed to measure
133 participants' ability to detect voices in an ongoing stream of vocal and non-vocal sounds.
134 This task was inspired by a visual detection task for faces to investigate an individual with
135 severe face recognition impairments (prosopagnosia; Duchaine, Yovel, Butterworth, &
136 Nakayama, 2006). Our task was adapted to address the inherent differences between the
137 visual domain (faces) and the analysis of auditory information as it unfolds over time. While
138 Duchaine and colleagues embedded their target stimuli (faces) in a noisy background, we
139 chose an ongoing stream of auditory, undistorted stimuli.

140 To examine a later module of voice perception, we included the Bangor Voice
141 Matching Test (BVMT; Mühl, Sheil, Jarutytė, & Bestelmeyer, 2017). This task involves
142 listening to two different utterances and then deciding whether these stem from the same or
143 different speakers. It thereby requires the extraction of identity information from a voice
144 before making a same/different judgment. Belin et al.'s (2004) model proposes that voice
145 identity cues are processed after the structural configuration of a voice has been extracted.
146 In contrast, Kreiman and Sidtis' (2013) model proposes that for this particular task,
147 participants have to extract the features of both unfamiliar voices and then compare these to
148 a template of an average voice.

149 Both the voice detection task and the BVMT will be examined in conjunction with the
150 performance on a third task, a voice distractor task. Here, participants have to make an
151 old/new judgment following initial exposure to a target voice. Crucially, a distractor voice is
152 introduced between hearing the first target voice and the same/different judgment needed for
153 the second target voice. We propose that the complexity of this distractor task should require
154 both of the processing stages we aim to tap into using the detection task and the BVMT. The
155 voice detection task depends on an earlier perception stage in which the signal is processed
156 as a vocal (as opposed to a non-vocal) sound. The BVMT, on the other hand, requires a
157 more complex analysis of the vocal signal. In fact, we assume that the BVMT and the

158 distractor task require the extraction of the same kind of vocal cues (voice identity
159 information/feature-based processing and comparison to an average voice). This reflects the
160 proposed succession of voice perception modules in Belin and colleagues' model (2004).
161 We therefore predict that both the voice detection task and the BVMT should correlate with
162 the distractor task [as they all rely on the analysis of a sound as a vocal object](#), but that the
163 correlation with the BVMT should be higher. In order to complete the distractor task
164 accurately, both an intact ability to detect voices and an intact ability to extract identity cues
165 from voices are necessary. We therefore also expect that performance in the voice detection
166 task and in the BVMT will both be predictors for the performance in the distractor task.
167 However, given the proposed similar, later processing stages necessary for the BVMT and
168 distractor task, we assume that the BVMT will be a better predictor.

169 Finally, we plan to revisit the issue of distractor similarity as initially explored by
170 Stevenage and colleagues (2013). They chose an arguably lenient criterion for their
171 manipulation of vocal similarity as it was solely based on speaker sex. A more fine-tuned
172 approach to voice similarity (relative proximity vs. relative distance in voice space) will
173 determine whether we classify distractors as similar or different. It has been proposed that
174 we perceive different voice identities by comparing them to a prototypical, average voice
175 (Latinus & Belin, 2011; Lavner, Rosenhouse, & Gath, 2001). Specifically, the existence of a
176 two-dimensional voice space based on two acoustic parameters (fundamental frequency,
177 F0, and first formant frequency, f1) has been suggested. Different vocal identities are located
178 within this voice space according to their vocal characteristics. The closer two voices are
179 within this voice space, the more likely it is that they are judged to belong to the same
180 person (Baumann & Belin, 2010). Therefore, our prediction is that the closer a distractor
181 voice is in terms of physical voice distance (i.e. the more similar it is in its physical
182 characteristics to a given target voice), the more distracting it will be. [We chose this
183 particular design, including the similarity manipulation, to incorporate both the concept of
184 voice recognition after interfering information \(as in previous forensic studies\), and the](#)

185 increased control over the nature of the distracting information afforded by the lab-based
186 conditions.

187

188

Method

Participants

189 The sample consisted of 100 native-English speakers (25 male; $M_{age} = 21.2$,
190 $SD_{age} = 6.5$) who took part in exchange for course credit. All participants reported normal
191 hearing. Written informed consent was obtained from all participants. The study was
192 approved by the Ethics committee of the School of Psychology at Bangor University.
193

194

Stimuli and Materials

195 Voice recordings for both the Bangor Voice Matching Test and the distractor task
196 consisted of non-sense syllables (different combinations of vowels and consonants like 'aga'
197 or 'hed') spoken by young female and male British-English native speakers. Sounds were
198 recorded in a sound attenuated booth using Audacity (16-bit, 44.1 kHz sampling rate, mono).
199 All speakers were between 18 – 28 years of age. All test stimuli were root-mean square
200 normalised and edited in Cool Edit Pro to start with onset of phonation and end with the offset
201 of phonation (mean duration = .51s; S.D. = .11). For each speaker gender, the distance
202 between each individual speaker and every other speaker was calculated using Pythagoras
203 theorem. This distance was defined as the distance in a two-dimensional voice space between
204 F0 and F1 (see Baumann & Belin, 2010). The smaller this distance, the more similar the
205 speakers are perceived to sound (Baumann & Belin, 2010). For a more detailed explanation
206 of this concept, see Figure S4 in the supplementary online material (SOM). Further detail on
207 the audio recordings as well as selection of voice pairs is provided in the stimulus details
208 described in Mühl et al., 2017.
209

210 Voice Detection Task

211 For this task, a total of 144 high quality sounds were chosen from a number of
212 different sources, including the Multimodal Stimulus Set (Schneider, Engel, & Debener,
213 2008). Sounds belonged to one of three categories: (1) human vocalisations like laughter or
214 singing (72 sounds; 32 male, 32 female, 8 children's voices), (2) inanimate environmental
215 sounds like telephone ringing (36 sounds), or (3) animate environmental sounds like a cat
216 meowing (36 sounds). Each stimulus was edited to include a 10 ms ramp up and down at its
217 start and end, respectively, using Cool Edit Pro, version 2.00 to avoid clipping. Sounds were
218 then RMS normalised using Matlab (R2013a). To ensure sufficient task difficulty, several
219 pilot versions of the detection task were run with differing stimulus lengths between 75 ms
220 and 250 ms. To avoid ceiling or floor effects we decided on a stimulus duration of 150 ms
221 which revealed an average performance of 77.36% during pilot testing ($n = 8$).

222 In the main part of the experiment, participants listened to the 144 sounds described
223 above. These sounds were either presented to the right or left ear, to follow the structure of
224 the face detection task used in Duchaine et al. (2006) where an intact face, presented within
225 an array of detached facial features, had to be spotted either on the left or the right side of
226 the picture. Ear assignments of sounds were counterbalanced across participants.
227 Participants had to indicate via keypress in which ear a human sound appeared ('x' for left
228 ear, 'm' for right ear). No response was necessary for the environmental sounds.
229 Participants had 2 seconds to react before the next sound was presented. During stimulus
230 presentation, participants saw a fixation cross centred on the screen as well as a reminder of
231 the key assignments in the upper half of the screen. Test duration was roughly 7 minutes.

232

233 Bangor Voice Matching Test

234 The Bangor Voice Matching Test is a computerised voice matching test in which
235 participants make a same/different identity decision after hearing 2 different syllables per

236 trial. Syllables were either articulated by the same speaker (40 trials) or by two different
237 speakers (another 40 trials; for further details on item selection for the Bangor Voice
238 Matching Test see Mühl et al., 2017). Speaker sex was balanced, with half of the trials
239 presenting male or female speakers, respectively. Instructions were given on the screen and
240 testing was self-paced. For each trial, participants saw two red speaker icons on the screen
241 and, below them, two response boxes, one for same and one for different speakers. Clicking
242 on the speaker icons led to the audio for each item being played. Responses were then
243 given by clicking on either of the response boxes. Participants could listen to each item
244 multiple times if they wished. Between trials, participants saw a centred fixation cross for
245 800 ms. On average, completion of the BVMT took less than 10 minutes.

246

247 Distractor Task

248 For the distractor task, each trial consisted of 3 voices: a first target voice (T1)
249 followed by a distractor voice (D) which, in turn, was followed by a second target voice (T2).
250 Voices were separated by a 0.8 s interval. Speaker sex throughout each trial was consistent
251 with 32 trials presenting male speakers and 32 trials presenting female speakers (64 trials in
252 total). For half of the items for each speaker block (male/female), T1 and T2 were the same
253 speaker. For the other half, T1 and T2 speaker identity differed. These formed the
254 same/different items. For all of those items, T1-D combinations represented the voice pairs
255 mentioned above. Items were formed in such a way that T1-D distances were either small (<
256 .020), representing similar speakers, or large (between .204 and .936), representing
257 speakers that were not similar and thus more easily distinguishable. This was done to allow
258 for an analysis of whether the similarity of a distractor D influences the recognisability of a
259 target voice T1. Half of the 'same' items and half of the 'different' items presented small T1-D
260 distances. For all different items, similarity between T1 and T2 was also balanced so that
261 half of the 'different items' consisted of similar T1 and T2. Similarity between distractor
262 voices and Target 2 voices (D-T2 similarity) could not be fully balanced due to the limited

263 number of voice pairings available, and were therefore not considered in our predictions.
264 Nevertheless, we tried to keep the distribution of D-T2 distances comparable for male and
265 female trials with 13 small and 19 larger D-T2 distances each. All syllables uttered within an
266 item were different (e.g. aba – hed – ubu, and not aba – hed - aba), and T2 syllable type
267 (consonant-vowel-consonant or vowel-consonant-vowel) either matched only T1 syllable
268 type (13 items), D syllable type (13 items), both T1 and D (18 items), or was different to T1
269 and D (20 items).

270 Independent t-tests between the female and male voices that were used in the
271 distractor task revealed no significant difference between the mean T1-D distance overall
272 ($t[62] = -.068, p = .946$). Additionally, there was no significant difference between either
273 similar T1-D voice pairings for female and male speakers, $t(30) = -.681, p = .541$, or
274 different T1-D voice pairings for female and male speakers, $t(30) = -.087, p = .931$. The
275 same was the case when considering the D-T2 similarities instead (all $p > .602$).

276 Participants' task was to listen to the three voices per trial, and then decide whether
277 the first and the third speaker were the same or not. Decisions were made using the 'f' and 'j'
278 key for same or different voices (key assignment counterbalanced across participants). The
279 next trial started following a button press. During stimulus presentation, participants saw a
280 fixation cross in the centre of the screen. After the third voice (T2) had been played, the key
281 assignment was displayed on the upper half of the screen. Completion of this task took
282 about 20 minutes.

283

284 *Procedure*

285 All tasks were implemented in Psychtoolbox-3 (Brainard, 1997; Kleiner, Brainard, &
286 Pelli, 2007) for Matlab (R2013a). Stimuli were presented via Beyerdynamic DT770 Pro
287 headphones (250 Ω). Up to 2 participants were tested at the same time. The order of the
288 three tasks was randomised across all participants. After being given general information

289 about the nature of the experiments, participants filled in a consent form before starting the
290 tasks. Each task was introduced by the experimenter, and both spoken and written
291 instructions were provided. Both voice detection task and distractor task included practice
292 blocks (8 trials/4 trials, respectively). Stimuli presented in those practice trials were not used
293 in the main parts of the experiments. Moreover, participants were encouraged to ask
294 questions in case of uncertainty about a task. After completion of all three tasks, participants
295 were debriefed and given contact details in case of further questions.

296

297 *Data analysis and design*

298 Data was analysed using Matlab (R2013a) and SPSS (version 22). Performance in
299 detection and distractor tasks were calculated as sensitivity A' , using signal detection theory,
300 to control for possible response bias in tasks that require detection of a signal within noise.
301 Accuracy in percentage correct, where reported, were calculated based on the corrected hit
302 and miss rates for detection and distractor task. These calculations followed the steps
303 proposed in Stanislaw and Todorov (1999) for use in SPSS packages (see [equation SE1](#) in
304 the SOM). Only valid trials with reaction times over 250ms were included. Bivariate
305 Pearson's correlations were used to determine the relationship between all three tasks.
306 Following that, a hierarchical linear regression analysis was performed to understand
307 whether the general ability for voice matching (BVMT score) and performance in the
308 detection task predicted the performance in the distractor task. Finally, paired t-tests on the
309 overall percentage correct in the distractor task were used to determine whether the
310 similarity of distractor voices influences the similarity decision for T1 and T2.

311 Two participants were identified as outliers for their performance on the distractor
312 task (studentised residuals ± 3 SDs), and excluded from subsequent analysis to meet the
313 assumptions for the regression analysis. Sample size for both the hierarchical linear
314 regression and the t-tests was $N = 98$. Inclusion of both outliers did not affect conclusions.

315 [Supplementary Figures S2 and S3](#) further illustrate the standardised residuals of the
316 regression analysis.

317

318

Results

319

320

321

322

323

324

325

326

327 Table 1

328 *Descriptive statistics (% correct), and bivariate correlations (Pearson's r) for percentage*

329 *correct in BVMT, and A' measures for voice detection task and distractor task*

330

	<i>M</i>	<i>SD</i>	Correlation with Detection Task	Correlation with Distractor Task
BVMT	85.14	7.13	.399**	.570**
Detection	87.31	5.11	-	.437**
Distractor	77.67	7.78	-	-

331 *Note.* *N* = 98. BVMT = Bangor Voice Matching Test. *M* is mean, *SD* is standard deviation. ***p* < .001.

332

333

334

335

336 A two-stage hierarchical multiple regression analysis was calculated to predict the
337 overall accuracy score (A') in the distractor task based on performance on the BVMT (BVMT
338 score; voice-specific, “high-level” voice perception task) and on performance on the voice
339 detection task (A' ; more general, “low-level” voice perception task). At stage one,
340 performance on the voice matching task (BVMT score) served as a significant predictor for
341 accuracy in the distractor task, $F(1,96) = 46.30, p < .001$, adjusted $R^2 = .318$. The addition of
342 performance on a “low-level” voice perception task (A' of voice detection task) to the
343 prediction of how vulnerable voice matching is to distraction (stage two) lead to a statistically
344 significant increase in R^2 (change statistics: $F[1,95] = 7.91, p = .006$). In the full model, both
345 BVMT score and A' of the voice detection task are significant predictors of performance on
346 the distractor task, $F(2,95) = 28.77, p < .001$, adjusted $R^2 = .364$. To test whether the BVMT
347 score was a significantly better predictor than performance in the detection task, we
348 estimated the 95% confidence intervals for both standardised beta weights (calculated after
349 z-transformation of all variables) following bias corrected bootstrap (10000 iterations).
350 Confidence intervals overlapped by more than 50%, suggesting that the difference between
351 both predictors ($\Delta\beta = .223$) is not significant, and that the BVMT score was not a statistically
352 significant better predictor of resilience against distraction. Table 2 gives full details of each
353 regression stage, and Figure 1 illustrates both predictors. [Supplementary](#) Figure S1 shows
354 the relationship between both predictors.

355

356

357 Table 2
 358 Hierarchical Multiple Regression Predicting Performance on distractor task from BVMT score
 359 and voice detection task (A')
 360

Accuracy in Distractor Task				
Variable	Stage 1		Stage 2	
	B	β	B	β
Constant	0.394**		-.043	
BVMT	0.007**	0.570	.006**	.471
Detection			.556*	.248

361 Note. $N = 98$. B is unstandardised coefficients, β is standardised coefficients after z-scoring of
 362 variables. * $p < .05$, ** $p < .001$.

363
 364
 365
 366
 367

[Insert Figure 1 here]

368 Finally, paired t-tests did not reveal a difference in accuracy between trials in which
 369 T1 voice and the distractor voice were similar vs. different, neither in overall percentage
 370 correct, $t(97) = 1.31$, $p = .195$, nor in reaction times, $t(97) = .70$, $p = .484$.

371

372 Discussion

373 The experiment was designed to engage two different stages of the voice perception
 374 hierarchy through a more general voice detection task and a more voice-specific, “higher-
 375 level” voice matching task (BVMT), and investigate how both relate to the ability to tolerate
 376 interference from distractor voices (distractor task). As predicted, task performance on the
 377 BVMT correlated more highly with resilience against distraction than performance on the
 378 voice detection task. Nevertheless, both correlations were of medium to high strength (voice

379 detection: $r = .44$, BVMT: $r = .57$). A hierarchical regression analysis further explored these
380 relationships and revealed that both voice detection and voice matching task (BVMT) are
381 significant predictors of the ability to resist distraction in a voice line-up task (distractor task).
382 Including the voice detection task as an additional predictor in the model led to a significant
383 change of variance explained, and although BVMT performance was descriptively a better
384 predictor than detection task performance, further analysis revealed that the difference
385 between both predictors was not significant. In terms of variance explained, though, BVMT
386 performance alone accounted for 31.8% of the variance (stage 1), whereas the inclusion of
387 detection task performance led to 36.4% of the variance explained in the full model. We
388 suggest that this is due to both voice matching (BVMT) and voice discrimination in the
389 distractor task occurring at later processing stages along the voice perception pathway
390 whereas detecting a human voice in an array of sounds represents an earlier voice
391 perception task.

392 Face perception research has tried to explore the different processing stages in face
393 recognition and their interactions systematically (e.g. Bate & Bennetts, 2015; Calder &
394 Young, 2005). One possible approach is to thoroughly assess the range of deficits in
395 individuals with known impairments in face perception. Developmental prosopagnosia, a
396 deficit to recognise faces since childhood, has been reported in a number of case studies
397 (e.g. de Haan, 1999; Duchaine et al., 2006), and several possible explanations for these
398 deficits, including non-face specific theories, have been suggested (e.g. Farah, 1990;
399 Moscovitch, Winocur, & Behrmann, 1997). Duchaine and colleagues (2006) give a thorough
400 account of these competing alternative explanations. They also tested these alternatives
401 against each other by having an individual (Edward) with developmental prosopagnosia
402 complete a vast array of face and object perception tasks, and comparing his performance to
403 that of suitable control groups. While most face perception tasks were indeed impaired (e.g.
404 recognition of famous faces, recognition of gender or affect in faces), Edward showed
405 normal scores in a face detection task. Duchaine and colleagues therefore concluded that

406 Edward's deficits must arise at some point after the initial, low-level processing of faces as a
407 distinct category of stimuli, namely at the stage of structural encoding (as defined by Bruce &
408 Young, 1986). This would explain Edward's ability to correctly detect faces while the
409 analyses of more complex facial cues (e.g. facial affect, face identity) are disrupted. Given
410 the highly similar proposed structure of face and voice perception (Belin et al., 2004;
411 Campanella & Belin, 2007), this supports our interpretation of voice detection being one of
412 the earliest processing stages in the voice perception pathway.

413 One limitation of our findings lies in the different characteristics of each task. Of all
414 correlations, the ones with the voice detection task were the smallest, while BVMT and
415 distractor task showed the highest correlation. This could be due to the differences in
416 structure between all three tasks. Arguably, the nature of the stimuli as well as the memory
417 demands of the voice detection task (rapid presentation of human vocalisations/animate and
418 inanimate environmental sounds) differed to those of both BVMT and the distractor task
419 (judgment of two/three vocalisations per trial without time limits). The variances introduced
420 by each specific method could therefore partly drive the strength of the correlations reported
421 here. Similarly, the fact that the BVMT showed a higher correlation with the distractor task,
422 and explained more variance in the regression model than the detection task, could lie in the
423 similarity of stimuli used for both tasks (BVMT and distractor task). Both employ short non-
424 speech syllables for which speakers have to be matched. However, task demands still differ
425 considerably. Each trial in the distractor task consisted of three voices, one played shortly
426 after the other (interval between each voice: 0.8 s). Instructions then called for a
427 same/different decision regarding the first and the third voice. The BVMT, on the other
428 hand, is a task in which participants can replay the two voices per trial as often as they like
429 before making their same/different decision. As such, memory demands and time constraints
430 of both BVMT and distractor task differ considerably. In addition to that, the strength of the
431 correlation between BVMT and distractor task was only moderate to high (.57), suggesting
432 that both tasks are sufficiently different and engage overlapping but still specific abilities. In

433 order to fully address these issues in future research, an additional assessment of auditory
434 memory, as well as the inclusion of pre-ratings on all stimuli used (both in terms of physical
435 characteristics like F0, but also perceptual attributes like distinctiveness of sounds) could
436 prove helpful. Additionally, introducing a time limit on the completion of the BVMT (e.g. time
437 constraints on each trial) might help making both predictor tasks more comparable in future
438 studies, and therefore eliminate some of the variance introduced by mere task differences.

439 Distractor voices were controlled in a way that half of them showed high similarity to
440 the first target voice (T1) while the other half were markedly different. Surprisingly, we did
441 not find an effect of distractor similarity on target identification, neither in the overall
442 performance (percentage correct) nor in the reaction time data. This is in line with the
443 findings of Stevenage and colleagues (2013) who tested the resilience to distraction in both
444 face and voice perception and found that voice perception is more susceptible to distraction,
445 regardless of whether the distractor is similar or not. It is worth noting, though, that the
446 similarity manipulation in that study only matched speaker sex for target and distractor
447 voices (e.g. similar distractors being female speakers for female targets and different
448 distractors being male speakers for female targets). Stevenage and colleagues argued that
449 voice recognition was vulnerable in itself due to the relative weakness of voice perception
450 pathways. As our design used a more stringent approach to what constitutes as a similar
451 distractor (smaller distance in voice space) rather than just speaker sex, our findings support
452 the notion of voice recognition pathways being vulnerable in general.

453 Alternatively, Kreiman and Sidtis (2013) present evidence that voice identification in
454 line-up situations are always dependent on the specific listeners as well. They suggest that
455 listeners differ widely in respect to which specific voice features are attended to during voice
456 perception. It is possible that our similarity manipulations based on physical difference
457 cannot suitably account for all possible voice features that were used by the participants in
458 our particular sample. If that is the case, it could also explain our null-result for the impact of

459 distractor similarity. [For further discussion of our findings regarding distractor similarity, see](#)
460 [supplementary text ST1.](#)

461 Research into the vulnerability of voice perception and, indeed, the robustness of
462 voice identity representation over time, has mainly occurred in forensic contexts to ascertain
463 the credibility of earwitness testimony. A number of studies have tried to identify factors that
464 determine the reliability of earwitness accounts, including the duration and variability of the
465 voice sample, the number of voices that need to be identified, whether the target's face was
466 visible or not, and how much time has passed between initial exposure to a voice and
467 subsequent identification of a target from a line-up (e.g. Clifford, 1980; Cook & Wilding,
468 1997; Cook & Wilding, 2001; Legge et al., 1984; Yarmey, 1995). Our study differs from these
469 classical designs by only presenting very short voice samples without speech content and an
470 almost immediate same/different decision following voice exposure. While this design is not
471 suitable to use in forensic voice line-up situations, our findings can still contribute to our
472 insight into voice perception in general. This is relevant for our understanding of the neural
473 mechanisms underlying human voice perception on the one hand, but can ultimately also
474 lead to a better application of such findings in a more ecologically relevant setting. For
475 example, it has been proposed that a certain percentage of the population are super
476 recognisers for faces, that is, they are extremely good at using facial identity cues to
477 recognise a person (Bobak, Bennetts, Parris, Jansari, & Bate, 2016; Russell, Duchaine, &
478 Nakayama, 2009). Indeed, a special unit of UK police officers has been formed in which
479 such super-recognisers are employed to identify individuals in particularly demanding
480 identification tasks (Robertson, Noyes, Dowsett, Jenkins, & Burton, 2016). An equivalent for
481 such super-recognisers but for voices seems feasible. Having a better understanding of how
482 voice recognition at all its different stages works could therefore help in identifying such
483 voice super-recognisers.

484 The heightened interest in developmental impairments in voice perception
485 (Roswandowitz et al., 2014; Shilowich & Biederman, 2016) as well as recent research into

486 the more general question of individual differences in voice perception (Aglieri et al., 2016;
487 Mühl et al., 2017) underline the need for a better understanding of how we perceive people
488 by their voices. We propose that a more systematic approach to identifying and probing
489 possible distinct processes in the voice perception pathway will not only help our theoretical
490 understanding of voice perception, but will ultimately also impact its application in clinical
491 and, possibly, forensic settings.

492 Funding

493 Constanze Mühl is a Ph.D. student funded by the School of Psychology at Bangor

494 University.

495

- 497 Aglieri, V., Watson, R., Pernet, C., Latinus, M., Garrido, L., & Belin, P. (2016). The
498 Glasgow Voice Memory Test: Assessing the ability to memorize and
499 recognize unfamiliar voices. *Behavior Research Methods*. Advance Online
500 Publication. doi:10.3758/s13428-015-0689-6
- 501 Baumann, O., & Belin, P. (2010). Perceptual scaling of voice identity: common
502 dimensions for different vowels and speakers. *Psychological Research, 74*,
503 110-120.
- 504 Bate, S., & Bennetts, R. (2015). The independence of expression and identity in
505 face-processing: evidence from neuropsychological case studies. *Frontiers in*
506 *Psychology, 6*:770. doi:10.3758/BF03207704
- 507 Belin, P., Fecteau, S., & Bédard, C. (2004). Thinking the voice: neural correlates of
508 voice perception. *Trends in Cognitive Sciences, 8*(3), 129–135.
509 doi:10.1016/j.tics.2004.01.008
- 510 Bestelmeyer, P., Maurage, P., Rouger, J., Latinus, M., & Belin, P. (2014). Adaptation
511 to vocal expressions reveals multistep perception of auditory emotion. *Journal*
512 *of Neuroscience, 34*(24), 8098-8105.
- 513 Blank, H., Wieland, N., & von Kriegstein, K. (2014). Person recognition and the brain:
514 Merging evidence from patients and healthy individuals. *Neuroscience and*
515 *Biobehavioral Reviews, 47*, 717–734. doi:10.1016/j.neubiorev.2014.10.022
- 516 Blasi, A., Mercure, E., Lloyd-Fox, S., Thomson, A., Brammer, M., Sauter, D., Deeley,
517 Q., Barker, G. J., Renvall, V., Deoni, S., Gasston, D., Williams, S. C. R.,
518 Johnson, M. H., Simmons, A., & Murphy, D. G. M. (2011). Early specialization

519 for voice and emotion processing in the infant brain. *Current Biology*, 21(14),
520 1220-1224.

521 Bobak, A. K., Bennetts, R. J., Parris, B. A., Jansari, A., & Bate, S. (2016). An in-
522 depth cognitive examination of individuals with superior face recognition skills.
523 *Cortex*, 82, 48-62.

524 Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, 10, 443-436.

525 Bruce, V., & Young, A. (1986). Understanding face recognition. *British Journal of*
526 *Psychology*, 77, 305-327.

527 Calder, A. J., & Young, A. W. (2005). Understanding the recognition of facial identity
528 and facial expression. *Nature Reviews Neuroscience*, 6(8), 641-651.
529 doi:10.1038/nrn1724

530 Campanella, S., & Belin, P. (2007). Integrating face and voice in person perception.
531 *Trends in Cognitive Sciences*, 11(12), 535-543.

532 Cheng, Y., Lee, S-Y, Chen, H-Y., Wang, P-Y., & Decety, J. (2012). Voice and
533 emotion processing in the human neonatal brain. *Journal of Cognitive*
534 *Neuroscience*, 24(6), 1411-1419.

535 Clifford, B. R. (1980). Voice identification by human listeners: On earwitness
536 reliability. *Law and Human Behavior*, 4(4), 373-394.

537 Cook, S., & Wilding, J. (1997). Earwitness Testimony 2: Voices, faces and context.
538 *Applied Cognitive Psychology*, 11, 527-541.

539 Cook, S., & Wilding, J. (2001). Earwitness testimony: Effects of exposure and
540 attention on the Face Overshadowing Effect. *British Journal of Psychology*,
541 92(4), 617-629. doi:10.1348/000712601162374

542 De Haan, E. H. F. (1999). A familial factor in the development of face recognition
543 deficits. *Journal of Clinical and Experimental Neuropsychology*, 21(3), 312-
544 315. doi:10.1076/jcen.21.3.312.917

545 Duchaine, B. C., Yovel, G., Butterworth, E. J., & Nakayama, K. (2006).
546 Prosopagnosia as an impairment to face-specific mechanisms: Elimination of
547 the alternative hypotheses in a developmental case. *Cognitive*
548 *Neuropsychology*, 23(5), 714-747. doi:10.1080/02643290500441296

549 Farah, M. J. (1990). *Visual agnosia*. Cambridge, MA: MIT Press.

550 Gainotti, G. (2014). Familiar people recognition disorders: An introductory review.
551 *Frontiers in Bioscience*, 6, 58-64.

552 Garrido, L., Eisner, F., McGettigan, C., Stewart, L., Sauter, D., Hanley, J. R.,
553 Schweinberger, S. R., Warren, J. D., & Duchaine, B. (2009). Developmental
554 phonagnosia: A selective deficit of vocal identity recognition.
555 *Neuropsychologia*, 47, 123-131. doi:10.1016/j.neuropsychologia.2008.08.003

556 Grossman, T. (2010). The development of emotion perception in face and voice
557 during infancy. *Restorative Neurology and Neuroscience*, 28, 219-236.

558 Iredale, J., M., Rushby, J. A., McDonald, S., Dimoska-Di Marco, A., & Swift, J.
559 (2013). Emotion in voice matters: Neural correlates of emotional prosody
560 perception. *International Journal of Psychophysiology*, 89, 483-490.

561 Kleiner, M., Brainard, D., & Pelli, D. (2007). What's new in Psychtoolbox-3?
562 *Perception 36 ECVF Abstract Supplement*.

563 Kreiman, J., & Sidtis, D. (2013). *Foundations of voice studies: An interdisciplinary*
564 *approach to voice production and perception*. Malden, MA: Wiley-Blackwell.

565 Latinus, M., & Belin, P. (2011). Anti-voice adaptation suggests prototype-based
566 coding of voice identity. *Frontiers in Psychology*, 2. doi:
567 10.3389/fpsyg.2011.00175

568 Lavner, Y., Rosenhouse, J., & Gath, I. (2001). The prototype model in speaker
569 identification by human listeners. *International Journal of Speech Technology*,
570 4, 63-74.

571 Lee, I. A., & Preacher, K. J. (2013, September). Calculation for the test of the
572 difference between two dependent correlations with one variable in common
573 [Computer software]. Available from <http://quantpsy.org>

574 Legge, G. E., Grosman, C., & Pieper, C. M. (1984). Learning unfamiliar voices,
575 *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 10(2),
576 298-303.

577 Moscovitch, M., Winocur, G., & Behrmann, M. (1997). What is special about face
578 recognition? Nineteen experiments on a person with visual object agnosia and
579 dyslexia but normal face recognition. *Journal of Cognitive Neuroscience*, 9,
580 555–604.

581 Mühl, C., Sheil, O., Jarutytė, L., & Bestelmeyer, P. (2017). The Bangor Voice
582 Matching Test: A standardized test for the assessment of voice perception
583 ability. *Behavior Research Methods*. Retrieved from
584 <https://doi.org/10.3758/s13428-017-0985-4>

585 Neuner, F., & Schweinberger, S. R. (2000). Neuropsychological impairment in the
586 recognition of faces, voices, and personal names. *Brain and Cognition*. 44(3),
587 342-366. doi:10.1006/brcg.1999.1196

588 Peretz, I., Kolinsky, R., Tramo, M., Labrecque, R., Hublet, C., Demeurisse, G., &
589 Belleville, S. (1994). Functional dissociations following bilateral lesions of
590 auditory cortex. *Brain*, *117*, 1283-1301.

591 Robertson, D. J., Noyes, E., Dowsett, A. J., Jenkins, R., & Burton, A. M. (2016).
592 Face recognition by Metropolitan Police super-recognisers. *PLoS ONE*, *11*(2):
593 e0150036. doi:10.1371/journal.pone.0150036

594 Rossi-Katz, J., & Arehart, K. H. (2009). Message and talker identification in older
595 adults: Effects of task, distinctiveness of the talkers' voices, and
596 meaningfulness of the competing message. *Journal of Speech, Language,*
597 *and Hearing Research*, *52*(2), 435-453.

598 Roswadowitz, C., Mathias, S. R., Hintz, F., Kreitewolf, J., Schelinski, S., & von
599 Kriegstein, K. (2014). Two cases of selective developmental voice-recognition
600 impairments. *Current Biology*, *24*(19), 2348–2353.
601 doi:10.1016/j.cub.2014.08.048

602 Roswadowitz, C., Kappes, C., Obrig, H., & von Kriegstein, K. (2018). Obligatory
603 and facultative brain regions for voice-identity recognition. *Brain*, *141*, 234-
604 247. doi:10.1093/brain/awx313

605 Russell, R., Duchaine, B., & Nakayama, K. (2009). Super-recognizers: People with
606 extraordinary face recognition ability. *Psychonomic Bulletin & Review*, *16*(2),
607 252-257. doi:10.3758/PBR.16.2.252

608 Schirmer, A., & Kotz, S. A. (2006). Beyond the right hemisphere: brain mechanisms
609 mediating vocal emotional processing. *Trends in Cognitive Sciences*, *10*(1),
610 24-30.

611 Schneider, T. R., Engel, A. K., & Debener, S. (2008). Multisensory identification of
612 natural objects in a two-way crossmodal priming paradigm. *Experimental*
613 *Psychology*, 55(2), 121-132. doi:10.1027/1618-3169.55.2.121

614 Shilowich, B. E., & Biederman, I. (2016). An estimate of the prevalence of
615 developmental phonagnosia. *Brain & Language*, 159, 84-91.

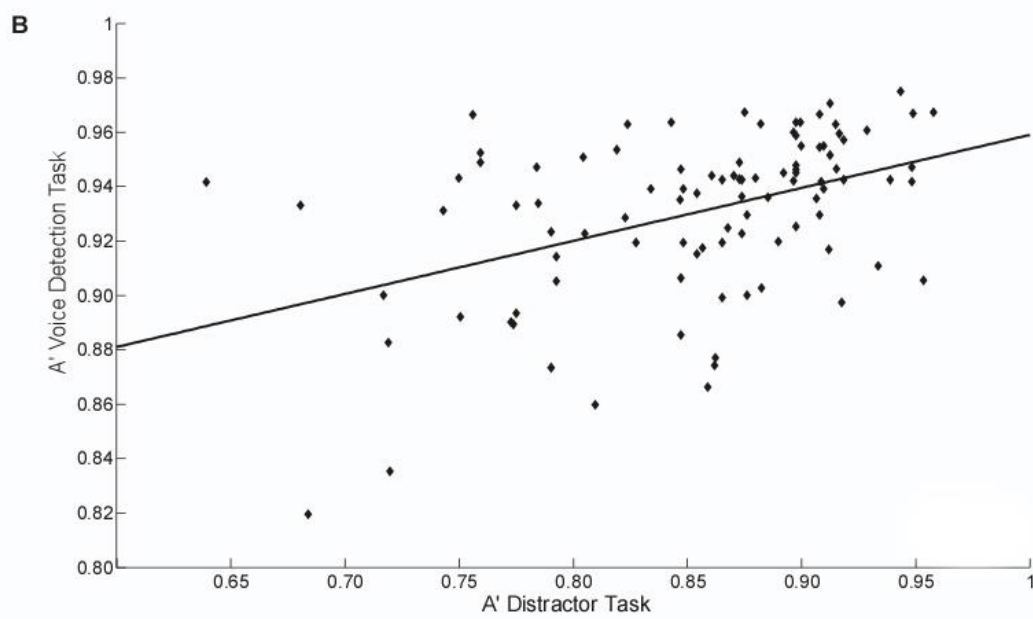
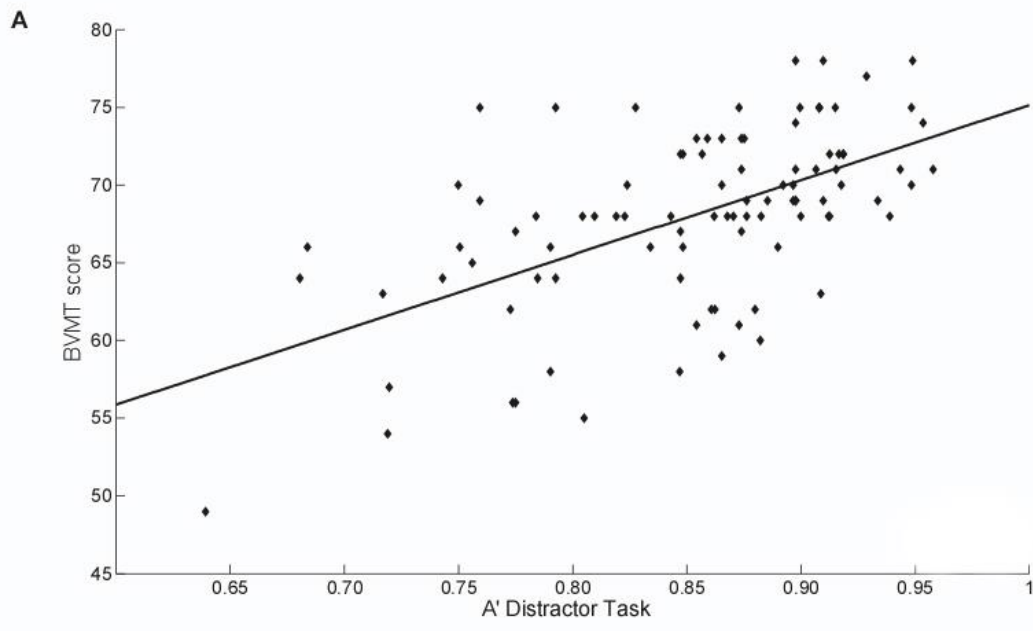
616 Stanislaw, H., & Todorov, N. (1999). Calculation of signal detection theory measures.
617 *Behaviour Research Methods, Instruments, & Computers*, 31(1), 137-149.
618 doi:10.3758/BF03207704

619 Stevenage, S. V., Neil, G. J., Barlow, J., Dyson, A., Eaton-Brown, C., & Parsons, B.
620 (2013). The effect of distraction on face and voice recognition. *Psychological*
621 *Research*, 77, 167-175. doi:10.1007/s00426-012-0450-z

622 van Lancker, D., & Kreiman, J. (1987). Voice discrimination and recognition are
623 separate abilities. *Neuropsychologia*, 25(5), 829-834.

624 Warren, J. D., Jennings, A. R., Griffiths, T. D. (2005). Analysis of the spectral
625 envelope of sounds by the human brain. *Neuroimage*, 24(4), 1052-1057.

626 Yarmey, A. D. (1995). Earwitness speaker identification. *Psychology, Public Policy,*
627 *and Law*, 1(4), 792-816.



628

629

630 *Figure 1. Relationship between (A) performance on BVMT (score) and distractor*
 631 *task (A') and (B)*

632 *performance on voice detection task (A') and distractor task (A'). Lines represent*

633 *linear regression fits to data points.*

634 Acknowledgements

635 We would like to thank Dr. Till Schneider and his team for providing us with the Multimost
636 materials which made up part of the stimulus set used in this study. The Multimodal Stimulus
637 Set was developed by T. R. Schneider, S. Debener and A. K. Engel at the Dept. of
638 Neurophysiology, University Medical Center Hamburg-Eppendorf, Germany.

639

640 Supplemental Online Material

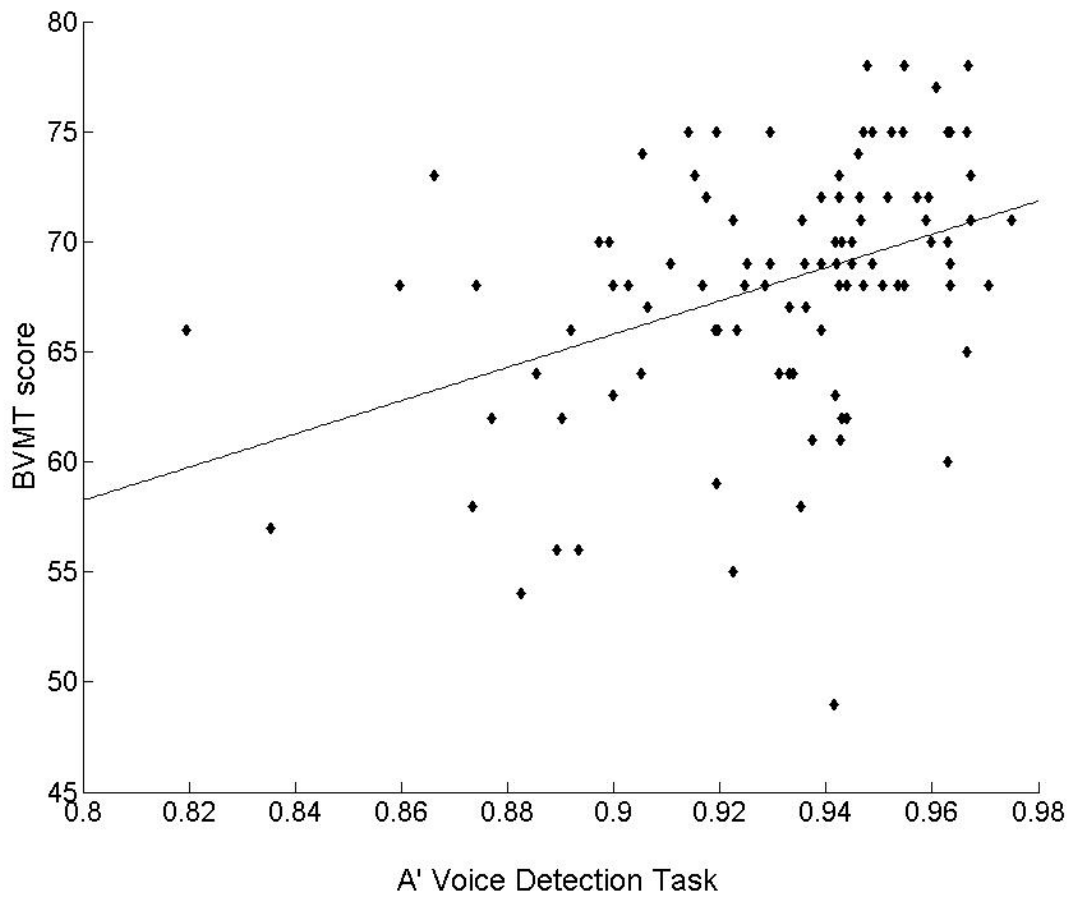
641 *SE1*. Equation for A' calculation for SPSS from Stanislaw & Todorov (1999). H denotes hit rate, F
642 denotes false alarm rate:

643

644
$$A' = 0.5 * \left(\frac{abs(H - F)}{H - F} \right) * \frac{(H - F)^2 + abs(H - F)}{4 * MAX(H, F) - 4 * H * F}$$

645

646

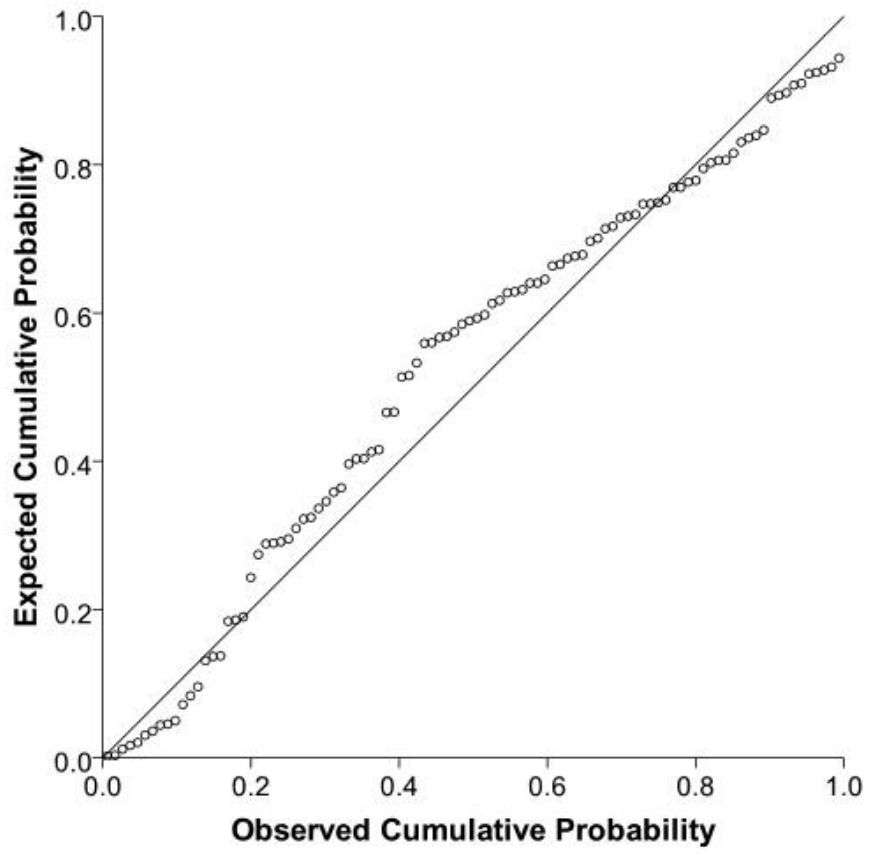


647

648 *Figure S1.* Relationship between performance on both independent variables (BMVT score and A' in
649 distractor task; $r = .40$). Line represents linear regression fit to data points.

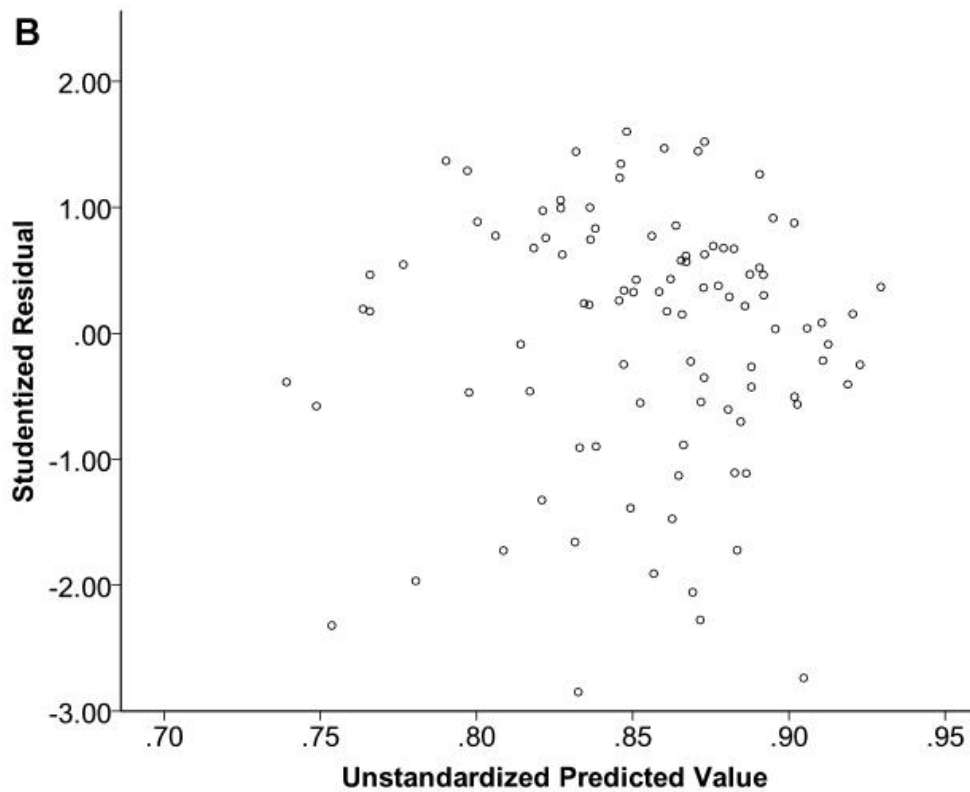
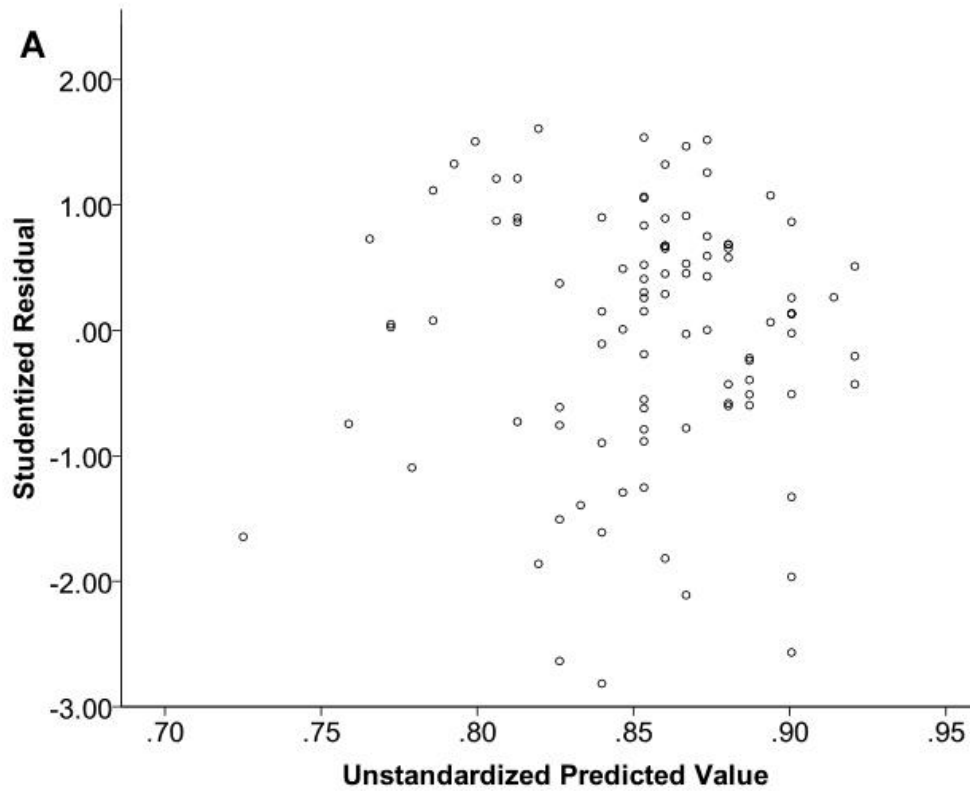
650

651



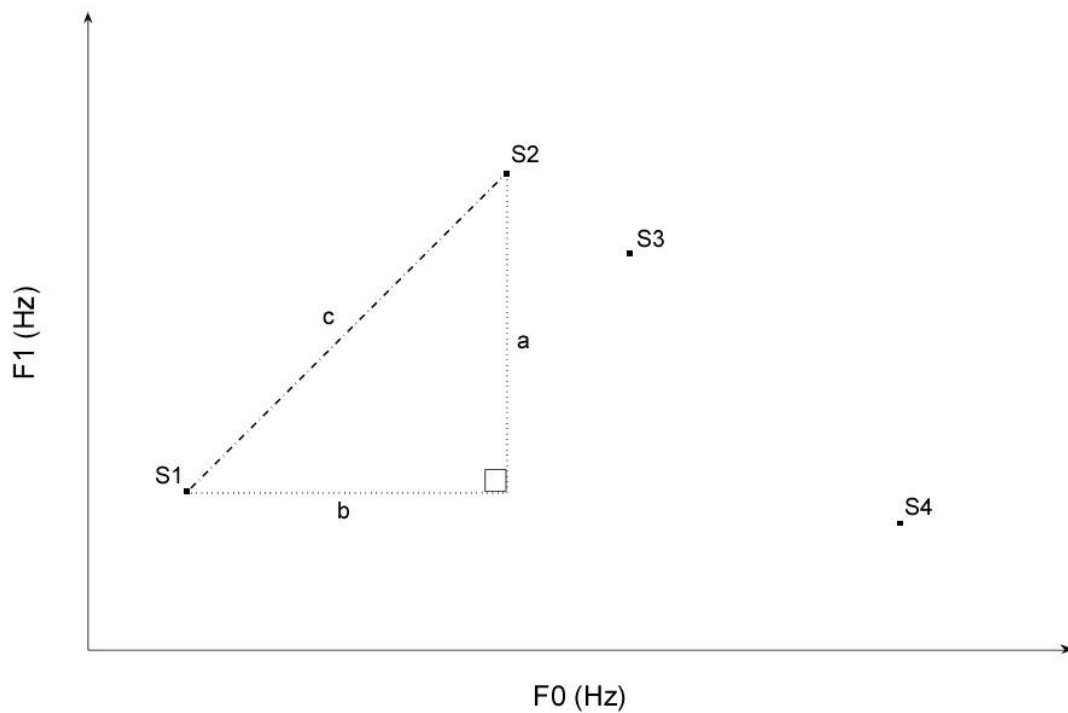
652
653 *Figure S2. P-P-Plot for regressions standardised residual.*

654



655

656 *Figure S3.* Distribution of standardised residuals against unstandardized predicted values for
 657 stage 1 of the model (**A**; BVMT score as sole predictor) and the full model (**B**; BVMT score
 658 and A' of detection task).



659

660 *Figure S4.* Schematic representation of voice space. Individual speakers (S1 to S4) are illustrated
 661 within a 2-dimensional voice space (Baumann & Belin, 2010), according to their fundamental
 662 frequency (F0) and their first formant frequency (F1). Voices that are close to each other (e.g. S2 and
 663 S3) sound more similar than those further apart (e.g. S2 and S4). Physical difference between S1 and
 664 S2 (alternating dashed line, hypotenuse c) is calculated using the Pythagoras theorem, given a right
 665 triangle with legs a and b (simple dashed lines), $c = \sqrt{a^2 + b^2}$.

666

667

668 ST1. Accuracy in distractor task based on similarity between distractor and second target
669 voice (T2)

670 Additional post-hoc analyses of accuracy for trials with similar vs. different distractor
671 and T2 voice pairings showed a significant difference in mean percentage correct, $t(97) = -$
672 2.53 , $p = .013$, with a higher accuracy for trials in which physical D-T2 distance was greater
673 ($M = 78.72\%$, $SD = 8.73$) compared to smaller D-T2 distances ($M = 76.26\%$, $SD = 10.90$).
674 However, this difference did not reach significance in the reaction time data ($t[97] = 1.79$, $p =$
675 $.077$). Our post-hoc analysis therefore revealed a significantly higher accuracy if the
676 distractor voice was markedly different to the T2 voice.

677 While this is in line with our initial prediction for the impact of distractor similarity, we
678 are cautious to interpret this finding. Unlike for the T1-D pairings, the number of
679 similar/different D-T2 pairings was not equal due to the limited availability of suitable voice
680 pairings. Consequently, as stated before, our predictions only considered the effect a
681 distractor voice could have for the accuracy of identifying a previously heard target voice
682 (T1). This issue needs to be revisited in future studies where the distractor similarity for both
683 target voices, T1 and T2, can be controlled more stringently (given a larger pool of initial
684 voice pairings).

685 Further indication of an effect of distractor similarity comes from research into
686 changes of our ability to identify speakers from different age ranges. Rossi-Katz and Arehart
687 (2009) manipulated distinctiveness of distractor voices via speaker sex, and investigated its
688 effect on the accuracies of (a) identifying a target message, that is, speech content, and (b)
689 identifying a target speaker identity. Both manipulations were tested in a group of young
690 adults (23 – 25 years of age) as well as in a group of older adults (> 65 years of age). While
691 the target message task profited from increased speaker distinctiveness (albeit to a lower
692 extent in the older group), target identification did not. Young adults showed high speaker
693 identification accuracy regardless of distractor distinctiveness whereas older adults showed
694 a decline of speaker identification accuracy for more distinct distractors (meaningful speech
695 condition). The null effect of distractor similarity/differences in Stevenage and colleagues'

696 study (2013) as well as in ours might therefore be due to the nature of the samples used
697 (young adults), and further investigation into different samples seems necessary.