



**HAL**  
open science

## Do Voice-Based Judgments of Socially Relevant Speaker Traits Differ Across Speech Types?

Agata Groyecka-Bernard, Katarzyna Pisanski, Tomasz Frąckowiak, Aleksander Kobylarek, Piotr Kupczyk, Anna Oleszkiewicz, Agnieszka Sabiniewicz, Monika Wróbel, Piotr Sorokowski

► **To cite this version:**

Agata Groyecka-Bernard, Katarzyna Pisanski, Tomasz Frąckowiak, Aleksander Kobylarek, Piotr Kupczyk, et al.. Do Voice-Based Judgments of Socially Relevant Speaker Traits Differ Across Speech Types?. *Journal of Speech, Language, and Hearing Research*, 2022, 65 (10), pp.3674-3694. 10.1044/2022\_JSLHR-21-00690 . hal-03799551

**HAL Id: hal-03799551**

**<https://hal.science/hal-03799551v1>**

Submitted on 5 Oct 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33

**Do voice-based judgments of socially relevant speaker traits differ across speech types?**

Agata Groyecka-Bernard<sup>1,2</sup>, Katarzyna Pisanski<sup>1,3,4</sup>, Tomasz Frąckowiak<sup>1</sup>, Aleksander Kobyłarek<sup>5</sup>, Piotr Kupczyk<sup>1</sup>, Anna Oleszkiewicz<sup>1,5</sup>, Agnieszka Sabiniewicz<sup>1,5</sup>, Monika Wróbel<sup>1</sup>, Piotr Sorokowski<sup>1\*</sup>

<sup>1</sup> University of Wrocław, Institute of Psychology, Wrocław, Poland

<sup>2</sup> Johannes Gutenberg University Mainz, Mainz, Germany

<sup>3</sup> ENES Bioacoustics Research Laboratory, University of Saint-Etienne, CRNL, CNRS UMR 5292, Inserm UMR\_S 1028, Saint-Etienne, France

<sup>4</sup> CNRS Centre National de la Recherche Scientifique, Laboratoire Dynamique du Langage, Université Lyon 2, Lyon, France

<sup>5</sup> University of Wrocław, Institute of Pedagogy, Wrocław Poland

\*Correspondence: [piotr.sorokowski@uwr.edu.pl](mailto:piotr.sorokowski@uwr.edu.pl), University of Wrocław, Institute of Psychology, ul. Dawida 1, 50-529, Wrocław, Poland

Words count: 8503

**Conflict of interest statement:** The authors declare no conflict of interest.

**Funding Statement:** The study was funded by the Polish National Science Center grant OPUS (2016/23/B/HS6/00771) awarded to PS.

34 ABSTRACT

35 **Purpose:** The human voice is a powerful and evolved social tool, with hundreds of studies  
36 showing that nonverbal vocal parameters robustly influence listeners' perceptions of socially  
37 meaningful speaker traits, ranging from perceived gender and age to attractiveness and  
38 trustworthiness. Yet these studies have utilized a wide variety of voice stimuli to measure  
39 listeners' voice-based judgments of these traits. Here, in the largest scale study known to date,  
40 we test whether listeners judge the same unseen speakers differently depending on the  
41 complexity of the neutral speech stimulus, from single vowel sounds to a full paragraph.

42 **Method:** In a playback experiment testing 2618 listeners, we examine whether commonly  
43 studied voice-based judgments of attractiveness, trustworthiness, dominance, likability,  
44 femininity/masculinity and health differ if listeners hear isolated vowels, a series of vowels,  
45 single words, single sentences (greeting), counting from 1 to 10, or a full paragraph recited  
46 aloud (Rainbow Passage), recorded from the same 208 men and women. Data were collected  
47 using a custom designed interface in which vocalizers and traits were randomly assigned to  
48 raters.

49 **Results:** Linear mixed models show that the type of voice stimulus does indeed consistently  
50 affect listeners' judgments. Overall, ratings of attractiveness, trustworthiness, dominance,  
51 likability, health, masculinity among men and femininity among women increase as speech  
52 duration increases. At the same time, speaker-level regression analyses show that inter-  
53 individual differences in perceived speaker traits are largely preserved across voice stimuli,  
54 especially among those of a similar duration.

55 **Conclusions:** Socially relevant perceptions of speakers are not wholly changed but rather  
56 moderated by the length of their speech. Indeed, the same vocalizer is perceived in a similar  
57 way regardless of which neutral statements they speak, with the caveat that longer utterances  
58 explain the most shared variance in listeners' judgments and elicit the highest ratings on all  
59 traits, possibly by providing additional nonverbal information to listeners.

60

61 *Key-words: voice, stimulus type, stimulus duration, voice perception, playback experiment*

62 Introduction

63 The human voice is a source of abundant information about the vocalizer, which listeners can  
64 use to make socially relevant decisions. Hundreds if not thousands of experimental studies  
65 have shown that the human voice plays a central role in predicting listeners' perceptions of  
66 the social and biological qualities of vocalizers, including psychological and physical traits  
67 (reviewed in: Aung & Puts, 2020; Kamiloğlu & Sauter, 2021; Kreiman & Sidtis, 2011;  
68 Pisanski & Bryant, 2019). For example, the nonverbal properties of a person's voice,  
69 particularly fundamental frequency ( $f_0$ ) perceived as voice pitch, and formant frequencies  
70 affecting voice timbre, can predict listeners' judgments of a vocalizer's personality traits  
71 (Stern et al., 2021), social dominance (Aung & Puts, 2020; Borkowska & Pawlowski, 2011;  
72 David A. Puts et al., 2016), attractiveness (reviewed in Pisanski & Feinberg, 2019), and health  
73 (Vukovic et al., 2010), to name only a few. In turn, these same voice features can predict who  
74 people vote for in an election (Klofstad et al., 2012; Mileva et al., 2020; Tigue et al., 2012),  
75 choose to hire in a job interview (Anderson et al., 2014), or choose as a romantic partner  
76 (Pisanski et al., 2018; Rosenfield et al., 2020).

77 Considering the increasing prevalence of voice research in the human behavioral sciences,  
78 and mounting empirical evidence that nonverbal voice parameters predict perceptions of the  
79 vocalizer in the 'ears' of the beholder, remarkably few studies have examined whether such  
80 voice-based perceptions depend on the type of voice stimulus being judged. Indeed, the  
81 abovementioned studies on human nonverbal voice production and perception utilized a  
82 variety of speech stimuli to test listeners' perceptions of vocalizer traits. Traditionally, these  
83 have often included affectively neutral speech, such as vowel sounds or scripted sentences  
84 and paragraphs, designed to control for linguistic content, but nevertheless varying in duration  
85 and complexity. The potential effects driven by these differences in duration or complexity  
86 across speech stimuli, and thus the amount of nonverbal information available to the listener,

87 remain largely unknown. There is hence a compelling need to gain knowledge about the  
88 validity and comparability of the diverse methods used to study vocal communication in  
89 humans. Testing the extent to which listeners' socially relevant judgments of speakers remain  
90 stable across different utterances produced by the same person also carries theoretical and  
91 social implications, namely regarding the underlying mechanisms driving voice-based  
92 judgments that are in turn known to predict important real-world outcomes for speakers. Here,  
93 to this aim, we directly test whether differences in the type of neutral speech stimulus used  
94 can affect listeners' judgements of the same vocalizers on a range of socially relevant traits.

95       Studies testing human voice perception often follow a similar basic protocol and are  
96 typically referred to as playback, psychoacoustic, or perception experiments. At the most  
97 basic level, the first step is to record vocalizers' voices, and the second is to subsequently play  
98 these speech stimuli to a sample of listeners who then evaluate them on various scales (e.g.,  
99 pertaining to personality, attractiveness, physical traits). Some studies utilize unchanged,  
100 natural voice samples (Cartei et al., 2014; McAleer et al., 2014; Pisanski et al., 2014;  
101 Sorokowski et al., 2019) while others manipulate the nonverbal acoustic properties of speech  
102 to causally test how changes in specific acoustic parameters (for instance voice pitch) affect  
103 listeners' judgements (e.g., Albert et al., 2021; Belin et al., 2017; Feinberg et al., 2005; Krahé  
104 et al., 2021; Pisanski et al., 2012).

105       The type of voice stimuli used in playback experiments vary greatly. Some researchers  
106 use standard phrases like vowels (for example a e i o u, Feinberg et al., 2008), counting from  
107 1 to 10 (e.g., Hughes et al., 2014), single words (e.g., greetings, see Apicella & Feinberg,  
108 2009), single sentences (e.g., Jones et al., 2008) or a phonetically balanced paragraph such as  
109 the Rainbow Passage that is recited aloud (Fairbanks, 1960) – a passage that contains a broad  
110 representation of vowel sounds (e.g., Pisanski et al., 2016; Pisanski, Anikin, et al., 2021; Puts  
111 et al., 2006; Schild et al., 2020). These stimulus types differ in numerous ways, most notably

112 in their duration and complexity. Preliminary evidence points to a negative relationship  
113 between the duration of a voice stimulus and within-listener stability in voice perception  
114 (Ohno et al., 2014). One possible explanation for this is that longer stimuli provide more  
115 nonverbal information than do shorter stimuli, particularly regarding stable or dynamic  
116 acoustic parameters of the speaker's voice such as fundamental and formant frequencies.  
117 Longer stimuli also offer more time for listeners to consider their responses. On the other  
118 hand, by increasing the amount of nonverbal information available to listeners, longer speech  
119 may introduce more variability in listener's judgments of the vocalizer owing to more  
120 opportunities for listeners to express their individual differences in voice preferences and  
121 perceptions, and in turn, may reduce inter-rater agreement.

122         A recent study has shown that individual differences in voice fundamental frequency  
123 ( $f_0$ , perceived as voice pitch) are preserved across speech types (Pisanski, Groyecka-bernard,  
124 et al., 2021). Pisanski and colleagues (2021) analyzed  $f_0$  in six different types of neutral  
125 speech utterances (from vowels to longer bouts of spontaneous speech) and showed that inter-  
126 individual differences in this salient voice property are highly robust, such that a person's  
127 voice pitch when speaking a vowel sound correlates strongly with their voice pitch when  
128 speaking a full paragraph of free speech. These results thus demonstrate the methodological  
129 validity of comparing  $f_0$  measures across different neutral speech stimulus types.  
130 Nevertheless, these results, based only on acoustic measures, do not necessitate that listeners'  
131 perceptions will likewise be robust across speech types, wherein stimuli can differ on a range  
132 of other spectrotemporal parameters that may affect listeners' voice-based judgments. Some  
133 voice stimuli also have higher ecological validity than do others (as in the case of greetings  
134 versus vowels), which may differentially influence listeners' judgements. For example, single  
135 vowel sounds, in addition to providing limited information regarding articulation compared to  
136 longer phrases (Kreiman, 1997), are rarely used in isolation in everyday real-life

137 conversations. Some researchers have thus opted to use standardized greetings resembling  
138 statements one might use in everyday life (e.g., “Hi, I’m a student at UCLA”: Bryant &  
139 Haselton, 2009; “Get out and be quiet”: Sell et al., 2010), or entirely unscripted speech  
140 produced in response to a given context (e.g., discussing one’s admirable traits in a  
141 competitive context: Puts et al., 2006), to increase the ecological validity of speech stimuli  
142 collected in the lab.

143         Very few studies have attempted to address how the use of different voice stimuli  
144 might affect the results of playback experiments. One study tested for differences in  
145 attractiveness ratings based on voice duration and type, but the researchers considered only  
146 very short (i.e., single vowel, three vowels, and ‘bonjour’) stimuli and artificial manipulation  
147 thereof (Ferdenzi et al., 2013). The authors showed that artificially shortening the voice  
148 samples decreased the perceived attractiveness of the same vocalizers and that words elicited  
149 higher attractiveness ratings compared to vowels. While Ferdenzi et al. focused solely on  
150 attractiveness ratings, Mahrholz and colleagues (2018) tested whether listeners’ judgements  
151 of attractiveness, dominance, and trustworthiness vary for one word versus one sentence. In  
152 the one-sentence stimulus, the authors also manipulated content (i.e., social relevance). In that  
153 study, listeners’ judgments were highly correlated for words and sentences produced by the  
154 same set of vocalizers, which could potentially mean that speech content did not influence  
155 judgments in that study. The findings thus suggest that, regardless of speech stimulus type,  
156 vocalizers are perceived similarly. However, these studies were limited to perceptions of only  
157 three vocalizer traits: attractiveness, dominance, and/or trustworthiness. They did not explore  
158 how the type of speech stimulus affects perceptions of other traits known to be highly relevant  
159 in human interpersonal relationships, and known to be perceptually linked to nonverbal  
160 parameters of the voice, including masculinity/femininity, likeability, and health (for review  
161 see Kreiman & Sidtis, 2011). Moreover, Mahrholz et al. (2018) used only two types of stimuli

162 (word, sentence), which, similar to Ferdenzi et al. (2013), cover only a small share of the  
163 methodological diversity observed in the literature, with both studies focusing on short  
164 stimuli.

165

166         The current study was designed to complement the scarce literature regarding the  
167 potential effects of the type of speech stimulus on listeners' perceptions of unseen vocalizers.  
168 In a large-scale playback experiment, we test whether a broad variety of utterances (consisting  
169 of isolated vowels, vowels pronounced in a series, single words, counting 1-10, greeting  
170 sentences, and reading aloud a phonetically balanced passage) elicit different assessments of  
171 the same vocalizers' perceived attractiveness, trustworthiness, dominance, likability,  
172 femininity/masculinity and health. The study examines the effects of specific stimulus types  
173 as well as the general effect of stimulus duration on perceptions of these traits as judged by a  
174 large sample of over 2000 male and female listeners. Because voice perception can heavily  
175 rely on inter-vocalizer or inter-listener factors, we controlled for variance due to individual  
176 differences in multilevel models.

177         The perceived vocalizer traits tested here were selected on the basis of countless  
178 studies that have shown that listeners' perceptions of these specific traits, while not  
179 exhaustive, are robustly influenced by nonverbal parameters of the voice, and are thus among  
180 the most intensively studied voice-based traits in the human voice sciences (for reviews see  
181 Kamiloğlu & Sauter, 2021; Kreiman & Sidtis, 2011; Pisanski & Bryant, 2019). Importantly,  
182 voice-based perceptions of these same traits are also known to predict a range of social  
183 decisions, such as mate preferences (Rosenfield et al., 2020) and election outcomes (Mileva et  
184 al., 2020).

185         The question of whether the type of speech uttered affects how a person is perceived is  
186 important from several perspectives. From a methodological perspective it can verify whether



187 results of past studies using myriad types of voice stimuli to test listeners' perceptions of the  
188 same trait are comparable. From a social perspective, it may provide novel insight into  
189 whether such perceptions are robustly preserved across different speech types. For example,  
190 would a speaker uttering a series of sentences be perceived as similarly attractive,  
191 trustworthy, or likeable if uttering only a single word? Likewise, are voice-based perceptions  
192 of masculinity or femininity, dominance, and health dependent on the complexity of a speech  
193 utterance? Given the importance of these perceived traits for interpersonal relationships, mate  
194 choice, and/or broader societal outcomes, the extent to which such judgments are made based  
195 on snippets of nonverbal vocal cues is of high public relevance. From a practical perspective,  
196 the research question is also relevant for voice-based technologies such as automated voice  
197 recognition and detection devices or voice-based clinical diagnostic tools, in which algorithms  
198 relying on artificial intelligence may be refined depending on the robustness of data obtained  
199 from various durations and types of speech.

200

## 201 **Method**

202 The study was conducted in accordance with the Declaration of Helsinki. Study protocols  
203 were accepted by the institutional ethics committee at the University of Wrocław. All  
204 participants (vocalizers and raters) provided informed consent prior to participation.

205 Vocalizers were informed that their voice recordings will be further used for the purposes of  
206 the present study and that they will be played to other participants.

207

### 208 **Participants**

209 The number of raters was based on the number of vocalizers with the aim that each voice  
210 stimulus be rated approximately 20 times, based on evidence that high inter-rated agreement  
211 ( $\alpha > 0.80$ ,  $p < 0.001$ ) among listeners is typically achieved with relatively small sample

212 sizes (e.g., less than 15 listeners per sex for voice-based judgements of dominance or  
213 attractiveness: Kordsmeyer et al., 2018; Schild et al., 2020). The number of speakers (ca. 200)  
214 was determined as a trade-off between a sample that translates into a reasonable number of  
215 raters and one that allows us to detect inter-individual differences in traits of interest and  
216 nonverbal aspects of voice communication.

#### 217 *Vocalizers*

218 208 vocalizers representing a broad age range and balanced sex ratio ( $M_{\text{age}} = 32.83$ ,  $SD_{\text{age}} =$   
219  $12.32$ , 48% women, 52% men) provided voice recordings for use as stimuli in playback  
220 experiments. Their native language was Polish and they were recruited through snowball  
221 sampling by researchers and research assistants who posted recruitment ads on their social  
222 media profiles, and around their city of residence, both inside and outside of the university.  
223 Vocalizers were not compensated for providing speech samples.

#### 224 *Voice raters*

225 2618<sup>1</sup> voice raters representing a broad age range and balance sex ratio ( $M_{\text{age}} = 32.51$ ,  $SD_{\text{age}} =$   
226  $13.01$ , 54% women, 46% men) judged the speech samples. All raters reported normal hearing.  
227 To increase our sample size and the diversity of our rater sample, raters were recruited via a  
228 combination of the snowball sampling method and through a dedicated research recruitment  
229 firm. Snowball sampling followed the same procedures as for vocalizers, with additional  
230 recruitment efforts targeting older and elderly individuals. All raters (lab and online) were  
231 residents of Poland and understood written Polish as verified by the recruitment firm and/or in  
232 the attention and hearing tests preceding the experiment. Participants were reimbursed in cash  
233 (for the cohort recruited via the recruitment firm) or through a lottery draw of small prizes  
234 such as pen drives (for all remaining participants).

235

---

<sup>1</sup> Data from more participants were collected, however, they were asked different questions for the purposes of a different study.

236 Materials

237 *Voice recording*

238 Participants were recorded at the Institute of Psychology at the University of [covered for  
239 blind review]. First, participants were audio recorded in private sessions in a small quiet room  
240 with a low level of external noises. We used a Zoom H4n professional digital recorder with  
241 X/Y stereo microphone array positioned 10 cm from the mouth. Participants familiarized  
242 themselves with a script and were then instructed to say aloud five items in their native  
243 language (Polish), listed below in English translation:

244 I. Vowels a-e-i-o-u (/a/, /ɛ/, /i/, /ɔ/, /u/, International Phonetic Alphabet)<sup>2</sup>

245 II. One-syllable word “lat” (containing the vowel /a/)

246 III. Counting from 1 to 10

247 IV. Greeting sentence “Dzień dobry, jestem z Polski” [Good morning, I am from Poland]

248 V. First paragraph of the Rainbow Passage (Fairbanks, 1960)

249 After recording the voice samples, vocalizers filled in a short demographic survey.

250 Participants were then thanked and debriefed. Voice recordings were saved as WAV files at  
251 96 kHz sampling frequency and 16-bit resolution, and then further divided into short  
252 fragments each containing one specific stimulus type. Vowels were saved together within a  
253 single segment and also saved separately<sup>2</sup>. Therefore, 6 different stimulus types from each  
254 vocalizer (1248 speech stimuli in total) were prepared for the playback experiments. We  
255 coded not only the type of speech but also its length, i.e., number of syllables (1, 5, 5, 16, 8,  
256 150 for isolated vowels, vowel series, word, counting, greeting and paragraph read aloud,  
257 respectively).

258 *Stimulus preparation*

---

<sup>2</sup> Vowels were recorded as a sequence. For playback experiments, they were then presented either in isolation, or as a sequence.

259 Voice samples were incorporated into a custom designed online web app. The survey  
260 included display of voice samples, assessment (rating) scales and demographic questions.

261

## 262 Procedure

263 The voice rating phase took place in person at the University lab, or online. Participants who  
264 took part at the lab conducted the playback experiment either separately, or in small groups (2  
265 to 5 individuals) in rooms prepared accordingly. To ensure privacy and independence in  
266 responses, participants who conducted the study in a small group were placed in partitioned  
267 computer booths to reduce visual contact, and always wore headphones. All lab participants  
268 used high quality professional headphones (Sennheiser HD 210). Participants who completed  
269 the playback online were instructed to use good quality headphones and to complete the study  
270 in a quiet environment without distractions. This was verified with hearing and attention tests.  
271 Before beginning the experiment, participants were exposed to a test voice sample to ensure  
272 they can hear the stimuli properly, and to set and standardize playback volume. Eighty  
273 participants (3%) failed this hearing test and their data were thus omitted from further  
274 analyses. At a random time during the playback, participants were also presented with an  
275 additional attention check item that read, “This is an attention checking question – please,  
276 mark 1”. Eleven additional participants (<1%) were excluded due to an incorrect answer.  
277 For each participant, a sample of 10 recorded individuals (5 men and 5 women) was randomly  
278 drawn. For each of these vocalizers, a set of 6 voice samples was presented (resulting in 60  
279 voice stimuli per listener). These stimuli were presented in a random order, each time  
280 followed by a request to evaluate the vocalizer on one given dimension. Thus, each speech  
281 stimulus was judged independently on a single trial and for a single trait. Each listener  
282 evaluated 60 voice samples in total on the same randomly assigned trait, i.e., one of the  
283 following questions appeared after each speech recording:

- 284 1. How attractive is this person? (1 = very unattractive, 7 = very attractive)
- 285 2. How dominant is this person? (1 = not dominant at all, 7 = very dominant)
- 286 3. How likeable is this person? (1 = not likable at all, 7 = very likable)
- 287 4. How trustworthy is this person? (1 = not trustworthy at all, 7 = very trustworthy)
- 288 5. How feminine/masculine is this person? (1 = very feminine, 7 = very masculine)
- 289 6. How healthy is this person? (1 = very unhealthy, 7 = very healthy)

290 To align closely with most studies in this domain (e.g., Cartei et al., 2014; Feinberg et al.,  
291 2012; Hughes et al., 2014; McAleer et al., 2014; Pisanski & Rendall, 2011), participants were  
292 not provided with any definitions of these concepts. Following the experiment, they reported  
293 their age and sex and were debriefed. Rating sessions lasted an average of 30 minutes  
294 including instructions, hearing and attention checks, playback (rating speech stimuli), brief  
295 survey, and debriefing.

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310 Statistical Analysis

311 In order to test the effects of stimulus type on perceived attractiveness, dominance, likability,  
312 trustworthiness, masculinity/femininity, and health, we proceeded with a series of Linear  
313 Mixed Models (LMMs). Due to sizeable sexual dimorphism in the nonverbal properties of the  
314 human voice (Titze, 1989) and known differences in perceived attractiveness, masculinity and  
315 dominance for male versus female voices (Borkowska & Pawlowski, 2011; Cartei et al.,  
316 2014; Pisanski & Rendall, 2011), we conducted separate LMMs for male and female  
317 vocalizers. The data fulfilled the assumption of residual normality which was identified by  
318 visual exploration of Q-Q plots of residuals. The only exception where the distribution was  
319 slightly different from normal was femininity-masculinity, however, this violation was not  
320 extreme.

321         The dependent variables (ratings based on a single voice stimulus) were nested within  
322 vocalizers and within listeners. The models differed only by the outcome variable (each time  
323 a different perceived trait as an outcome) and were estimated using a Restricted Maximum  
324 Likelihood (REML) estimator. We included type of voice stimulus as a factor, vocalizer's age  
325 and listener's sex (0-F, 1-M) and method of data collection during playback experiments (0-  
326 online raters, 1-lab raters) as covariates. These variables were treated as fixed effects.  
327 Random effects of vocalizer numeric ID and listener numeric ID were included to control for  
328 noise in the models owing to potential individual differences in speech stimuli across  
329 vocalizers or systematic biases in different listeners' ratings, i.e., random effects refer to  
330 speaker- and vocalizer-level variance. All continuous predictors were grand-mean centered.  
331 Each model was followed up by Bonferroni corrected pairwise comparisons. As a measure of  
332 effect size, we compared the amount of explained variance in trait ratings in models that  
333 included fixed effects of speech stimulus type against a null model (without predictors,  
334 accounted for clustering).

## 335 **Results**

336 Our results show that listeners' ratings on almost every trait generally increased with the  
337 duration of the voice stimulus. In general, single vowels and single words elicited the lowest  
338 ratings, longer utterances (counting, greetings and paragraph read aloud) elicited the highest  
339 ratings, and vowels pronounced in a series often obtained intermediate ratings (see dominance  
340 for an exception). The absolute differences between the lowest and highest rated stimulus  
341 types were significant, for all six dimensions, and ranged between 0.30 for masculinity among  
342 male vocalizers and 1.08 for trustworthiness in female vocalizers (on a 1-7 scale). The vast  
343 majority of differences in ratings across stimuli were statistically significant (for exceptions  
344 see Figure 1a-f).

### 345 *Attractiveness*

346 Overall average perceived voice attractiveness of female vocalizers, controlling for other  
347 predictors, was  $3.88 \pm 0.06$  (mean  $\pm$  standard error of the mean). The effect of voice speech  
348 stimulus type was significant  $F(5, 12515) = 219.14, p < .001$ . For the estimates of fixed and  
349 random effects see Table 1. The highest attractiveness ratings overall were those assessed  
350 after listening to the recited paragraph (longest utterance,  $4.24 \pm 0.07$ ), followed by counting  
351 ( $4.17 \pm 0.07$ ) and greetings ( $4.10 \pm 0.07$ ) (intermediate duration stimuli), and then by the  
352 series of vowels ( $4.04 \pm 0.07$ ). These same vocalizers were judged as least attractive when  
353 speaking only a single word ( $3.48 \pm 0.07$ ) or single vowel ( $3.34 \pm 0.07$ ), the shortest speech  
354 utterances. Estimated marginal means controlling for covariates are presented in Table S1 in  
355 Supplementary Materials. Pairwise comparisons with Bonferroni correction are presented in  
356 Figure 1a and Table S2. Perceptions of attractiveness were positively related to the age of the  
357 vocalizer ( $b = -0.03, p < .001$ ) but not to the rater's sex ( $b = 0.01, p = 0.99$ ). There were no  
358 significant differences in ratings between raters who completed the playback experiment in  
359 the lab versus online ( $b = -0.17, p = 0.06$ ).

360 Overall perceived attractiveness of male vocalizers, controlling for other predictors, was 3.53  
361  $\pm 0.06$ ). The effect of stimulus type was significant  $F(5, 12260) = 168.30, p < .001$ . The  
362 pattern of results was almost identical to that observed for female vocalizers with a slight  
363 difference in the most highly evaluated stimuli: greetings were evaluated as similarly  
364 attractive ( $3.84 \pm 0.07$ ) as counting ( $3.79 \pm 0.07$ ), and there was no significant difference  
365 between greetings and recited paragraphs ( $3.91 \pm 0.07$ ). Thus, for both sexes of vocalizers,  
366 attractiveness ratings tended to increase with the duration of the speech stimulus (See Figure  
367 1a, and Tables S3 and S4 for post-hoc comparisons and estimated marginal means). Younger  
368 vocalizers were evaluated as more attractive than were older vocalizers ( $b = -0.01, p = .018$ )  
369 and female listeners rated men's voices as more attractive than did male listeners ( $b = -0.48, p$   
370  $< .001$ ).

371 Comparing groups of raters who completed the playback experiment in the lab versus  
372 online, those who completed the study in the lab rated men's voices as less attractive overall  
373 ( $b = -0.22, p = 0.02$ ) compared to those who completed it online ( $3.68 \pm 0.02$  vs.  $3.47 \pm 0.02$ ,  
374 respectively). Adding speech stimulus type as a predictor did not substantially increase the  
375 amount of variance in attractiveness ratings explained by the model (0.08% and 0.06% for  
376 female and male voices, respectively). Due to a significant effect of rater group, we conducted  
377 additional analogous analyses separately for each group. The key pattern of results did not  
378 change, that is, the effect of speech stimulus type on attractiveness ratings was the same for  
379 both groups of raters. The results of these analyses are presented in Tables S5-S9 in the  
380 Supplementary Materials.

381  
382



383 *Dominance*

384 Overall average perceived dominance of female vocalizers, controlling for other predictors,  
385 was  $3.60 \pm 0.07$ . Speech stimulus type had a significant effect  $F(5, 11931.9) = 138.20, p$   
386  $<.001$ . Listeners rated the series of vowels ( $3.96 \pm 0.07$ ), followed by the greeting ( $3.81 \pm$   
387  $0.07$ ), counting ( $3.80 \pm 0.07$ ), the recited paragraph ( $3.70 \pm 0.07$ ), single words ( $3.40 \pm 0.07$ )  
388 and vowels ( $3.04 \pm 0.07$ ) as most dominant, respectively (Figure 1b, Table S10 and S11 for  
389 exact differences and their significance).

390 For male vocalizers, average perceived dominance, controlling for other predictors,  
391 was  $3.68 \pm 0.07$ . The type of voice stimulus also affected perceived dominance  $F(5, 11754) =$   
392  $80.33, p < .001$ . Listeners rated vocalizer's producing a greeting ( $3.89 \pm 0.07$ ), followed by  
393 the series of vowels ( $3.80 \pm 0.07$ ), counting ( $3.90 \pm 0.07$ ), recited paragraph ( $3.68 \pm 0.07$ ),  
394 single words ( $3.58 \pm 0.07$ ) and vowels ( $3.20 \pm 0.07$ ), as most dominant, respectively (see  
395 Tables S12 and S13 and Figure 1b). None of the covariates yielded a significant influence for  
396 dominance ratings of female nor male vocalizers (see Table 2).

397 There were no significant differences in ratings between raters who completed the  
398 playback experiment in the lab versus online for either vocalizer sex ( $b = -0.09, p = .36$  for  
399 female and  $b = -0.12, p = .21$  for male vocalizers). Adding speech stimulus type as a predictor  
400 did not significantly increase the amount of variance in dominance ratings explained by the  
401 model (0.05% and 0.03% for female and male voices, respectively).

402  
403

404 *Likability*

405 Overall average perceived likability of female vocalizers, controlling for other predictors, was  
406  $4.15 \pm 0.06$ . Speech stimulus type was a significant predictor  $F(5, 12363.7) = 174.11, p$   
407  $<.001$ . The order of likability ratings from highest to lowest was as follows: greetings ( $4.52 \pm$   
408  $0.06$ ), recited paragraph ( $4.45 \pm 0.06$ ), counting ( $4.32 \pm 0.06$ ), series of vowels ( $4.17 \pm 0.06$ ),  
409 single word ( $3.77 \pm 0.06$ ), and vowels pronounced individually ( $3.73 \pm 0.06$ ; see Figure 1c  
410 and corresponding Tables S14 and S15 in SOM). Older vocalizers were evaluated as  
411 relatively slightly less likeable ( $b = -0.01, p <.001$ , see Table 3 for estimates of all effects).

412 For male vocalizers, average perceived likability, controlling for other predictors, was  
413  $3.90 \pm 0.05$ . The effect of speech stimulus type was significant  $F(5, 12158) = 157.53, p <.001$ ,  
414 and almost identical for male and female vocalizers, with the only difference in male  
415 vocalizers indicating slightly higher ratings based on a series of vowels ( $3.80 \pm 0.07$ )  
416 compared to a single word ( $3.67 \pm 0.07$ ; Figure 2c). See Tables S16 and S17 for estimated  
417 marginal means and post-hoc tests. Male listeners judged voices as significantly less likeable  
418 than did female listeners ( $b = -0.27, p = .003$ ).

419 There were no significant differences in ratings between raters who completed the  
420 playback experiment in the lab versus online for either vocalizer sex ( $b = 0.04, p = .64$  for  
421 female and  $b = -0.01, p = .87$  for male vocalizers). Adding fixed effects of speech stimulus  
422 type into the models did not significantly increase the amount of variance in likability ratings  
423 explained by the model (0.07% and 0.06% for female and male voices, respectively).

424  
425  
426  
427

428 *Trustworthiness*

429 Overall average perceived trustworthiness of female vocalizers, controlling for other  
430 predictors, was  $3.79 \pm 0.06$ . Among female vocalizers, speech stimulus type also had a  
431 significant effect on perceived trustworthiness  $F(5, 11466.8) = 256.44, p < .001$ . The highest  
432 trustworthiness ratings were assigned to greetings ( $4.20 \pm 0.07$ ), then counting ( $4.12 \pm 0.07$ ),  
433 recited paragraph ( $4.10 \pm 0.07$ ), series of vowels ( $3.97 \pm 0.07$ ), single word ( $3.26 \pm 0.07$ ) and  
434 individual vowels ( $3.12 \pm 0.07$ ), respectively (Figure 1d, see also Tables S18 and S19 in SOM  
435 for estimated marginal means and pairwise comparisons). Covariates were not significant (see  
436 Table 4).

437 In male vocalizers, average perceived trustworthiness, controlling for other predictors,  
438 was  $3.68 \pm 0.06$ . Speech stimulus type had a significant effect on trustworthiness ratings  $F(5,$   
439  $11156) = 216.67, p < .001$ ). The highest trustworthiness ratings were assigned to vocalizers  
440 when reciting a paragraph ( $4.05 \pm 0.07$ ), followed by the greeting ( $4.01 \pm 0.07$ ), counting  
441 ( $4.00 \pm 0.07$ ), series of vowels ( $3.73 \pm 0.07$ ), single word ( $3.23 \pm 0.07$ ) and individual vowels  
442 ( $3.05 \pm 0.07$ ), respectively (see Table S20 and S21 for estimated marginal means and pairwise  
443 comparisons). None of the covariates yielded significant effects (see Table 4).

444 There were no significant differences in ratings between raters who completed the  
445 playback experiment in the lab versus online for either vocalizer sex ( $b = 0.01, p = .95$  for  
446 female and  $b = 0.04, p = .66$  for male vocalizers). Adding fixed effect of speech stimulus type  
447 into the model increased the amount of variance in trustworthiness ratings explained by the  
448 model only slightly (0.1% and 0.09% for female and male voices, respectively).

449

450 *Femininity-Masculinity*

451 The scale was coded such that higher scores indicate masculinity, and lower scores indicate  
452 femininity. In female vocalizers, overall average perceived femininity-masculinity,  
453 controlling for other predictors, was  $2.35 \pm 0.05$  (thus on the ‘feminine’ side of the scale).  
454 Speech stimulus type had a significant effect on perceived masculinity-femininity  $F(5,$   
455  $11704.9) = 75.31, p < .001$ . Single vowels ( $2.62 \pm 0.06$ ) followed by a word ( $2.52 \pm 0.06$ ) and  
456 vowels in a series ( $2.32 \pm 0.06$ ) elicited the highest masculinity ratings (least feminine).  
457 Longer utterances, including greetings ( $2.23 \pm 0.06$ ), counting ( $2.20 \pm 0.06$ ), and the recited  
458 paragraph ( $2.19 \pm 0.06$ ), respectively, were rated as least masculine (most feminine) (see  
459 Figure 1e and Tables S22 and S23 for means and pairwise comparisons). Higher ratings were,  
460 on average, assigned by slightly older raters ( $b = 0.01, p < .001$ , see Table 5).

461 The average femininity-masculinity rating in males, controlling for other predictors,  
462 was  $5.64 \pm 0.06$  (thus on the ‘masculine’ side of the scale). Speech stimulus type was  
463 significant  $F(5, 11459) = 35.35, p < .001$ . The pattern of results was opposite to that observed  
464 in female vocalizers. Specifically, longer utterances including the recited paragraph ( $5.68 \pm$   
465  $0.06$ ), counting ( $5.76 \pm 0.06$ ), and greeting ( $5.75 \pm 0.06$ ) yielded higher masculinity ratings  
466 whereas shorter speech samples yielded the lowest, including vowels in a series ( $5.66 \pm$   
467  $0.06$ ), single word ( $5.56 \pm 0.06$ ) and individual vowels ( $5.45 \pm 0.06$ ). For both sexes, the  
468 series of vowels elicited intermediate ratings (see Figure 1e). The opposing results show that  
469 longer utterances generally elicited higher sex-typical ratings (i.e., higher masculinity ratings  
470 for men and higher femininity ratings for women) (see Figure 1e; Tables S24 and S25). The  
471 effects of covariates were nonsignificant.

472 There were no significant differences in ratings between raters who completed the  
473 playback experiment in the lab versus online for either vocalizer sex ( $b = 0.18, p = .07$  for  
474 female and  $b = -0.08, p = .28$  for male vocalizers). Adding the fixed effect of speech stimulus

475 type into the model did not significantly increase the amount of variance in  
476 femininity/masculinity ratings explained by the model (0.03% and 0.01% for female and male  
477 voices, respectively).

478  
479

480 *Health*

481 Overall average perceived health of female vocalizers, controlling for other predictors, was  
482  $5.03 \pm 0.06$ . Among female vocalizers, stimulus type also predicted differences in voice-based  
483 health assessments  $F(5, 12691.1) = 143.40, p < .001$ . The perceived healthiness of vocalizers  
484 based on listening to their utterances increased in the following order: single vowels ( $4.53 \pm$   
485  $0.06$ ), single word ( $4.74 \pm 0.06$ ), series of vowels ( $5.23 \pm 0.06$ ), counting ( $5.25 \pm 0.06$ ),  
486 greetings ( $5.36 \pm 0.06$ ), and recited paragraph ( $5.19 \pm 0.06$ ; see Figure 1f and Tables S26 and  
487 S27 in Supplementary materials). The younger the vocalizers were, the higher their evaluated  
488 health ( $b = -0.02, p < .001$ ). The effect of the rater's group was not significant (see Table 6).

489 For male vocalizers, average perceived health, controlling for other predictors, was  $4.80 \pm$   
490  $0.06$ . Like female vocalizers, stimulus type significantly predicted judgements of men's health  
491  $F(5, 12355) = 94.66, p < .001$ . The perceived healthiness of male vocalizers based on listening  
492 to their utterances increased in a following order: single vowels ( $4.36 \pm 0.07$ ), single word  
493 ( $4.71 \pm 0.07$ ), series of vowels ( $4.78 \pm 0.07$ ), recited paragraph ( $4.95 \pm 0.07$ ), counting ( $5.04 \pm$   
494  $0.07$ ) and greetings ( $5.11 \pm 0.07$ ; see Tables S28 and S29 and Figure 1f for means and  
495 comparisons). Model summaries are given in Table 6. Younger male vocalizers were also  
496 evaluated as slightly healthier than were older male vocalizers ( $b = -0.01, p = .045$ ).

497 The effect of rater group was significant for male ( $b = -0.22, p = .023$ ) but not female  
498 vocalizers ( $b = -0.17, p = .06$ ), with online raters assessing male vocalizers' health as slightly  
499 higher than did lab raters ( $M/SEM = 4.93/0.02$  vs.  $4.72/0.02$ ). Adding fixed effects of speech  
500 stimulus type into the model did not substantially increase the amount of variance in health  
501 ratings explained by the model (0.05% and 0.04% for female and male voices, respectively).  
502 We present separate analyses of health judgments of male vocalizers made by online versus  
503 lab raters in the Supplementary Materials (Tables S30-S34). The key pattern of results did not

504 change, that is, the effect of speech stimulus type on health ratings was the same for both  
505 groups of raters.

506

507 *Additional analyses*

508 To corroborate our findings and to test whether inter-vocalizer differences in perceived  
509 speaker traits are preserved across stimuli, we proceeded with a series of Pearson's correlation  
510 analyses conducted on the vocalizer level, i.e., for each speaker, we averaged their ratings for  
511 each stimulus type separately and regressed these averaged scores on one another. The results  
512 are presented in Figures 2-8. All ratings were significantly and positively correlated at the  
513 inter-individual level, suggesting that individual differences in social judgments were  
514 preserved across stimulus types. For example, the average attractiveness rating that a given  
515 vocalizer received when producing one type of speech utterance explained 10%-74% of the  
516 variance in the attractiveness ratings given to that same vocalizer when producing any other  
517 utterance (Fig 2). The relationships were moderate to strong for all traits, wherein 9.6% to  
518 73% of variance in each rating dimension was explained across speech types. However, the  
519 fact that the correlations in many cases did not exceed  $r = 0.5$  implies that the ratings,  
520 although meaningfully correlated, are not identical.

521 In all cases, the weakest ( $r = 0.31-0.6$ ) correlations were observed between pairs  
522 consisting of one short (one syllable) utterance and one longer utterance (greeting, counting or  
523 recited paragraph; Table S25). In contrast, for almost all traits, the strongest correlations ( $r =$   
524  $0.64-0.86$ ) were observed between two longer utterances. In only two cases (dominance  
525 judgments of females and trustworthiness judgments of males), the strongest coefficient was  
526 found between a series of vowels (mid-length utterance) and a longer speech stimulus.

527  
528  
529  
530

531  
532  
533  
534  
535  
536



537 Discussion

538 Studies on perceptions of speakers based on the nonverbal properties of their voices vary  
539 methodologically owing to the use of a wide range of speech stimulus types, from single  
540 vowels to longer paragraphs of speech. This raises concerns regarding comparability of  
541 results across studies, namely whether different methodological choices are equally valid and  
542 whether listeners judge the exact same unseen speaker differently based on different types of  
543 neutral speech utterances, with a range of social and practical implications. The goal of the  
544 current study was to compare listeners' judgements of various traits of male and female  
545 vocalizers based on six content-neutral speech utterances, differing in duration and  
546 complexity, to test the degree to which variance in listeners' voice-based judgments is shared  
547 across speech types produced by the same person. In a large-scale playback experiment, we  
548 presented over 2000 raters with recordings of various utterances produced by the same  
549 approximately 200 vocalizers, and asked them to rate the unseen vocalizers on six different  
550 socially relevant dimensions. These included traits relevant to human mate selection and to  
551 other non-sexual social contexts (Pisanski & Feinberg, 2019). All of the considered traits have  
552 been of high interest to researchers studying the evolution and social outcomes of human  
553 vocal communication for decades (for reviews of the literature see: Aung & Puts, 2020;  
554 Kamiloğlu & Sauter, 2021; Kreiman & Sidtis, 2011; Oleszkiewicz et al., 2017; Pisanski &  
555 Bryant, 2019; Pisanski & Feinberg, 2019; Stern et al., 2021).

556 Our results show that while listeners' assessments of socially relevant traits are highly  
557 correlated across different neutral speech utterances produced by the same vocalizers,  
558 listeners' ratings generally increase for utterances of longer duration. Generally, vocalizers of  
559 both sexes are likely to be evaluated as more attractive, likeable, trustworthy, healthy, and  
560 dominant when producing longer than shorter utterances, particularly compared to single  
561 vowels or single words each comprised of only one syllable. In the case of perceived

562 femininity-masculinity, longer utterances elicited higher femininity ratings for women and  
563 higher masculinity ratings for men. At the same time, ratings among different types of longer  
564 utterances (reciting a paragraph aloud, counting, greeting, and in some cases, series of  
565 vowels; 150, 16, 8 and 5 syllables, respectively) were not remarkably different from one  
566 another. For instance, although the multilevel models show significant differences in ratings  
567 of dominance among men when comparing counting/greeting versus reciting a full paragraph,  
568 the differences in mean scores between these speech stimulus types are very small (0.22/0.21)  
569 and a high degree of variance is still shared between ratings based on these stimuli (37-42%).

570         Indeed, listeners' judgments shared a high degree of variance across stimulus types,  
571 within vocalizers. For example, a vocalizer who received high attractiveness ratings based on  
572 her or his vowel series was likely to also be judged as highly attractive when producing a full  
573 paragraph. However, based on vocalizer-level analyses, relatively more variance was shared  
574 between pairs of two long utterances than between two short utterances, or than between a  
575 short and long utterance. Importantly, then, while we show that longer utterances elicit  
576 relatively higher ratings of attractiveness, trustworthiness, likability, health, dominance,  
577 masculinity in males, and femininity in females, we also show moderate to strong  
578 relationships between ratings of the same individuals across speech types, indicating stability  
579 intra-individual stability in speaker perception. The strongest relationships were observed  
580 between stimulus types of a similar length, wherein a given speech type could explain  
581 upwards of 73% of the variance in ratings of the same vocalizers based on another speech  
582 type.

583         Although the relationships between pairs of short-long utterances were weaker than  
584 between pairs of longer utterances, they were still significant and never lower than  $r = .26$ .  
585 While acoustic parameters were not measured in the present study, this finding raises the  
586 possibility that listeners' social judgments are tapping into the same underlying acoustic

587 properties in vocalizer's voices, whether those vocalizers are producing vowels or a  
588 paragraph, however that longer utterances provide additional information that listeners readily  
589 utilize. These results are in line with previous research reporting that listeners can judge a  
590 range of personality traits from nothing more than a single utterance (e.g., "hello", McAleer et  
591 al., 2014) and that ratings of trustworthiness, dominance, and attractiveness are highly  
592 correlated when based on a word vs. sentence (Mahrholz et al., 2018). At the same time,  
593 corroborating a previous study showing higher attractiveness ratings for words compared to  
594 single vowels (Ferdenzi et al., 2013), our results provide novel insight into the importance of  
595 stimulus duration on listeners' perceptions.

596         There is growing evidence from acoustic analyses of nonverbal vocal parameters of  
597 the human voice that certain vocal parameters are stable across speech types. For example,  
598 Pisanski and colleagues (2021) found that fundamental frequency ( $f_0$ ), perceived as voice  
599 pitch and one of the most extensively studied and socially meaningful acoustic parameters in  
600 the human voice (Aung & Puts, 2020), is highly stable across neutral utterances of different  
601 lengths produced by the same vocalizers. At least half and up to 80% of the variance in  $f_0$   
602 measured from one utterance was explained by the  $f_0$  of any other utterance within speakers.  
603 In another study, Pisanski and colleagues show that these inter-individual differences in  $f_0$   
604 also extend to emotional speech and nonverbal vocalisations such as screams, roars and cries  
605 produced by the same men and women (Pisanski et al., 2020). This suggests that individual  
606 differences in voice pitch (often measured by researchers to test for relationships between  $f_0$   
607 and vocalizer traits such as body size, testosterone levels, attractiveness, or dominance, to  
608 name a few) are robust regardless of speech type or duration, within the same group of  
609 vocalizers. However, long and short speech recordings can differ in a number of other  
610 characteristics like formant patterns, prosody, articulation, or speed of speech (Leung et al.,  
611 2018). Longer utterances are likely to convey more information than shorter utterances

612 regarding not only the nonverbal parameters of the vocalizer's voice but also their prosody or  
613 the vocalizer's way of speaking, which might be perceptually related to personality (Zellner  
614 Keller, 2005). Indeed, voice pitch only partly explains the personality judgements of unseen  
615 vocalizers (Stern et al., 2021).

616         Short utterances such as vowels and counting have traditionally been used in voice  
617 perception studies due to their contextual neutrality, standardized nature, and thus high  
618 experimental control. Despite the brevity and neutrality of these voice stimuli, studies have  
619 generally shown that listeners can gauge various social traits from short utterances such as  
620 vowels or a single word (Apicella & Feinberg, 2009; McAleer et al., 2014; Pisanski et al.,  
621 2014). Nevertheless, there exists an important difference between very short utterances (like  
622 single vowels) and longer utterances: the former are typically less ecologically valid. In real  
623 life conversations, people are unlikely to base judgements solely on single vowels, as they  
624 rarely hear them in separation from the rest of the statement. Therefore, another possibility is  
625 that voice recordings consisting of one vowel, or a neutral one-syllable word, may be  
626 relatively unfamiliar to listeners and, in turn, may thus be evaluated less positively.

627

#### 628 *Limitations and future research recommendations*

629         Our research design does not allow to draw direct conclusions regarding the  
630 'accuracy' of listeners' judgments, as no objective measures of speaker traits were obtained.  
631 As such, the results cannot speak to the question of superiority of one methodological  
632 stimulus choice over the other. However, a close look at the result patterns shows that short  
633 utterances are in almost all cases closer to the scale mean (a rating of 4; for an exception see  
634 attractiveness among male vocalizers). This suggests that listeners might have been less  
635 certain about their ratings of shorter versus longer stimuli, resulting in judgments closer to  
636 what they might have perceived as "average". The geographic homogeneity of our vocalizer

637 and listener samples, paired with a lack of data about various socially relevant participant  
638 variables including their gender identity, sexual orientation and socioeconomic status,  
639 represent limitations of our study that, if known, could broaden its generalizability. Regarding  
640 the study design, we observed a tendency for raters to make judgements closer to the scale  
641 mean for shorter than for longer utterances, possibly reflecting a level of uncertainty for short  
642 speech utterances. Providing raters with the option to omit the question, to mark “I don’t  
643 know”, or to rate the confidence of their ratings, could help to clarify the mechanisms driving  
644 the small observed differences in ratings across speech stimulus types. Moreover, although  
645 online raters were instructed to use high quality headphones, and headphone use was verified  
646 with hearing tests, we cannot be certain that the quality of their headphones was comparable  
647 to that of those used by lab participants. Our results suggest, however, that the online sample  
648 of raters produced qualitatively similar ratings to the lab sample that used professional  
649 headphones.

650       Our results may not necessarily generalize to judgments of other traits not tested here,  
651 for example, the Big Five personality traits (McCrae & Costa, 1989). Traits like extraversion  
652 or neuroticism may be especially closely perceptually interrelated with the prosody of speech  
653 (Feldstein & Sloan, 1984). Future studies could therefore extend the list of evaluated traits.  
654 We also limited our analyses to judgments of traits described by a single word, such as  
655 ‘trustworthiness’, wherein people may differ in the way they understand such terms. For  
656 masculinity and femininity, for example, some listeners may associate these terms with  
657 biological sex while others may judge masculinity and femininity independently of the  
658 perceived gender of the vocalizer. While using single undefined terms on rating scales is  
659 common practice in this research domain (Cartei et al., 2014; Feinberg et al., 2012; Hughes et  
660 al., 2014; McAleer et al., 2014), researchers may provide a working definition of these  
661 constructs or assess them through behavioral measures (e.g., probability of certain behaviors

662 or implicit association tasks). Moreover, femininity and masculinity were represented on a  
663 single rating scale, rather than representing two qualitatively independent constructs. We  
664 acknowledge that femininity and masculinity do not necessarily represent two ends of one  
665 continuum (Bem, 1974; Donnelly & Twenge, 2017) and that in future research may be treated  
666 separately.

667 It may also be of practical interest to test the ‘accuracy’ of social judgements. This is  
668 methodologically challenging because of the imperfection of measurement scales and the  
669 social approval factor that biases self-ratings toward more desired profiles. However,  
670 researchers may be able to commit to reliable proxies (behavioral or self-assessed) to test how  
671 accurately social traits can be inferred from voice samples depending on speech stimulus type.  
672 Another important avenue for future studies will be to test the effect of stimulus type in more  
673 ecologically valid or multi-modal conditions, i.e., where a speaker is seen (or smelled).  
674 Different modalities (e.g., voice, physical appearance or body odor) interact with one other to  
675 jointly affect impressions we form of people (for reviews and discussion see Feinberg, 2008;  
676 Groyecka et al., 2017; Krumpholz et al., 2021). Multi-modal studies can help to clarify  
677 whether speech complexity affects ratings of person dimensions to the same extent when  
678 voice stimuli are perceived alongside other modalities as when they are perceived in isolation.

679

#### 680 *Conclusion: Methodological, social and practical implications*

681 We show that inter-individual differences in voice-based judgments are relatively stable  
682 across neutral speech stimuli, particularly relatively longer utterances, for the same vocalizers.  
683 Interestingly, while single-syllable voice stimuli elicit significantly lower ratings than do  
684 multi-syllable stimuli, a stimulus with 16 syllables of length appears to be comparable with a  
685 150-syllable one. This suggests that in terms of playback experiments, using lengthy excerpts  
686 from passages like the Rainbow Passage, although well established and commonly used

687 (Cartei et al., 2014; Pisanski et al., 2016; Schild et al., 2020; Tigue et al., 2012), can  
688 unnecessarily extend the time of the experimental procedure while being similarly  
689 informative to listeners as one short sentence or counting from one to ten. Judgments based on  
690 five vowels are also quite comparable to full sentences and paragraphs, at least in their effects  
691 on listeners' social judgments of vocalizers.

692         Our results also corroborate earlier acoustic analyses and playback studies showing  
693 that a great deal of information about a person can be encoded in very short speech segments.  
694 Vowel sounds, for example, encode enough acoustic information, namely in voice  
695 fundamental and formant frequencies, for listeners to reliably judge static speaker traits such  
696 as body size, age, and biological sex, owing largely to anatomical or physiological constraints  
697 on vocal production (for reviews see Aung & Puts, 2020; Charlton et al., 2020; Pisanski &  
698 Bryant, 2019). Listeners also show high agreement on judgments of social or personality traits  
699 from a single word (McAleer et al., 2014). Our results suggest that listeners' judgments do not  
700 change drastically if they are presented with longer bouts of neutral speech from the same  
701 speaker, indicating that voice-based person impressions, which can meaningfully impact  
702 social and societal outcomes, may be formed early during first interactions and/or may rely  
703 heavily on low-level and relatively static acoustic features that do not vary greatly across  
704 speech utterances.

705         Finally, paired with emerging evidence from acoustic analyses that show remarkable  
706 intra-individual stability in people's voices across time and context (Fouquet et al., 2016;  
707 Levvero et al., 2018; Pisanski et al., 2020; 2021), voice-based human perception studies can  
708 inform practical technologies, including automated voice recognition devices in mobile  
709 phones and computers that increasingly rely on large amounts of human user response data to  
710 build high-performance predictive algorithms using artificial intelligence and deep learning  
711 (Deng, 2018).

712

713 Data Availability Statement

714 Dataset and measurement instruments are available at

715 [https://osf.io/cevpd/?view\\_only=e3ecf2aad04d4e01be5314e13c3446b1](https://osf.io/cevpd/?view_only=e3ecf2aad04d4e01be5314e13c3446b1).

716

717 References

718 Albert, G., Arnocky, S., Puts, D. A., & Hodges-Simeon, C. R. (2021). Can listeners assess

719 men's self-reported health from their voice? *Evolution and Human Behavior*, 42(2), 91–

720 103. <https://doi.org/10.1016/j.evolhumbehav.2020.08.001>

721 Anderson, R. C., Klofstad, C. A., Mayew, W. J., & Venkatachalam, M. (2014). Vocal fry may

722 undermine the success of young women in the labor market. *PLoS ONE*, 9(5), 1–8.

723 <https://doi.org/10.1371/journal.pone.0097506>

724 Apicella, C. L., & Feinberg, D. R. (2009). Voice pitch alters mate-choice-relevant perception

725 in hunter–gatherers. *Proceedings of the Royal Society B: Biological Sciences*, 276(1659),

726 1077–1082.

727 Aung, T., & Puts, D. (2020). Voice pitch: a window into the communication of social power.

728 *Current Opinion in Psychology*, 33, 154–161.

729 Belin, P., Boehme, B., & McAleer, P. (2017). The sound of trustworthiness: Acoustic-based

730 modulation of perceived voice personality. *PLoS ONE*, 12(10), 4–12.

731 <https://doi.org/10.1371/journal.pone.0185651>

732 Bem, S. L. (1974). The measurement of psychological androgyny. *Journal of Consulting and*

733 *Clinical Psychology*, 42(2), 155.

734 Borkowska, B., & Pawlowski, B. (2011). Female voice frequency in the context of dominance

735 and attractiveness perception. *Animal Behaviour*, 82(1), 55–59.

736 <https://doi.org/10.1016/j.anbehav.2011.03.024>



- 737 Bryant, G. A., & Haselton, M. G. (2009). Vocal cues of ovulation in human females. *Biology*  
738 *Letters*, 5(1), 12–15.
- 739 Cartei, V., Bond, R., & Reby, D. (2014). What makes a voice masculine: Physiological and  
740 acoustical correlates of women's ratings of men's vocal masculinity. *Hormones and*  
741 *Behavior*, 66(4), 569–576. <https://doi.org/10.1016/j.yhbeh.2014.08.006>
- 742 Charlton, B. D., Pisanski, K., Raine, J., & Reby, D. (2020). Coding of static information in  
743 terrestrial mammal vocal signals. In *Coding strategies in vertebrate acoustic*  
744 *communication* (pp. 115–136). Springer.
- 745 Deng, L. (2018). Artificial intelligence in the rising wave of deep learning: The historical path  
746 and future outlook [perspectives]. *IEEE Signal Processing Magazine*, 35(1), 177–180.
- 747 Donnelly, K., & Twenge, J. M. (2017). Masculine and feminine traits on the Bem Sex-Role  
748 Inventory, 1993–2012: A cross-temporal meta-analysis. *Sex Roles*, 76(9), 556–565.
- 749 Fairbanks, G. (1960). *Voice and articulation drillbook*. Addison-Wesley Educational  
750 Publishers.
- 751 Feinberg, D. R., Jones, B. C., Little, A. C., Burt, D. M., & Perrett, D. I. (2005). Manipulations  
752 of fundamental and formant frequencies influence the attractiveness of human male  
753 voices. *Animal Behaviour*, 69(3), 561–568.  
754 <https://doi.org/10.1016/j.anbehav.2004.06.012>
- 755 Feinberg, D. R. (2008). Are human faces and voices ornaments signaling common underlying  
756 cues to mate value? *Evolutionary Anthropology: Issues, News, and Reviews: Issues,*  
757 *News, and Reviews*, 17(2), 112–118.
- 758 Feinberg, D. R., DeBruine, L. M., Jones, B. C., Little, A. C., O'Connor, J. J. M., & Tigue, C.  
759 C. (2012). Women's self-perceived health and attractiveness predict their male vocal  
760 masculinity preferences in different directions across short-and long-term relationship  
761 contexts. *Behavioral Ecology and Sociobiology*, 66(3), 413–418.

762 Feinberg, D. R., DeBruine, L. M., Jones, B. C., & Perrett, D. I. (2008). The role of femininity  
763 and averageness of voice pitch in aesthetic judgments of women's voices. *Perception*,  
764 37(4), 615–623.

765 Feldstein, S., & Sloan, B. (1984). Actual and stereotyped speech tempos of extraverts and  
766 introverts. *Journal of Personality*, 52(2), 188–204.

767 Ferdenzi, C., Patel, S., Mehu-Blantar, I., Khidasheli, M., Sander, D., & Delplanque, S. (2013).  
768 Voice attractiveness: Influence of stimulus duration and type. *Behavior Research*  
769 *Methods*, 45(2), 405–413. <https://doi.org/10.3758/s13428-012-0275-0>

770 Fouquet, M., Pisanski, K., Mathevon, N., & Reby, D. (2016). Seven and up: individual  
771 differences in male voice fundamental frequency emerge before puberty and remain  
772 stable throughout adulthood. *Royal Society Open Science*, 3(10), 160395.

773 Groyecka, A., Pisanski, K., Sorokowska, A., Havlíček, J., Karwowski, M., Puts, D., Craig  
774 Roberts, S., & Sorokowski, P. (2017). Attractiveness is multimodal: Beauty is also in the  
775 nose and ear of the beholder. *Frontiers in Psychology*, 8(MAY).  
776 <https://doi.org/10.3389/fpsyg.2017.00778>

777 Hughes, S. M., Mogilski, J. K., & Harrison, M. A. (2014). The perception and parameters of  
778 intentional voice manipulation. *Journal of Nonverbal Behavior*, 38(1), 107–127.

779 Jones, B. C., Feinberg, D. R., DeBruine, L. M., Little, A. C., & Vukovic, J. (2008).  
780 Integrating cues of social interest and voice pitch in men's preferences for women's  
781 voices. *Biology Letters*, 4(2), 192–194. <https://doi.org/10.1098/rsbl.2007.0626>

782 Kamiloğlu, R. G., & Sauter, D. A. (2021). Voice Production and Perception. In *Oxford*  
783 *Research Encyclopedia of Psychology*.

784 Klothstad, C. A., Anderson, R. C., & Peters, S. (2012). Sounds like a winner: Voice pitch  
785 influences perception of leadership capacity in both men and women. *Proceedings of the*  
786 *Royal Society B: Biological Sciences*, 279(1738), 2698–2704.

787 <https://doi.org/10.1098/rspb.2012.0311>

788 Kordsmeyer, T. L., Hunt, J., Puts, D. A., Ostner, J., & Penke, L. (2018). The relative  
789 importance of intra-and intersexual selection on human male sexually dimorphic traits.  
790 *Evolution and Human Behavior*, 39(4), 424–436.

791 Krahé, B., Uhlmann, A., & Herzberg, M. (2021). The voice gives it away: Male and female  
792 pitch as a cue for gender stereotyping. *Social Psychology*, 52(2), 101–113.  
793 <https://doi.org/10.1027/1864-9335/a000441>

794 Kreiman, J. (1997). Listening to voices: theory and practice in voice perception  
795 research. *Talker variability in speech processing*, 85-108.

796 Kreiman, J., & Sidtis, D. (2011). *Foundations of voice studies: An interdisciplinary approach*  
797 *to voice production and perception*. John Wiley & Sons.

798 Krumholz, C., Quigley, C., Little, A. C., Zäske, R., & Riebel, K. (2021). Multimodal  
799 signalling of attractiveness. *Proceedings of the Annual Meeting of the Cognitive Science*  
800 *Society*, 43(43).

801 Leung, Y., Oates, J., & Chan, S. P. (2018). Voice, articulation, and prosody contribute to  
802 listener perceptions of speaker gender: A systematic review and meta-analysis. *Journal*  
803 *of Speech, Language, and Hearing Research*, 61(2), 266–297.

804 Levrero, F., Mathevon, N., Pisanski, K., Gustafsson, E., & Reby, D. (2018). The pitch of  
805 babies’ cries predicts their voice pitch at age 5. *Biology Letters*, 14(7), 20180065.

806 Mahrholz, G., Belin, P., & McAleer, P. (2018). Judgements of a speaker’s personality are  
807 correlated across differing content and stimulus type. *PLoS ONE*, 13(10), 1–22.  
808 <https://doi.org/10.1371/journal.pone.0204991>

809 McAleer, P., Todorov, A., & Belin, P. (2014). How do you say “hello”? Personality  
810 impressions from brief novel voices. *PLoS ONE*, 9(3), 1–9.  
811 <https://doi.org/10.1371/journal.pone.0090779>

812 McCrae, R. R., & Costa, P. T. (1989). The structure of interpersonal traits: Wiggins's  
813 circumplex and the five-factor model. *Journal of Personality and Social Psychology*,  
814 56(4), 586.

815 Mileva, M., Tompkinson, J., Watt, D., & Burton, A. M. (2020). The role of face and voice  
816 cues in predicting the outcome of student representative elections. *Personality and Social  
817 Psychology Bulletin*, 46(4), 617–625. <https://doi.org/10.1177/0146167219867965>

818 Ohno, R., Masanori, K., & Tetsuro, M. (2014). Relationship between perception of cuteness  
819 in female voices and their durations. *International Conference on Speech and Computer*,  
820 8773. <https://doi.org/10.1007/978-3-319-11581-8>

821 Oleszkiewicz, A., Pisanski, K., Lachowicz-Tabaczek, K., & Sorokowska, A. (2017). Voice-  
822 based assessments of trustworthiness, competence, and warmth in blind and sighted  
823 adults. *Psychonomic Bulletin & Review*, 24(3), 856–862.

824 Pisanski, K., Anikin, A., & Reby, D. (2021). Static and dynamic formant scaling conveys  
825 body size and aggression. *Royal Society Open Science*, *The Royal Society*.  
826 <https://doi.org/10.1098/rsos.211496> . hal-03501112

827 Pisanski, K., & Bryant, G. A. (2019). The evolution of voice perception. *The oxford handbook  
828 of voice studies*, 269-300.

829 Pisanski, K., & Feinberg, D. R. (2019). Vocal attractiveness. *Oxford Handbooks. The Oxford  
830 Handbook of Voice Perception*, 606–626.

831 Pisanski, K., Fraccaro, P. J., Tigue, C. C., Connor, J. J. M. O., David, R., Pisanski, K.,  
832 Fraccaro, P. J., Tigue, C. C., Connor, J. J. M. O., & Feinberg, D. R. (2014). Return to  
833 Oz: Voice pitch facilitates assessments of men's body size. *Journal of Experimental  
834 Psychology: Human Perception and Performance*, 40(4), 1316.

835 Pisanski, K., Groyecka-Bernard, A., & Sorokowski, P. (2021). Human voice pitch measures  
836 are robust across a variety of speech recordings: methodological and theoretical

837 implications. *Biology Letters*, 17, 20210356.  
838 <https://doi.org/https://doi.org/10.1098/rsbl.2021.0356>

839 Pisanski, K., Mishra, S., & Rendall, D. (2012). The evolved psychology of voice: Evaluating  
840 interrelationships in listeners' assessments of the size, masculinity, and attractiveness of  
841 unseen speakers. *Evolution and Human Behavior*, 33(5), 509–519.  
842 <https://doi.org/10.1016/j.evolhumbehav.2012.01.004>

843 Pisanski, K., Nowak, J., & Sorokowski, P. (2016). Individual differences in cortisol stress  
844 response predict increases in voice pitch during exam stress. *Physiology and Behavior*,  
845 163, 234–238. <https://doi.org/10.1016/j.physbeh.2016.05.018>

846 Pisanski, K., Oleszkiewicz, A., Plachetka, J., Gmiterek, M., & Reby, D. (2018). Voice pitch  
847 modulation in human mate choice. *Proceedings of the Royal Society B: Biological*  
848 *Sciences*, 285(1893). <https://doi.org/10.1098/rspb.2018.1634>

849 Pisanski, K., Raine, J., & Reby, D. (2020). Individual differences in human voice pitch are  
850 preserved from speech to screams, roars and pain cries. *Royal Society Open Science*,  
851 7(2), 191642.

852 Pisanski, K., & Rendall, D. (2011). The prioritization of voice fundamental frequency or  
853 formants in listeners' assessments of speaker size, masculinity, and attractiveness. *The*  
854 *Journal of the Acoustical Society of America*, 129(4), 2201–2212.  
855 <https://doi.org/10.1121/1.3552866>

856 Puts, D. A., Hill, A. K., Bailey, D. H., Walker, R. S., Rendall, D., Wheatley, J. R., Welling, L.  
857 L. M., Dawood, K., Cárdenas, R., Burriss, R. P., Jablonski, N. G., Shriver, M. D., Weiss,  
858 D., Lameira, A. R., Apicella, C. L., Owren, M. J., Barelli, C., Glenn, M. E., & Ramos-  
859 Fernandez, G. (2016). Sexual selection on male vocal fundamental frequency in humans  
860 and other anthropoids. *Proceedings of the Royal Society B: Biological Sciences*,  
861 283(1829), 0–7. <https://doi.org/10.1098/rspb.2015.2830>

862 Puts, D. Andrew, Gaulin, S. J. C., & Verdolini, K. (2006). Dominance and the evolution of  
863 sexual dimorphism in human voice pitch. *Evolution and Human Behavior*, 27(4), 283–  
864 296. <https://doi.org/10.1016/j.evolhumbehav.2005.11.003>

865 Rosenfield, K. A., Sorokowska, A., Sorokowski, P., & Puts, D. A. (2020). Sexual selection  
866 for low male voice pitch among Amazonian forager-horticulturists. *Evolution and*  
867 *Human Behavior*, 41(1), 3–11.

868 Schild, C., Stern, J., & Zettler, I. (2020). Linking men’s voice pitch to actual and perceived  
869 trustworthiness across domains. *Behavioral Ecology*, 31(1), 164–175.  
870 <https://doi.org/10.1093/beheco/arz173>

871 Sell, A., Bryant, G. A., Cosmides, L., Tooby, J., Sznycer, D., Von Rueden, C., Krauss, A., &  
872 Gurven, M. (2010). Adaptations in humans for assessing physical strength from the  
873 voice. *Proceedings of the Royal Society B: Biological Sciences*, 277(1699), 3509–3518.

874 Sorokowski, P., Puts, D., Johnson, J., Żółkiewicz, O., Oleszkiewicz, A., Sorokowska, A.,  
875 Kowal, M., Borkowska, B., & Pisanski, K. (2019). Voice of authority: Professionals  
876 lower their vocal frequencies when giving expert advice. *Journal of Nonverbal Behavior*,  
877 43(2), 257–269. <https://doi.org/10.1007/s10919-019-00307-0>

878 Stern, J., Schild, C., Jones, B. C., DeBruine, L. M., Hahn, A., Puts, D. A., Zettler, I.,  
879 Kordsmeyer, T. L., Feinberg, D., Zamfir, D., Penke, L., & Arslan, R. C. (2021). Do  
880 voices carry valid information about a speaker’s personality? *Journal of Research in*  
881 *Personality*, 104092. <https://doi.org/10.1016/j.jrp.2021.104092>

882 Tigue, C. C., Borak, D. J., O’Connor, J. J. M., Schandl, C., & Feinberg, D. R. (2012). Voice  
883 pitch influences voting behavior. *Evolution and Human Behavior*, 33(3), 210–216.  
884 <https://doi.org/10.1016/j.evolhumbehav.2011.09.004>

885 Titze, I. R. (1989). Physiologic and acoustic differences between male and female voices.  
886 *Journal of the Acoustical Society of America*, 85(4), 1699–1707.

887 <https://doi.org/10.1121/1.397959>

888 Vukovic, J., Feinberg, D., Debruine, L., Smith, F., & Jones, B. (2010). Women's voice pitch

889 is negatively correlated with health risk factors. *Journal of Evolutionary Psychology*,

890 8(3), 217–225. <https://doi.org/10.1556/JEP.8.2010.3.2>

891 Zellner Keller, B. (2005). Speech prosody, voice quality and personality. *Logopedics*

892 *Phoniatrics Vocology*, 30(2), 72–78. <https://doi.org/10.1080/14015430500256543>

893

894

895 **Figure 1.** Listeners' judgments generally increased with speech stimulus duration. Differences in  
896 perceived vocalizer traits across speech stimulus types for listeners' judgments of (a) Attractiveness;  
897 (b) Dominance; (c) Likability; (d) Trustworthiness; (e) Femininity-masculinity; (f) Health. For exact  
898 differences see Tables S2, S4, S11, S13, S15, S17, S19, S21, S23, S25, S29 and S29 in SOM. Error  
899 bars depict 95% CI. Dashed lines depict nonsignificant differences based on pairwise comparisons  
900 following Bonferroni correction for multiple comparisons (see supplementary tables listed above). All  
901 remaining pairwise comparisons (those without dashed lines joining them) are statistically significant  
902 at  $p < .05$ .

903 **Figure 2.** Distributions and correlations (Pearson's  $r$  coefficient) between ratings of **attractiveness**  
904 based on six different stimulus types. Red – female vocalizers, green – male vocalizers.

905 **Figure 3.** Distributions and correlations (Pearson's  $r$  coefficient) between ratings of **dominance**  
906 based on six different stimulus types. Red – female vocalizers, green – male vocalizers.

907  
908 **Figure 4.** Distributions and correlations (Pearson's  $r$  coefficient) between ratings of **likability** based  
909 on six different stimulus types. Red – female vocalizers, green – male vocalizers.

910  
911 **Figure 5.** Distributions and correlations (Pearson's  $r$  coefficient) between ratings of **trustworthiness**  
912 based on six different stimulus types. Red – female vocalizers, green – male vocalizers.

913  
914 **Figure 6.** Distributions and correlations (Pearson's  $r$  coefficient) between ratings of **femininity-**  
915 **masculinity** based on six different stimulus types. Higher scores denote higher masculinity (lower  
916 femininity). Red – female vocalizers, green – male vocalizers.

917 **Figure 7.** Distributions and correlations (Pearson's  $r$  coefficient) between ratings of **health** based on  
918 six different stimulus types. Red – female vocalizers, green – male vocalizers.

919  
920



Table 1. Estimated fixed and random effects of the model with perceived **attractiveness** as an outcome variable

## FEMALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.8816	0.06044	3.7630	4.0002	313	64.142	< .001
Stimulus type	2 – 1	0.6995	0.03659	0.6278	0.7712	12515	19.118	< .001
	3 – 1	0.1418	0.03659	0.0701	0.2135	12515	3.875	< .001
	4 – 1	0.8296	0.03659	0.7579	0.9013	12515	22.675	< .001
	5 – 1	0.7620	0.03659	0.6903	0.8337	12515	20.827	< .001
	6 – 1	0.8949	0.03659	0.8232	0.9666	12515	24.459	< .001
Vocalizer's age		-0.02237	0.00340	-0.0290	-0.0157	100	-6.592	< .001
Listener's sex (0 – F, 1 – M)		-0.0129	0.08954	-0.1770	0.1744	432	-0.0144	0.988
Rater group (0 – online, 1 – lab)		-0.17260	0.09041	-0.3498	0.00460	432	-1.9090	0.057
Random effects			Variance	ICC				
Listener's ID			0.808	0.355				
Vocalizer's ID			0.170	0.105				
Residuals			1.457					
N <sub>observations</sub> = 13055, N <sub>listeners</sub> = 435, N <sub>vocalizers</sub> = 101								

## MALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.52695	0.06037	3.4086	3.64527	353	58.43	< .001
Stimulus type	2 – 1	0.41500	0.03543	0.3456	0.48446	12260	11.71	< .001
	3 – 1	0.19720	0.03543	0.1278	0.26664	12260	5.57	< .001
	4 – 1	0.68385	0.03543	0.6144	0.75329	12260	19.30	< .001
	5 – 1	0.73959	0.03543	0.6701	0.80903	12260	20.87	< .001
	6 – 1	0.80375	0.03543	0.7343	0.87320	12260	22.69	< .001
Vocalizer's age		-0.00860	0.00359	-0.0156	-0.00157	104	-2.40	0.018
Listener's sex (0 – F, 1 – M)		-0.50296	0.09230	-0.6639	-0.30206	432	-5.45	< .001
Rater group (0 – online, 1 – lab)		-0.21543	0.09311	-0.3979	-0.03295	432	-2.31	0.021
Random effects			Variance	ICC				
Listener's ID			0.857	0.390				
Vocalizer's ID			0.165	0.110				
Residuals			1.338					
N <sub>observations</sub> = 12805, N <sub>listeners</sub> = 435, N <sub>vocalizers</sub> = 106								

Note. 1 = single vowels, 2 = series of vowels, 3 = single word, 4 = counting, 5 = greeting, 6 = recited paragraph (Rainbow Passage)

Table 2. Estimated fixed and random effects of the model with perceived **dominance** as an outcome variable

FEMALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.60406	0.06661	3.4735	3.73463	257.8	54.101	< .001
Stimulus type	2 – 1	0.91370	0.04053	0.8343	0.99315	11932.0	22.542	< .001
	3 – 1	0.36286	0.04053	0.2834	0.44230	11932.0	8.952	< .001
	4 – 1	0.76045	0.04053	0.6810	0.83989	11932.0	18.761	< .001
	5 – 1	0.76623	0.04053	0.6868	0.84568	11932.0	18.904	< .001
	6 – 1	0.65394	0.04053	0.5745	0.73339	11932.0	16.134	< .001
Vocalizer's age		-0.00255	0.00400	-0.0104	0.00528	98.4	-0.638	0.525
Listener's sex (0 – F, 1 – M)		-0.09743	0.09306	-0.2688	0.09732	411.6	-0.918	0.359
Rater group (0 – online, 1 – lab)		-0.08572	0.09339	-0.2688	0.09732	411.6	-0.918	0.359
Random effects			Variance	ICC				
Listener's ID			0.798	0.319				
Vocalizer's ID			0.239	0.123				
Residuals			1.706					
N <sub>observations</sub> = 12455, N <sub>listeners</sub> = 415, N <sub>vocalizers</sub> = 101								

MALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.68345	0.06855	3.54908	3.8178	241	53.730	< .001
Stimulus type	2 – 1	0.59922	0.04153	0.51783	0.6806	11755	14.429	< .001
	3 – 1	0.38194	0.04153	0.30054	0.4633	11755	9.197	< .001
	4 – 1	0.70030	0.04153	0.61890	0.7817	11755	16.863	< .001
	5 – 1	0.68955	0.04153	0.60816	0.7709	11755	16.604	< .001
	6 – 1	0.47569	0.04153	0.39429	0.5571	11755	11.454	< .001
Vocalizer's age		0.00154	0.00471	-0.00769	0.0108	104	0.327	0.745
Listener's sex (0 – F, 1 – M)		-0.03377	0.09083	-0.21180	0.1443	412	-0.372	0.710
Rater group (0 – online, 1 – lab)		0.11532	0.09112	-0.06328	0.2939	412	1.266	0.206
Random effects			Variance	ICC				
Listener's ID			0.755	0.300				
Vocalizer's ID			0.289	0.141				
Residuals			1.763					
N <sub>observations</sub> = 12283 N <sub>listeners</sub> = 415, N <sub>vocalizers</sub> = 106								

Note. 1 = single vowels, 2 = series of vowels, 3 = single word, 4 = counting, 5 = greeting, 6 = recited paragraph (Rainbow Passage)

Table 3. Estimated fixed and random effects of the model with perceived **likability** as an outcome variable  
 FEMALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		4.1529	0.05688	4.0414	4.26435	323.8	73.01	<.001
Stimulus type	2 – 1	0.4423	0.03652	0.3707	0.51391	12363.7	12.11	<.001
	3 – 1	0.0372	0.03652	-0.0344	0.10880	12363.7	1.02	0.308
	4 – 1	0.5898	0.03652	0.5182	0.66135	12363.7	16.15	<.001
	5 – 1	0.7944	0.03652	0.7228	0.86601	12363.7	21.75	<.001
	6 – 1	0.7181	0.03652	0.6466	0.78973	12363.7	19.66	<.001
Vocalizer's age		-0.0116	0.00313	-0.0177	-0.00546	99.0	-3.71	<.001
Listener's sex (0 – F, 1 – M)		-0.1563	0.08708	-0.3270	0.01437	426.7	-1.795	0.073
Rater group (0 – online, 1 – lab)		0.0414	0.08764	-0.1304	0.21320	426.7	0.473	0.637
Random effects			Variance	ICC				
Listener's ID			0.737	0.3395				
Vocalizer's ID			0.142	0.0901				
Residuals			1.434					
N <sub>observations</sub> = 12900, N <sub>listeners</sub> = 430, N <sub>vocalizers</sub> = 101								

#### MALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.89720	0.05298	3.79337	4.00104	412	73.563	<.001
Stimulus type	2 – 1	0.30043	0.03547	0.23091	0.36994	12158	8.470	<.001
	3 – 1	0.17241	0.03547	0.10290	0.24193	12158	4.861	<.001
	4 – 1	0.55598	0.03547	0.48646	0.62549	12158	15.676	<.001
	5 – 1	0.76996	0.03547	0.70044	0.83947	12158	21.709	<.001
	6 – 1	0.74020	0.03547	0.67068	0.80971	12158	20.870	<.001
Vocalizer's age		0.00120	0.00283	-0.00434	0.00674	105	0.424	0.673
Listener's sex (0 – F, 1 – M)		-0.26678	0.08859	-0.44041	-0.09315	426	-3.011	0.003
Rater group (0 – online, 1 – lab)		-0.01412	0.08923	-0.18901	-0.16077	427	-0.158	0.874
Random effects			Variance	ICC				
Listener's ID			0.7676	0.3657				
Vocalizer's ID			0.0966	0.0676				
Residuals			1.3315					
N <sub>observations</sub> = 12702, N <sub>listeners</sub> = 430, N <sub>vocalizers</sub> = 106								

Note. 1 = single vowels, 2 = series of vowels, 3 = single word, 4 = counting, 5 = greeting, 6 = recited paragraph (Rainbow Passage)

Table 4. Estimated fixed and random effects of the model with perceived **trustworthiness** as an outcome variable

FEMALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.7917	0.05966	3.67465	3.90871	399.0	63.556	< .001
Stimulus type	2 – 1	0.8481	0.04185	0.76609	0.93015	11466.8	20.265	< .001
	3 – 1	0.1383	0.04185	0.05632	0.22037	11466.8	3.306	< .001
	4 – 1	0.9935	0.04185	0.91145	1.07551	11466.8	23.738	< .001
	5 – 1	1.0722	0.04185	0.99015	1.15421	11466.8	25.618	< .001
	6 – 1	0.9754	0.04185	0.89341	1.05747	11466.8	23.307	< .001
Vocalizer's age		-3.32e-4	0.00283	-0.00588	0.00521	98.7	-0.118	0.906
Listener's sex (0 – F, 1 – M)		-0.07119	0.10060	-0.26836	0.12598	396.9	-0.7077	0.480
Rater group (0 – online, 1 – lab)		0.00690	0.10256	-0.19412	0.20792	396.3	0.0673	0.946
Random effects			Variance	ICC				
Listener's ID			0.924	0.3466				
Vocalizer's ID			0.110	0.0592				
Residuals			1.747					
N <sub>observations</sub> = 11970, N <sub>listeners</sub> = 399, N <sub>vocalizers</sub> = 101								

MALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		3.67739	0.05999	3.55974	3.79490	397	61.299	< .001
Stimulus type	2 – 1	0.67541	0.04183	0.59342	0.75740	11156	16.145	< .001
	3 – 1	0.17490	0.04183	0.09291	0.25689	11156	4.181	< .001
	4 – 1	0.94444	0.04183	0.86245	1.02644	11156	22.577	< .001
	5 – 1	0.95936	0.04183	0.87737	1.04135	11156	22.933	< .001
	6 – 1	0.99691	0.04183	0.91492	1.07890	11156	23.831	< .001
Vocalizer's age		0.00208	0.00322	-0.00424	0.00840	105	0.646	0.520
Listener's sex (0 – F, 1 – M)		-0.07742	0.04970	-0.17482	0.01999	396	-1.558	0.120
Rater group (0 – online, 1 – lab)		0.04457	0.10159	-0.15455	0.24369	396	0.439	0.661
Random effects			Variance	ICC				
Listener's ID			0.907	0.3482				
Vocalizer's ID			0.124	0.0679				
Residuals			1.701					
N <sub>observations</sub> = 11664 N <sub>listeners</sub> = 399, N <sub>vocalizers</sub> = 106								

Note. 1 = single vowels, 2 = series of vowels, 3 = single word, 4 = counting, 5 = greeting, 6 = recited paragraph (Rainbow Passage)

Table 5. Estimated fixed and random effects of the model with perceived **femininity-masculinity** as an outcome variable

FEMALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		2.3539	0.05384	2.24838	2.4594	459.3	43.72	< .001
Stimulus type	2 – 1	-0.2972	0.02931	-0.35461	-0.2397	11705.0	-10.14	< .001
	3 – 1	-0.1036	0.02931	-0.16100	-0.0461	11705.0	-3.53	< .001
	4 – 1	-0.4181	0.02931	-0.47550	-0.3606	11705.0	-14.27	< .001
	5 – 1	-0.3915	0.02931	-0.44896	-0.3341	11705.0	-13.36	< .001
	6 – 1	-0.4240	0.02931	-0.48139	-0.3665	11705.0	-14.47	< .001
Vocalizer's age		0.0108	0.00217	0.00653	0.0151	97.6	4.96	< .001
Listener's sex (0 – F, 1 – M)		-0.1601	0.09488	-0.34608	0.0258	403.8	-1.69	0.090
Rater group (0 – online, 1 – lab)		0.1755	0.09547	-0.01162	0.3626	403.9	1.84	0.067
Random effects			Variance	ICC				
Listener's ID			0.8856	0.5028				
Vocalizer's ID			0.0664	0.0705				
Residuals			0.8757					
N <sub>observations</sub> = 12220, N <sub>listeners</sub> = 407, N <sub>vocalizers</sub> = 101								

MALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		5.63938	0.05638	5.5289	5.7499	279	100.03	< .001
Stimulus type	2 – 1	0.20407	0.02837	0.1485	0.2597	11459	7.19	< .001
	3 – 1	0.10107	0.02838	0.0455	0.1567	11460	3.56	< .001
	4 – 1	0.30717	0.02837	0.2516	0.3628	11459	10.83	< .001
	5 – 1	0.29966	0.02837	0.2441	0.3553	11459	10.56	< .001
	6 – 1	0.22459	0.02837	0.1690	0.2802	11459	7.92	< .001
Vocalizer's age		0.00719	0.00369	-5.85e-5	0.0144	104	1.94	0.055
Listener's sex (0 – F, 1 – M)		-0.10305	0.07717	-0.2543	0.0482	404	-1.34	0.182
Rater group (0 – online, 1 – lab)		-0.08457	0.07768	-0.2368	0.0677	405	-1.09	0.277
Random effects			Variance	ICC				
Listener's ID			0.574	0.417				
Vocalizer's ID			0.180	0.183				
Residuals			0.802					
N <sub>observations</sub> = 11977, N <sub>listeners</sub> = 407, N <sub>vocalizers</sub> = 106								

Note. Higher values indicate higher masculinity. 1 = single vowels, 2 = series of vowels, 3 = single word, 4 = counting, 5 = greeting, 6 = recited paragraph (Rainbow Passage)

Table 6. Estimated fixed and random effects of the model with perceived **health** as an outcome variable

FEMALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		5.0332	0.05752	4.9205	5.14605	312.5	87.357	< .001
Stimulus type	2 – 1	0.7026	0.03937	0.6254	0.77976	12691.4	17.845	< .001
	3 – 1	0.2096	0.03937	0.1325	0.28679	12691.4	5.324	< .001
	4 – 1	0.7203	0.03937	0.6431	0.79745	12691.4	18.295	< .001
	5 – 1	0.8305	0.03937	0.7533	0.90765	12691.4	21.094	< .001
	6 – 1	0.6582	0.03937	0.5810	0.73532	12691.4	16.717	< .001
Vocalizer's age		-0.0155	0.00321	-0.0218	-0.00924	98.6	-4.842	< .001
Listener's sex (0 – F, 1 – M)		-0.1010	0.08721	-0.2720	0.06989	437.4	-1.16	0.247
Rater group (0 – online, 1 – lab)		-0.1667	0.08885	-0.3408	0.00743	437.1	-1.88	0.061
Random effects			Variance	ICC				
Listener's ID			0.746	0.3036				
Vocalizer's ID			0.149	0.0802				
Residuals			1.712					
N <sub>observations</sub> = 13240, N <sub>listeners</sub> = 441, N <sub>vocalizers</sub> = 101								

MALE VOCALIZERS

Fixed effects		Estimate	SE	95% CI		Df	t	p
				Lower	Upper			
Intercept		4.80384	0.06304	4.6803	4.9274	325	76.209	< .001
Stimulus type	2 – 1	0.41164	0.03963	0.3340	0.4893	12355	10.387	< .001
	3 – 1	0.34336	0.03963	0.2657	0.4210	12355	8.664	< .001
	4 – 1	0.67963	0.03963	0.6020	0.7573	12355	17.149	< .001
	5 – 1	0.74280	0.03963	0.6651	0.8205	12355	18.743	< .001
	6 – 1	0.58674	0.03963	0.5091	0.6644	12355	14.805	< .001
Vocalizer's age		-0.00788	0.00388	-0.0155	-2.87e-4	103	-2.034	0.045
Listener's sex (0 – F, 1 – M)		-0.05583	0.09500	-0.2420	0.1304	438	-0.1588	0.557
Rater group (0 – online, 1 – lab)		-0.21749	0.09678	-0.4072	-0.0278	438	-2.247	0.025
Random effects			Variance	ICC				
Listener's ID			0.895	0.347				
Vocalizer's ID			0.191	0.102				
Residuals			1.686					
N <sub>observations</sub> = 12908 N <sub>listeners</sub> = 441, N <sub>vocalizers</sub> = 106								

Note. 1 = single vowels, 2 = series of vowels, 3 = single word, 4 = counting, 5 = greeting, 6 = recited paragraph (Rainbow Passage)

