



BIROn - Birkbeck Institutional Research Online

Tierney, Adam and Patel, A. and Jasmin, Kyle and Breen, M. (2021) Individual differences in perception of the speech-to-song illusion are linked to musical aptitude but not musical training. *Journal of Experimental Psychology: Human Perception and Performance* 47 (12), pp. 1681-1697. ISSN 0096-1523.

Downloaded from: <https://eprints.bbk.ac.uk/id/eprint/47395/>

Usage Guidelines:

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html>
contact lib-eprints@bbk.ac.uk.

or alternatively

1 **Title:** Individual Differences in Perception of the Speech-to-Song Illusion are Linked to Musical
2 Aptitude but not Musical Training

3

4 **Authors:** Adam Tierney^{1†}, Aniruddh D. Patel^{2,3}, Kyle Jasmin¹, Mara Breen⁴

5

6 **Affiliations**

7 ¹Department of Psychological Sciences, Birkbeck, University of London, London, UK

8 ²Department of Psychology, Tufts University, Medford, MA

9 ³Azrieli Program in Brain, Mind, & Consciousness, Canadian Institute for Advanced Research (CIFAR),

10 Toronto

11 ⁴Department of Psychology, Mount Holyoke College, South Hadley, MA

12

13 †Corresponding author

14 Adam Tierney, Ph.D.

15 Birkbeck, University of London

16 Malet Street, London, WC1E 7HX

17 United Kingdom

18 Email: a.tierney@bbk.ac.uk

19 Phone: +44-020-7631-6368

20

21 Word count: 10346

22

23 **Abstract**

24 In the speech-to-song illusion certain spoken phrases are perceived as sung after repetition. One
25 possible explanation for this increase in musicality is that, as phrases are repeated, lexical activation
26 dies off, enabling listeners to focus on the melodic and rhythmic characteristics of stimuli and assess
27 them for the presence of musical structure. Here we tested the idea that perception of the illusion
28 requires implicit assessment of melodic and rhythmic structure by presenting individuals with
29 phrases that tend to be perceived as song when repeated, as well as phrases that tend to continue
30 to be perceived as speech when repeated, measuring the strength of the illusion as the rating
31 difference between these two stimulus categories after repetition. Illusion strength varied widely
32 and stably between listeners, with large individual differences and high split-half reliability,
33 suggesting that not all listeners are equally able to detect musical structure in speech. Although
34 variability in illusion strength was unrelated to degree of musical training, participants who
35 perceived the illusion more strongly were proficient in several musical skills, including beat
36 perception, tonality perception, and selective attention to pitch. These findings support models of
37 the speech-to-song illusion in which experience of the illusion is based on detection of musical
38 characteristics latent in spoken phrases.

39 **Public Significance Statement**

40 In the speech-to-song illusion certain spoken phrases begin, after repeated presentation, to
41 sound as if they were sung. People vary widely in how vividly they experience this illusion,
42 but the sources of this variability are still poorly understood. We find that people show
43 stable individual differences in the perception of the illusion and that this variability is
44 largely unrelated to factors such as degree of musical training and language background.
45 Instead, we find that individuals who vividly perceive the illusion have high proficiency in
46 music perception skills related to timing and pitch perception. These findings contribute to
47 understanding the relations between language and music processing in the human mind.

49 1. Introduction

50 1.1. Domain-level classification

51 One of the most common perceptual tasks accomplished by individuals is categorization, in
52 which complex input varying continuously in many dimensions is binned into one of several discrete
53 categories. In speech perception, for example, word recognition involves accurate speech sound
54 categorization, while in music perception the recognition of meter involves perceiving strong vs.
55 weak beats (Patel, 2008). For decades researchers have investigated the cues that listeners use to
56 distinguish between categories in speech (Toscano & McMurray, 2010) and music (Prince, 2014).
57 However, comparatively little is known about the processes by which listeners decide to which
58 domain a stimulus belongs—to speech or to music, for example. Prior research on domain-level
59 categorization (i.e. speech, music, or environmental sounds) has largely focused on classification of
60 sounds that differ in the physical source used to produce the sound (for example, tones produced by
61 musical instruments, vowels produced by the human vocal tract, and environmental sounds
62 produced by a variety of sources). This research has found, for example, that listeners can use a
63 variety of acoustic cues to distinguish between sounds from these three categories within tens of
64 milliseconds, as reflected in both behavioral responses and magnetoencephalographic patterns
65 (Bigand, Delbé, Gérard, & Tillmann, 2011; Ogg, Slevc, & Idsardi, 2017; Ogg, Carlson, & Slevc, 2020).
66 However, sounds from different domains do not always differ in their physical source, as in the case
67 of speech and song, both of which are produced by the human vocal tract. Another way to
68 investigate domain-level classification is to make use of ambiguous spoken stimuli which can be
69 heard as spoken or sung, depending on context (Deutsch, Honthorn, & Lapidis, 2011). For these
70 ambiguous cross-domain stimuli, domain-level categorization can be a much more protracted
71 process lasting tens of seconds rather than tens of milliseconds, suggesting that listeners'
72 categorization may be influenced by acoustic patterns on a longer time scale.

73 1.2. Sound-to-music illusions

74 Music often has acoustic patterns which distinguish it from speech or other sounds (Ding et
75 al., 2017), and also elicits distinct neural responses compared to other sounds (Norman-Haignere et
76 al. 2015; Zuk et al., 2020). Nevertheless, studies over the past decade have demonstrated that
77 listeners sometimes report perceiving non-musical sounds (i.e., sounds not originally intended to be
78 heard as music) as sounding like music. An important phenomenon in this regard is the speech-to-
79 song illusion. In this illusion, first published in recorded form in 2003 (Deutsch 2003) and in the
80 academic literature ten years ago (Deutsch et al., 2011), a phrase which is perceived as spoken when
81 presented once tends to be heard as sung after repetition. (The key phrase occurs in an even earlier
82 recording from 1995 (Deutsch 1995), but it was not isolated and presented as an illusion until 2003.)
83 Building on this work, Simchy-Gross & Margulis (2018) and Rowland, Kasdan, & Poeppel (2019)
84 demonstrated that repeated environmental sounds are rated as more musical after repetition, and
85 Tierney, Patel, & Breen (2018a) found that pitch contours extracted from speech and resynthesized
86 as complex tones increased in perceived musicality with repetition. However, not all non-musical
87 stimuli sound musical after repetition, with some undergoing dramatic perceptual changes and
88 others continuing to be heard as non-musical (Tierney, Patel, & Breen, 2018b). This raises the key
89 question of why listeners perceive certain non-musical stimuli as music when they are played
90 repeatedly.

91 *1.3. Proposed mechanisms of the speech-to-song illusion*

92 To explain the speech-to-song illusion (henceforth, “song illusion”), Deutsch et al. (2011)
93 suggested that pitch salience is inhibited by default during speech perception, given that pitch plays
94 a secondary role in communicating lexical information relative to other acoustic dimensions (in non-
95 tonal languages). Based on this hypothesis, they predicted that an area anterolateral to primary
96 auditory cortex which had previously been linked to pitch salience (Patterson, Uppenkamp,
97 Johnsrude, & Griffiths, 2002; Penagos, Melcher, & Oxenham, 2004; Schneider et al., 2005) would
98 show greater activation for stimuli which perceptually transform into song when repeated. This
99 prediction was confirmed in an fMRI study by Tierney, Dick, Deutsch, & Sereno (2013), suggesting

100 that perception of the illusion is indeed linked to increased pitch salience. Deutsch et al. (2011) also
101 hypothesized that perception of speech as song involves perceptual distortion of pitch, such that the
102 perceived pitches of syllables are warped to fit a musical scale template. This hypothesis was
103 supported by Vanden Bosch der Nederlanden, Hannon, & Snyder (2015a), who found that when
104 participants perceived the song illusion in repeated spoken phrases, pitch changes which departed
105 from the perceived melodic template were easier to detect than pitch changes which moved
106 towards it.

107 The explanation for the song illusion advanced by Deutsch et al. (2011), however, does not
108 specify why pitch perception should be disinhibited by stimulus repetition. A potential explanation
109 was advanced by Castro, Mendoza, Tampke, & Vitevich (2018), who proposed that Node Structure
110 Theory (MacKay, Wulf, Yin, & Abrams, 1993) could account for the song illusion. According to this
111 theory, repeated activation of lexical nodes temporarily reduces their ability to become activated.
112 Castro et al. (2018) suggest that while deactivation of lexical nodes diminishes the perception of
113 speech, continued activation of syllable nodes leads to the emergence of a song percept. Supporting
114 a role for lexical deactivation in the illusion, they find that words from denser phonological
115 neighborhoods and words from an unfamiliar language are rated as more song-like after repetition.
116 Similarly, Vitevitch, Ng, Hatley, & Castro (2020) report that words with a high phonological clustering
117 coefficient (which measures the tendency for phonological neighbors of a word to also be neighbors
118 of one another) are perceived as more song-like after repetition.

119 The Node Structure Theory explanation for the song illusion, however, struggles to explain
120 why deactivation of lexical nodes along with continued activation of syllable nodes should lead to
121 perception of song, rather than speech leeched of semantic content (as in semantic satiation; Smith,
122 1984). The explanation advanced in Castro et al. (2018) is that since syllables are the unit of rhythmic
123 structure in speech, and rhythm is also an important aspect of music, then syllable node activation
124 without lexical node activation should lead to a song-like percept. According to this explanation, all
125 spoken stimuli should transform into song when repeated, as lexical nodes are satiated and syllable

126 nodes continue to be active. Moreover, according to this explanation, the extent to which stimuli do
127 or do not give rise to the illusion should be primarily tied to phonological rather than acoustic or
128 musical characteristics. Finally, this theory cannot explain increases in musicality with repetition of
129 non-verbal stimuli.

130 These predictions, however, are not borne out by the literature. First, there is ample
131 evidence that not all spoken phrases transform into song when repeated. Tierney et al. (2013), for
132 example, assembled a corpus of "illusion" phrases, which listeners report transform into song when
133 repeated, and "control" phrases, which continue to be perceived as speech when repeated. (These
134 two types of phrases were matched for speakers, syllable rate, and duration.) The relative lack of
135 transformation in the control phrases was replicated in Graber, Simchy-Gross, & Margulis (2017),
136 Tierney et al. (2018a), and Tierney et al. (2018b). Second, several studies have demonstrated that a
137 range of acoustic and musical stimulus characteristics modulate the strength of the illusion. Tierney
138 et al. (2013) showed that the illusion stimuli featured flatter pitch contours than the control stimuli.
139 Falk, Rathcke, & Dalla Bella (2014) manipulated the stability of pitch contours within syllables and
140 found that song percepts were more common for the more stable contours. Tierney et al. (2018b)
141 found that higher song ratings for song illusion stimuli were linked to more isochronous rhythmic
142 structure and a greater tendency for syllable pitches to follow musical melodic statistics. Rathcke,
143 Falk, & Dalla Bella (in press) found that recordings with greater sonority led to earlier, stronger, and
144 more frequent transformations to song. Third, Tierney et al. (2018a) found that eliminating linguistic
145 content from the illusion/control stimulus sets (Tierney et al., 2013) preserves the relative
146 differences across these stimuli in the extent to which they give rise to the sound-to-music illusion,
147 suggesting that the primary factors driving cross-stimulus differences in illusion strength cannot be
148 phonological. Fourth, several studies have demonstrated that environmental sounds and non-verbal
149 complex tone sequences increase in musicality after repetition (Simchy-Gross & Margulis, 2018;
150 Tierney et al., 2018a; Rowland et al., 2019), a finding that cannot be explained by lexical processes.

151 Thus Node Structure Theory alone is not a sufficient explanation of the song illusion. Given
152 this observation, what theoretical framework could account for both the influence of phonological
153 neighborhood and the influence of acoustic and musical characteristics on the perception of song in
154 speech? One possibility is a hybrid of the accounts advanced in Deutsch et al. (2011) and Castro et al.
155 (2018) in which deactivation of lexical nodes is *necessary but not sufficient* for perception of the
156 illusion. According to this account, deactivation of lexical nodes frees up attention which can then be
157 focused on acoustic characteristics of the stimuli, including pitch and timing patterns. The degree to
158 which a stimulus transforms into song after lexical node deactivation, therefore, will depend on the
159 melodic and rhythmic characteristics of the stimulus. One implication of this hybrid model is that
160 there may exist stable individual differences across participants in the extent to which the vividness
161 of the illusion varies across stimuli, and these individual differences may relate to musical training
162 and musical aptitude. In two experiments we tested this implication by presenting participants with
163 repeated phrases drawn from one of two stimulus sets, a stimulus set which listeners consistently
164 report strongly transforms into song after repetition and a stimulus set which transforms to a much
165 lesser degree ("Illusion" and "Control" stimuli, as documented in Tierney et al., 2013; Graber et al.,
166 2017; Tierney et al., 2018b). We predicted that participants would reliably vary in the difference in
167 perceived musicality between Illusion and Control stimuli after repetition, and that this variability
168 would relate to musical training (Experiment 1) and musical aptitude (Experiment 2).

169 2. Experiment 1

170 2.1. Introduction

171 Musical training is one potential factor which could underly individual differences in the
172 song illusion: if perception of the illusion requires detection of musical characteristics latent in non-
173 musical stimuli, then musicians may be better able to detect these characteristics. Yet prior research
174 has reported that musicians are somewhat less susceptible to auditory illusions in which stimulus
175 details are lost in a top-down percept or gestalt (Craig, 1979; Davidson, Power, & Michie, 1987;
176 Brennan & Stevens, 2002; Pressnitzer, Graves, Chambers, de Gardelle, & Egré, 2018), perhaps

177 because musicians are better able to attend to low-level acoustic characteristics of stimuli. The
178 existing evidence regarding musical training and the speech-to-song illusion is inconsistent. Falk,
179 Rathcke, & Dalla Bella (2014) found that participants with more years of musical training were no
180 more likely to perceive the illusion, and in fact that their perception of the illusion was delayed
181 relative to musically untrained participants. However, Vanden Bosch der Nederlanden et al. (2015b)
182 found that musicians reported hearing all stimulus presentations as more musical than non-
183 musicians, regardless of stimulus repetition or transposition. Given that Vanden Bosch der
184 Nederlanden et al. (2015b) used only the original example from Deutsch et al. (2011) as a stimulus,
185 this finding could indicate either that musicians are better able to detect the latent musical
186 characteristics of this stimulus or that musicians are biased to perceive musicality regardless of
187 stimulus characteristics.

188 In addition to our primary research question regarding musical training, we also investigated
189 several other possible predictors of individual differences in perception of the illusion. For example,
190 another possible factor underlying variability in the song illusion is age. As described above, Castro
191 et al. (2018) speculated that one important precondition for perception of the illusion is deactivation
192 of lexical nodes due to satiation. A similar explanation has also been advanced for the verbal
193 transformation effect, in which repetition of a word eventually leads to an unstable percept that
194 swaps between different words with different meanings (Warren & Gregory, 1958): MacKay et al.
195 (1993) suggested that this illusion occurs when lexical nodes are satiated but phonological nodes
196 remain activated. The verbal transformation effect is weaker in older participants (Warren &
197 Warren, 1966; Pilotti & Khurshid, 2004; Pilotti, Simcox, Baldy, & Schauss, 2011), suggesting that
198 satiation of lexical nodes due to repetition may be less rapid or less robust in older age. If so, then
199 older participants may demonstrate less robust perception of the speech-to-song illusion as well.
200 However, Mullin, Norkey, Kodwani, Vitevitch, & Castro (2021) reported no relationship between age
201 and the strength of the speech-to-song effect, as perceived in the stimulus from Deutsch et al.
202 (2011).

203 A third possible factor underlying individual differences in the song illusion is gender¹. There
204 is some limited prior evidence for sex/gender differences in the perception of auditory illusions. For
205 example, Irwin, Whalen, & Fowler (2006) reported that female participants perceived the McGurk
206 illusion more strongly (i.e. greater influence of visual stimuli on auditory perception), but only when
207 the visual stimuli were relatively brief. Women were also found to be more likely to report the
208 illusory percept (i.e. "yanny") after hearing the Yanny/Laurel stimuli (Pressnitzer et al. 2018), while
209 Gwilliams & Wallisch (2020) similarly found more reports of "yanny" among female participants.
210 Given this prior evidence for sex/gender differences in perception of auditory illusions, as well as
211 evidence for sex differences in subcortical auditory processing (Krizman, Bonacina, & Kraus, 2019),
212 we compared perception of the illusion between participants identifying as male versus as female.

213 Finally, there is some prior evidence that language background can affect perception of the
214 speech-to-song illusion. Jaisin, Suphanchaimat, Candia, & Warren (2016) reported that native
215 speakers of tonal languages did not perceive the speech-to-song illusion, although given the small
216 sample sizes used (10 participants in each group) this finding needs replication. Castro et al. (2018)
217 reported that repeated speech in an unfamiliar language was perceived as more musical than
218 repeated speech in a familiar language. Rathcke et al. (2021) reported that bilingual participants
219 reported more speech-to-song transformations in their second language. Finally, Margulis, Simchy-
220 Gross, & Black (2015) reported that repeated speech in an unfamiliar language was perceived as
221 particularly musical for languages that would be difficult for a participant to pronounce. Given these
222 somewhat mixed prior findings, we investigated effects of language experience by determining both
223 whether participants spoke English as their dominant language as well as whether participants were
224 monolingual (only able to speak one language) or bilingual (able to speak at least two languages).

225

226 2.2. Methods

¹ Sex and gender are highly correlated but dissociable. In the current study we focused on gender, and will use that term throughout, but we describe prior results using the language that was used by the study authors.

227 2.2.1. Participants

228 One hundred seventy-one (171) participants were initially recruited from a mixture of
229 undergraduate psychology students and Mechanical Turk. Twenty-three (23) participants were
230 removed from the dataset because they failed catch trials, not reporting an increase in song
231 perception after stimuli switched from actual speech to actual song (see below for details), or
232 because they failed to produce at least one response after the first through fourth repetitions and
233 after the fifth through eighth repetitions during the speech-to-song illusion rating experiment.
234 Responses from 148 remaining participants were analyzed. A series of power analyses were
235 conducted in GPower to determine the smallest effect sizes which could be detected at a power of
236 0.8, given this sample size and an alpha of 0.05. This sample size resulted in a power of 0.8 to detect
237 a correlation of $r = 0.23$. For the analysis of the speech-to-song ratings, given a 2 (stimulus set) X 8
238 (repetition) RMANOVA design, the sample size resulted in a power of 0.80 to detect main effects and
239 interactions of $F = 0.078$. Ninety-three (93) participants identified as male, while 55 identified as
240 female. (Participants were also given the additional options "Other" and "Prefer not to say", but no
241 participants selected these options.) Participants reported a mean age of 29.89 (SD 10.14) years old
242 (range 19-61). 120 participants reported that their dominant language was English, while 28
243 participants reported an alternate dominant language. These other dominant languages were Arabic
244 ($n = 1$), Bulgarian ($n = 1$), French ($n = 1$), German ($n = 1$), Greek ($n = 2$), Italian ($n = 6$), Polish ($n = 3$),
245 Portuguese ($n = 3$), Romanian ($n = 2$), Russian ($n = 2$), Slovak ($n = 2$), Spanish ($n = 3$), and Swedish ($n =$
246 1). Sixty-seven (67) participants reported being bilingual (i.e. able to speak at least two languages),
247 while 81 participants reported being monolingual. Participants reported a mean of 3.19 (SD 4.24)
248 years of musical training (range 0-17). Study procedures were approved by the ethics board of the
249 Department of Psychological Sciences at Birkbeck College.

250 2.2.2. Speech-to-song illusion experiment

251 Participants were presented with 48 spoken phrases drawn from the stimulus set first
252 reported in Tierney et al. (2013). Phrases were produced by a mix of American and British English

253 speakers, but illusion and control datasets were perfectly matched for speakers. Prior studies
254 (Tierney et al., 2013; 2018a; 2018b; Graber et al., 2017) have indicated that 24 of these stimuli
255 tended to be heard as song after repetition (henceforth “illusion” stimuli), while 24 tended to be
256 heard as speech after repetition (henceforth “control” stimuli). The mean length of stimuli was 6.6
257 (SD 1.5) syllables and 1.38 (SD 0.41) seconds. Stimuli were spoken by three different male talkers,
258 represented in equal portions among illusion and control stimuli. These two stimulus categories did
259 not differ significantly in syllable rate (illusion stimuli 5.05 syllables/second, control stimuli 4.99
260 syllables/second, $t(46) = 0.16$, $p = 0.877$) or duration (illusion stimuli 1.32 seconds, control stimuli
261 1.44 seconds, $t(46) = -1.05$, $p = 0.301$).

262 In each trial, eight repetitions of each phrase were presented. After each repetition,
263 participants pressed buttons labelled 1 through 10 to indicate whether the phrase sounded like
264 speech or song, where 1 indicated completely speech-like, while 10 indicated completely song-like.
265 Participants were given 2 seconds to respond to each phrase, after which time the program
266 automatically went on to the next repetition. Any missing rating (due to a non-response) was
267 replaced by the mean of the previous and following ratings. Missing first/last repetitions were
268 replaced by the second/seventh repetitions.

269 In addition, four catch trials were presented. In each catch trial, a spoken phrase was
270 presented during the first four repetitions, while a matched sung phrase (with the same words,
271 roughly same rate, and similar pitch contour) was presented during the second four repetitions.
272 Catch trials were produced by the first author. Analysis was limited to data from those participants
273 whose average rating of the last four repetitions exceeded that of the first four (by any amount) on
274 these catch trials. As noted above, the exclusion criteria resulted in the removal of 23 out of 171
275 participants from the dataset, leaving 148 for analysis. The experiment lasted around 20 minutes.

276 To establish whether an overall speech-to-song illusion effect was present, a repeated
277 measures ANOVA was conducted with two within-subjects factors, repetition (eight levels) and
278 stimulus set (illusion versus control). Next, metrics summarizing illusion perception were compared

279 to demographic characteristics. With eight ratings for each illusion and control stimulus, this is a
280 high-dimensional dataset, potentially exacerbating the multiple comparisons problem. To reduce the
281 dimensionality of the data, therefore, we summarized each participant's responses in two ways for
282 the purpose of investigating individual differences in perception of the illusion. First, the difference
283 between the ratings of illusion and control stimuli after the eighth repetition was calculated and will
284 be referred to as *illusion strength*. (In other words, for each participant we averaged their final rating
285 for all illusion stimuli and all control stimuli, and then took the difference between these two mean
286 values to compute illusion strength.) Second, the average rating across all repetitions and all stimuli
287 was calculated and will be referred to as *musical prior*.

288 Any variables not normally distributed (Jarque & Bera, 1980) were transformed prior to
289 statistical analysis. A square root transformation was used for musical prior, and age and years of
290 musical training were converted to ranks. Pearson's correlations were used to assess the strength of
291 the relationship between musical prior and illusion strength versus age and years of musical training.
292 Musical prior and illusion strength were also compared between dominant English speakers and
293 non-dominant English speakers, as well as between bilinguals and monolinguals and between male
294 and female participants, using unpaired t tests. Analyses predicting illusion strength and musical
295 prior were corrected separately for multiple comparisons using false discovery rate (Benjamini &
296 Hochberg, 1995).

297 2.2.3. Available materials

298 Data from both Study 1 and Study 2 and test materials from Study 2 are available at
299 <https://osf.io/hegxs/>. All stimuli from the speech-to-song illusion corpus are available at
300 <https://osf.io/t4pqj/>.

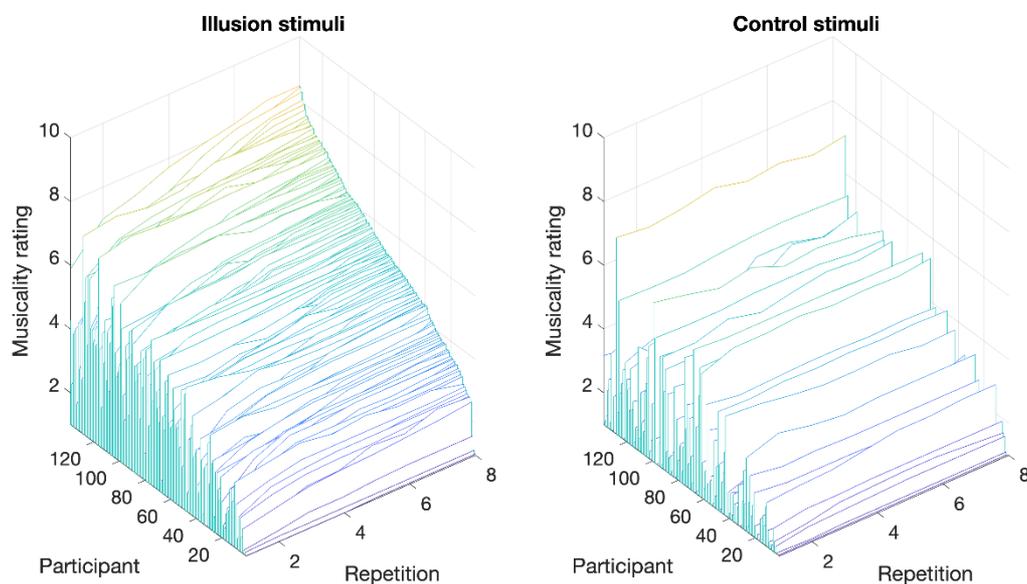
301 2.3. Results

302 Overall, across participants the illusion was reliably experienced, as seen in the strong
303 interaction between repetition and stimulus set ($F(7,1029) = 216.3, p < 0.001$). There was, in
304 addition, a main effect of repetition ($F(1,1029) = 146.8, p < 0.001$) and of stimulus set ($F(1,137) =$

305 426.4, $p < 0.001$). Overall, participants rated that the illusion stimuli sounded much more musical
306 than the control stimuli after repetition; in fact, after averaging across participants, there was no
307 overlap in the ratings between the two stimulus sets (in other words, no control stimulus was rated
308 higher than any illusion stimulus). However, there were very large individual differences across
309 participants in the strength of the illusion, ranging from participants who reported no difference
310 between illusion and control stimuli whatsoever, to participants who reported that the illusion
311 stimuli sounded exactly like song when repeated while the control stimuli continued to sound
312 exactly like speech. These individual differences are displayed in Figure 1, which plots average
313 ratings of illusion and control stimuli across all eight repetitions from all 148 participants, sorted so
314 that participants with larger ratings of illusion stimuli after the eighth repetition are closer to the
315 top. (A color plot version of the same figure is presented in the supplementary information, see
316 Figure S1. Additional plots showing individual differences for individual illusion and control stimuli
317 are also available in supplementary information, see Figures S2 and S3. Also see the supplementary
318 information for sound examples of the illusion stimuli and control stimuli with the least and most
319 variability in rating across participants, see Sound Examples S1 through S4.) In addition, Figure 2
320 displays the difference between the average rating of the last repetition and the average ratings of
321 the first repetition in illusion and control stimuli, as well as the difference in average rating between
322 illusion and control stimuli after the first repetition and the last repetition. The difference between
323 the average last and first repetition rating for illusion stimuli ranged from -0.17 to 8 (mean 1.91, std
324 1.45), while this difference for control stimuli ranged from -1.17 to 4.92 (mean 0.38, std 0.90). The
325 difference in average rating between illusion and control stimuli after the first repetition ranged
326 from -0.83 to 4.75 (mean 1.33, std 1.04), while this difference after the last repetition ranged from -
327 0.25 to 7.75 (mean 2.85, std 1.61).

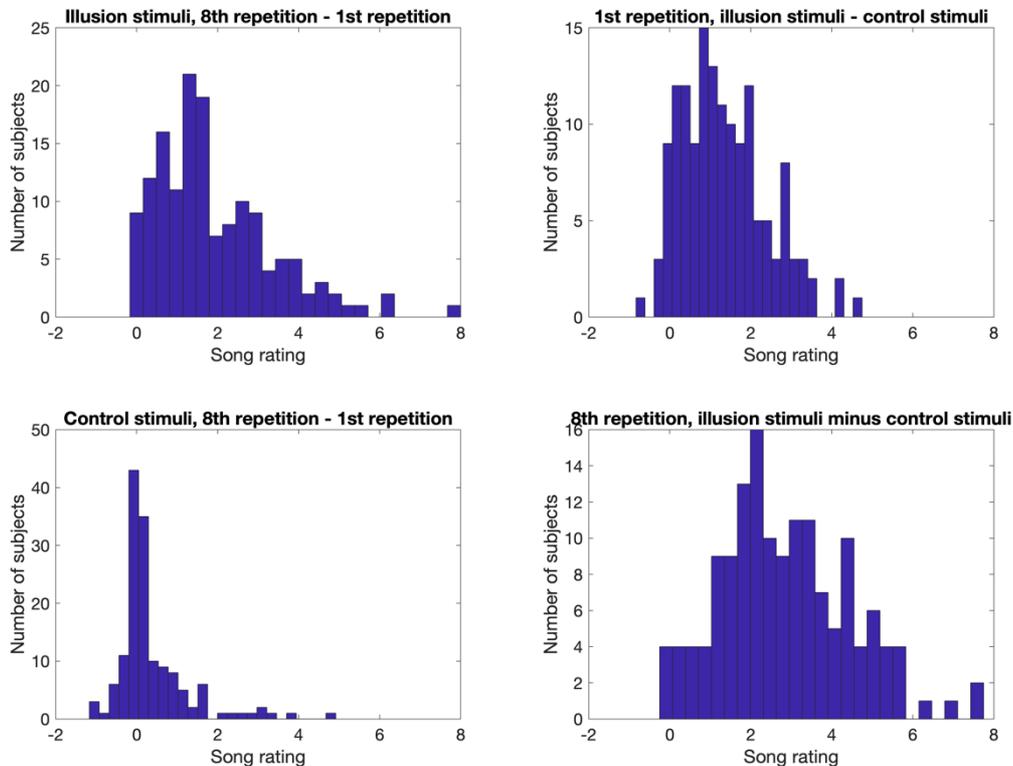
328 To assess the reliability of this speech-to-song testing battery as a measure of individual
329 differences in musicality perception, a Monte Carlo procedure was used. Across 100 iterations, the
330 illusion and control stimuli were randomly divided into two halves. The average rating after the

331 eighth repetition was then calculated separately for the illusion and control stimuli. The difference
332 between ratings of the illusion and control stimuli was then calculated separately for each half of the
333 trials (the *illusion strength* measure described above). The estimates from the two halves were
334 correlated using a Pearson's correlation. The resulting correlation was corrected using the
335 Spearman-Brown correction, to account for the fact that reliability would be expected to be slightly
336 higher for a test with twice the number of trials. Finally, the median split half reliability across the
337 100 iterations was calculated. This procedure revealed a reliability of 0.89, suggesting that this
338 testing battery provides a highly reliable measure of individual differences in the perception of
339 musicality in spoken stimuli. A similar procedure revealed that reliability for the *musical prior*
340 measure was also high, reaching 0.96. Interestingly, however, illusion strength and musical prior only
341 weakly correlated ($r(146) = 0.19, p = 0.023$), suggesting that these two measures, although highly
342 reliable, index dissociable aspects of the speech-to-song illusion.
343



344
345 **Figure 1:** Waterfall plot displaying ratings of illusion and control stimuli across all 148 participants,
346 sorted by the final rating of illusion stimuli. Each row shows the averaged ratings of a participant
347 across all 24 stimuli in that category. The illusion stimuli matrix is sorted from highest to lowest final

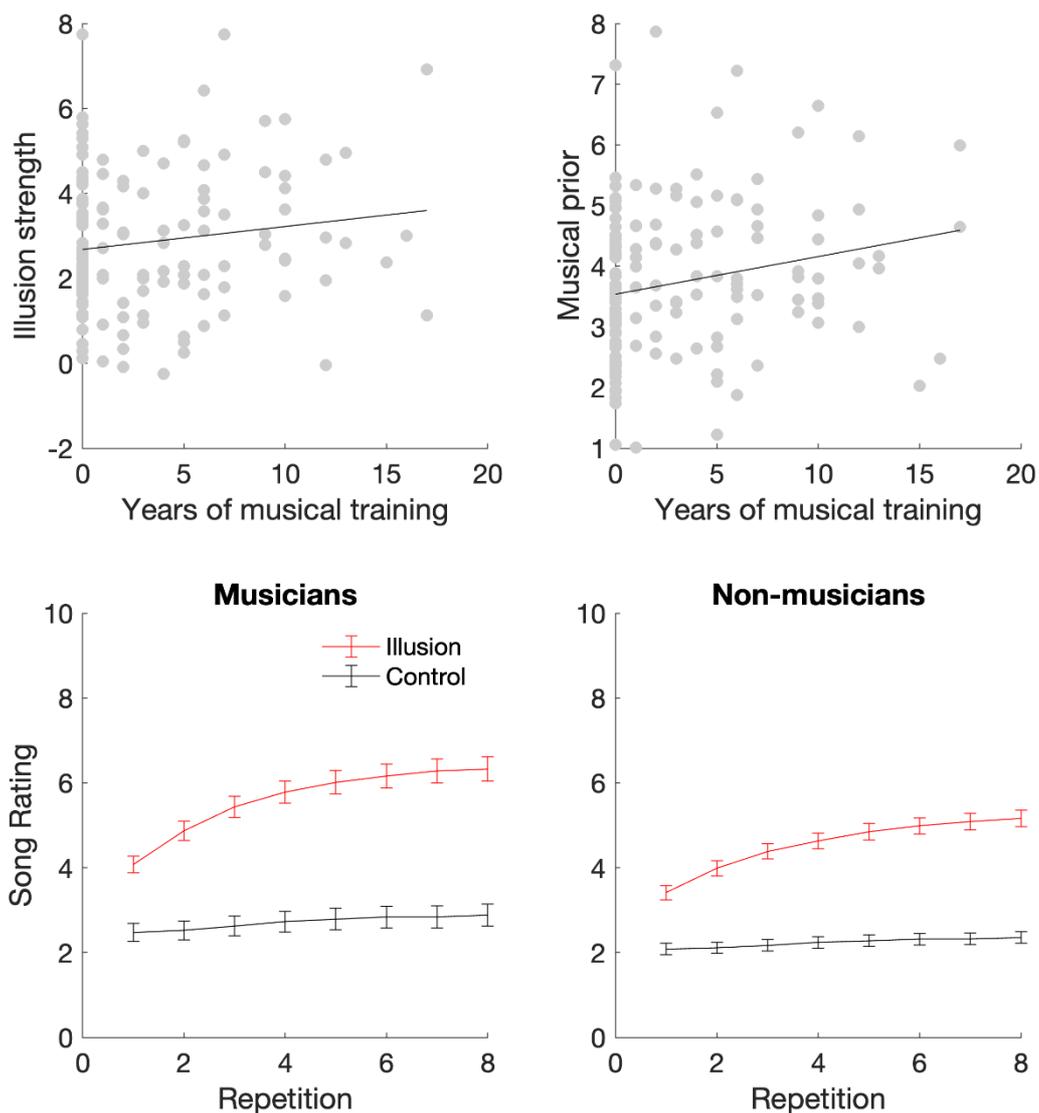
348 rating, and the control stimuli matrix shows each participant's corresponding average rating profile
 349 for control stimuli.
 350



351
 352 **Figure 2:** Top-left: histogram displaying the difference between the average rating after the last
 353 repetition minus the average rating after the first repetition for illusion stimuli. Bottom-left:
 354 histogram displaying similar data for control stimuli. Top-right: histogram displaying the difference
 355 between the average initial rating for illusion stimuli minus the average initial rating for control
 356 stimuli. Bottom-right: histogram displaying similar data for final ratings of illusion and control
 357 stimuli.

358
 359 Years of musical training did not correlate with illusion strength ($r(146) = 0.10$, $p(\text{corrected}) =$
 360 0.531). However, years of musical training did correlate with musical prior ($r(146) = 0.25$,
 361 $p(\text{corrected}) = 0.013$), such that participants with a greater degree of musical training rated *all*
 362 stimuli (both illusion and control) as more song-like across all repetitions. Figure 3 displays

363 scatterplots relating degree of musical training to both illusion strength and musical prior, as well as
364 average responses to each repetition of illusion and control stimuli in non-musicians and musicians.
365 (For display purposes musicians were defined as having at least six years of training, following Zhang
366 et al. (2020), while non-musicians were defined as having no years of training.)
367



368
369 **Figure 3:** Top-left: scatterplot displaying years of musical training versus illusion strength. Illusion
370 strength was calculated as the rating difference between the illusion and control stimulus sets after
371 the eighth repetition. Top-right: scatterplot displaying years of musical training versus musical prior.

372 Musical prior was calculated as the grand average of ratings across all experimental trials over both
373 stimulus sets. Bottom-left: average responses to illusion and control stimuli in musicians (n = 35).
374 Bottom-right: average responses to illusion and control stimuli in non-musicians (n = 64). Error bars
375 indicate standard error of the mean.

376

377 There was no correlation between age and illusion strength ($r(146) = -0.08$, $p(\text{corrected}) =$
378 0.589). There was similarly no correlation between age and musical prior ($r(146) = -0.02$,
379 $p(\text{corrected}) = 0.823$). There was no difference in illusion strength ($t(146) = 2.19$, $p(\text{corrected}) =$
380 0.152) between participants who identified as female (mean = 3.07, std = 1.54) versus participants
381 who identified as male (mean = 2.48, std = 1.68). Similarly, there was no difference in musical prior
382 ($t(146) = 1.25$, $p(\text{corrected}) = 0.533$) between participants who identified as female (mean = 3.83, std
383 = 1.25) compared to participants who identified as male (mean = 3.59, std = 1.28). There was no
384 difference in illusion strength ($t(146) = 0.09$, $p(\text{corrected}) = 0.927$) between bilinguals (mean = 2.87,
385 SD = 1.69) versus monolinguals (mean = 2.84, SD = 1.56). Similarly, there was no difference in
386 musical prior ($t(146) = 0.235$, $p(\text{corrected}) = 0.823$) between bilinguals (mean = 3.77, SD = 1.30)
387 versus monolinguals (mean = 3.71, SD = 1.24). There was no difference in illusion strength ($t(146) = -$
388 0.57 , $p(\text{corrected}) = 0.710$) between participants whose dominant language was English (mean =
389 2.82 , std = 1.68) compared to participants whose dominant language was not English (mean = 3.01,
390 std = 1.31). Similarly, there was no difference in musical prior ($t(146) = -0.22$, $p(\text{corrected}) = 0.823$)
391 between participants whose dominant language was English (mean = 3.72, std = 1.25) compared to
392 participants whose dominant language was not English (mean = 3.78, std = 1.32).

393

394 *2.4. Discussion*

395 We find large individual differences in the magnitude of the song illusion across a diverse
396 sample of participants drawn from the general population. There was widespread agreement across
397 participants that the musicality of the control stimuli did not change with repetition. However, there

398 was a high degree of variability across participants in the extent to which the illusion and control
399 stimuli differed in musicality, and this variability increased as the stimuli were repeated: some
400 participants reported large increases in musicality in the illusion stimuli with repetition, whereas
401 others reported more modest increases.

402 Importantly, our split-half reliability calculations showed a reliability of 0.88 for our measure
403 of illusion strength, which was calculated as the difference in musicality between the illusion and
404 control stimuli after repetition. This confirms a key prediction of the hybrid model of the speech-to-
405 song illusion (see Section 1.2), which suggests that deactivation of lexical nodes after repetition frees
406 up cognitive resources, enabling participants to extract the melodic and rhythmic characteristics of
407 the phrases and thus assess the presence of musical structure. This model predicts that there will be
408 strong differences across stimuli in the extent to which they give rise to the illusion (due to
409 differences in degree of musical structure), but that the magnitude of the differences across stimuli
410 will vary across participants (due to individual differences in the ability to detect musical structure).
411 The first prediction was confirmed by previous research (Tierney et al., 2013; 2018b; Graber et al.,
412 2017). The second prediction is confirmed here, as we find that there are highly stable individual
413 differences in the disparity in musicality between the illusion and control stimuli. Interestingly,
414 however, we also find that there are stable individual differences in musicality ratings across the
415 entire stimulus set, and across all repetitions (our "musical prior" measure). Nonetheless, although
416 we find both illusion strength and musical prior are reliable measures, they correlated only very
417 weakly. This suggests that an individual's perception of the musicality of a phrase after repetition is
418 determined by at least two main factors: their overall tendency to perceive musicality in sound, and
419 their ability to assess musical structure in complex sounds.

420 We predicted that musical training would relate to the difference in musicality ratings
421 between illusion and control stimuli after repetition (i.e. to illusion strength). This prediction was
422 based on the idea that individuals with musical training would be better able to extract melodic and
423 rhythmic characteristics from stimuli to assess the presence of musical structure. However, we

424 found that while degree of musical training did not relate to illusion strength, it did relate to musical
425 prior, such that participants with more musical training were more likely to produce high musicality
426 ratings across all stimuli after all repetitions. This difference can be clearly seen in Figure 3, in which
427 musicians' musicality ratings are higher than those of non-musicians even for control stimuli after
428 the first repetition. This result aligns with and clarifies Vanden Bosch der Nederlanden et al.'s
429 (2015b) finding that when musicians and non-musicians provided ratings of the original Deutsch et
430 al. (2011) speech-to-song stimulus, musicians provided higher ratings overall, i.e. they heard all
431 stimulus presentations more musically regardless of repetition or transposition. One possible
432 explanation for the relationship between musical training and musical prior is that musicians have
433 sufficiently extensive knowledge of a wide variety of musical styles and examples that they can map
434 any random sequence of pitches and durations onto a possible musical framework to some degree.²
435 The lack of a relationship between illusion strength and musical training has two possible
436 explanations. One possibility is that perception of the illusion does not require extraction of the
437 melodic and rhythmic characteristics of stimuli. A second possibility is that individual differences in
438 melody and rhythm perception (i.e. music aptitude) exist independently of and prior to engagement
439 in formal musical training (Kragness et al., 2020), possibly driven in part by genetic differences
440 (Niarchou et al., 2021). Experiment 2 was designed to examine the relationship between music
441 aptitude and illusion strength (see below).

442 Turning to the effects of age, we did not find any differences between younger and older
443 participants in illusion strength. Prior research has suggested that lexical node satiation decreases
444 with age, as can be seen in decreased verbal transformation effects in older participants (Warren &
445 Warren, 1966; Pilotti & Khurshid, 2004; Pilotti et al., 2011). This result suggests that individual
446 differences in the readiness with which lexical nodes satiate is not a primary factor driving
447 differences in perception of the song illusion in these stimuli. However, this conclusion may not

² An extreme example of this ability can be seen when jazz accompaniment brings out the musical characteristics of natural speech, as in Henry Hey's orchestration of political speeches and interviews (<https://www.youtube.com/watch?v=9nlwwFZdXck>).

448 generalize beyond the current stimulus set, and stimulus sets designed to elicit cross-stimulus
449 differences in lexical satiation (by manipulating, for example, phonological neighborhood density;
450 see Castro et al., 2018) could lead to larger age effects. We also found no differences in either
451 illusion strength or musical prior between participants who identified as male versus female,
452 suggesting that perception of the speech-to-song illusion is relatively unaffected by gender.

453 In terms of the effect of language background, we found no differences in illusion strength
454 between monolinguals and bilinguals, nor between dominant and non-dominant speakers of English.
455 This result suggests that detecting musical characteristics in speech is not enhanced by exposure to
456 multiple languages, and that lack of proficiency in a spoken language does not necessarily modulate
457 perception of the speech-to-song illusion in that language. Initially this may seem to contradict prior
458 research suggesting that language history can modulate the speech-to-song illusion, with song
459 perception diminished for stimuli drawn from unfamiliar languages (Castro et al., 2018), especially
460 difficult-to-pronounce languages (Margulis et al., 2015), but enhanced in familiar second languages
461 in bilinguals (Rathcke et al. 2021). Our findings, however, do not necessarily contradict these
462 findings, as our participants were all somewhat familiar with English, even when English was not
463 their dominant language (although we do not have detailed information about participants'
464 proficiency, limiting somewhat the conclusions that can be drawn about language experience and
465 perception of the illusion). Moreover, these prior studies were not designed to investigate individual
466 differences in the variability in illusion perception between strongly versus weakly transforming
467 stimuli. Our finding of a null effect of language background, however, should be qualified by the lack
468 of tone language speakers in our sample, given previous evidence that tone language speakers
469 perceive the illusion differently (Jaisin et al. 2016) and that tone language ability correlates with
470 melody perception skills (Swaminathan, Kragness, Schellenberg, 2021). Our stimuli could be useful
471 for studying effects of tone language experience on perception of the song illusion in future
472 research, provided the participants are bilingual in English.

473 Overall, we find that although our paradigm can reliably measure individual differences in
474 the strength of the speech-to-song illusion, none of the demographic predictors we measured
475 (including musical training) were significantly associated with illusion strength. What factors, then,
476 might drive individual differences in illusion strength? We hypothesized that individual differences in
477 the ability to extract rhythmic and melodic information from speech could drive variability in illusion
478 strength and investigated this possibility in Experiment 2.

479

480 **3. Experiment 2**

481 *3.1. Introduction: investigating musical aptitude and selective attention to pitch as drivers of the*
482 *song illusion*

483 One factor driving individual differences in perception of the song illusion could be musical aptitude,
484 i.e., musical abilities which vary between individuals independently of any effect of formal musical
485 training (Kragness et al., 2020) and which may partly reflect genetic differences (Wesseldijk, Mosing,
486 & Ullén, 2021; Niarchou et al., 2021). Indeed, prior work has suggested that certain perceptual and
487 cognitive abilities are more strongly tied to musical aptitude than to musical training (Swaminathan,
488 Schellenberg, & Khalil, 2017; Swaminathan & Schellenberg, 2019; Wesseldijk, Gordon, Mosing, &
489 Ullén, 2021). The extent to which speech stimuli sound musical after repetition has been linked to a
490 variety of stimulus characteristics, including the presence of a steady beat and tonal structure
491 (Tierney et al., 2018b) as well as the flatness of pitch within syllables (Tierney et al., 2013; Falk et al.,
492 2014). Detection of these characteristics may be more successful in participants with greater musical
493 aptitude and more precise auditory perception, leading to more robust perception of the speech-to-
494 song illusion. Another possibility is that individual differences in the ability to direct attention to
495 pitch information help determine the strength of the illusion. Prior research has indicated a link
496 between perception of the illusion and increased pitch salience, as shown via increased activation in
497 pitch-sensitive cortical areas (Tierney et al., 2013). That song perception is linked to increased pitch
498 salience is also supported by the finding that imitation of pitch content is enhanced for song relative

499 to speech stimuli (Pfordresher, Mantell, & Pruitt, 2021). Participants who can more readily direct
500 attention to pitch information in speech when instructed to do so, therefore, may be better able to
501 detect latent musical structure in stimulus pitch patterns and therefore more robustly perceive the
502 illusion.

503 To investigate these ideas, in Experiment 2 we examined whether illusion strength was
504 linked to performance on six tests of auditory perception: two tests each of musical aptitude,
505 dimension-selective attention, and psychophysics (Table 1). The musical aptitude tests were the Beat
506 Alignment Test (BAT; Iversen & Patel, 2008), which asks participants to judge whether a series of
507 tone pips is aligned with the beat or shifted away from the beat of clips of music, and the Tonality
508 Alignment Test (TAT), a novel test constructed for the current study which asks participants to judge
509 whether a sung melody is aligned with a tonal grid or misaligned (by being either compressed or
510 expanded). The dimension-selective attention tests (first introduced in Jasmin, Sun, & Tierney, 2021)
511 presented participants with pairs of spoken words which varied orthogonally in relative pitch and
512 relative loudness. Participants were either asked to attend to pitch, indicating which of the two
513 words was higher, or attend to loudness, indicating which of the two words was louder. We
514 predicted that performance on the attention-to-pitch test, but not the attention-to-loudness test,
515 would be linked to the robustness of perception of the speech-to-song illusion. The psychoacoustic
516 tests were pitch and amplitude rise time discrimination (using parameters taken from Kachlicka,
517 Saito, & Tierney, 2019), as we reasoned that participants with more accurate auditory perception
518 may be better able to detect subtle cues to musicality such as slight differences in pitch contour
519 flatness or subtle amplitude envelope cues to beat location. (Amplitude rise time is one of the
520 primary cues conveying the exact timing of musical beats (Danielsen et al., 2019)).

521

522

523

524

Test	Source
Attention to Amplitude	Jasmin et al., 2021
Attention to Pitch	Jasmin et al., 2021
BAT	Iversen & Patel, 2008
TAT	New
Frequency Discrimination	Kachlicka et al., 2019
Rise Time Discrimination	Kachlicka et al., 2019

525 **Table 1.** Summary of measures included as possible predictors of illusion strength.

526 *3.2. Methods*

527 *3.2.1. Participants*

528 Ninety-five (95) participants were initially recruited from the online recruitment service
529 Prolific. Because catch trials were not included in the speech-to-song illusion perception test in this
530 experiment, participant inclusion was based on performance exceeding at least 55% across both
531 dimension-selective attention tests, the BAT, and the TAT, as well as having at least one response for
532 the first and second half of the repetitions during each trial of the speech-to-song illusion paradigm.
533 One additional participant was excluded due to receiving the lowest possible score on both the
534 frequency and amplitude rise time discrimination tests. Seventy-six (76) participants passed these
535 criteria. Forty-eight of these identified as female, while 28 identified as male. A series of power
536 analyses were conducted in GPower to determine the smallest effect sizes which could be detected
537 at a power of 0.8, given this sample size and an alpha of 0.05. This sample size resulted in a power of
538 0.8 to detect a correlation of $r = 0.31$. For the analysis of the speech-to-song ratings, given a 2
539 (stimulus set) X 8 (repetition) RMANOVA design, the sample size resulted in a power of 0.80 to
540 detect main effects and interactions of $F = 0.109$. For the regression analysis, the sample size
541 resulted in a power of 0.80 to detect effect sizes of $F = 2.23$ across six predictors. Participants
542 reported a mean age of 30.7 years (SD 10.5, range 18-67). 74 of the participants reported that their
543 dominant language was English, while 2 participants reported other dominant languages (1 Hindi
544 and 1 Dutch). 62 participants reported being monolingual while 14 participants reported being
545 bilingual (i.e. being able to speak at least two languages). Participants reported a mean of 5.22 years
546 of musical training (std 7.72, range 0-36). Participants were tested using the Gorilla platform for

547 online testing (Anwyl-Irvine et al., 2020). The entire experiment lasted approximately 45 minutes. All
548 participants completed the tasks in the same order, as follows: attention-to-amplitude, attention-to-
549 pitch, song illusion experiment, TAT, BAT, frequency discrimination, and amplitude rise time
550 discrimination. Participants were allowed to rest between tasks, but these breaks were optional
551 rather than mandated. Study procedures were approved by the ethics board of the Department of
552 Psychological Sciences at Birkbeck College.

553 3.2.2. Dimension-Selective Attention Tests

554 The ability to attend to individual acoustic dimensions in speech was assessed with a pair of
555 tests first presented in Jasmin et al. (2021). Stimuli in this test are drawn from a pair of recordings of
556 sentence fragments which are identical lexically but differ in the position of word emphasis: “STUDY
557 music” and “study MUSIC”. These recordings were morphed onto one another using STRAIGHT
558 (Kawahara & Irino, 2005) so that the extent to which individual acoustic dimensions resembled one
559 versus the other recording could be precisely controlled, while all other acoustic characteristics were
560 set to be constant across stimuli. For these tests, a 4 X 4 grid of stimuli was constructed in which
561 pitch and amplitude varied in the extent to which they resembled the pattern in the initial-emphasis
562 recording (“STUDY music”) versus the final-emphasis recording (“study MUSIC”). (Note that phonetic
563 content was fixed across trials, enabling participants to focus on the target dimensions.) Specifically,
564 the pitch levels used were 0%, 33%, 67%, and 100%, where 0% indicates patterns identical to the
565 initial-emphasis recording, 100% indicates patterns identical to the final-emphasis recording, and
566 33% and 67% indicate intermediate values. Similarly, the amplitude levels used were 0%, 33%, 67%,
567 and 100%. To summarize, a set of 16 stimuli were constructed in which the extent to which pitch
568 patterns versus amplitude patterns implied the existence of initial versus final word emphasis was
569 independently varied (See Figure S4 for a schematic of the stimulus design).

570 On each trial, participants were presented with a single stimulus. The two tests involved
571 presentation of stimuli drawn from the same set of 16 possible stimuli but differed in the
572 instructions given to participants. For the Attention to Amplitude test, they were asked to press a

573 button to indicate whether the first or second word was louder, ignoring any differences in pitch. For
574 the Attention to Pitch test, they were asked to press a button to indicate whether the first or second
575 word was higher in pitch, ignoring any differences in amplitude. Performance on each test was
576 summarized as portion correct. For each test 40 trials were presented, with each trial
577 pseudorandomly drawn from the set of 16 possible stimuli (a single randomization was performed
578 for each test and used across all participants). This test has previously been shown to be sensitive to
579 individual differences in dimensional salience: Mandarin speakers were found to display better
580 performance on the Attention to Pitch test and worse performance on the Attention to Amplitude
581 test, compared to native English speakers, suggesting that their experience speaking a tone language
582 led to increased pitch salience across languages (Jasmin et al. 2021). Stimuli from these tests are
583 available at <https://osf.io/hegxs/>.

584

585 3.2.3. Song Illusion Experiment

586 The song illusion experiment was conducted in a manner almost identical to that of
587 Experiment 1, the only difference being that the catch trials were not included. Instead, the first four
588 trials of the test consisted of 8 repetitions of practice items. **These practice items were taken from**
589 **the same speakers as the illusion/control stimuli but were different recorded phrases.** (None of
590 these four practice stimuli were drawn from the main stimulus set, and so cannot be defined as
591 strictly being included among either the illusion or control stimuli. However, the first author, who
592 assembled the original stimulus set (Tierney et al., 2013), perceived two of them as clearly
593 transforming into song and two of them as continuing to be perceived as speech when repeated.)
594 The inclusion of these items was meant to help participants get accustomed to rating the musicality
595 of stimuli; the ratings of these stimuli were not analyzed.

596 3.2.4. TAT

597 To assess participants' ability to determine whether a stream of sung pitches was aligned
598 with a tonal template, we used a new tonality perception test devised for this study, the TAT. (We

599 elected not to use the mistuning perception test (Larrouy-Maestri, Harrison, & Müllensiefen, 2019)
600 because it assesses whether listeners can detect a pitch shift between a vocal line and an
601 instrumental accompaniment, whereas the illusion induces melody perception in a single
602 unaccompanied vocal recording.) Stimuli were drawn from the DSD100 (Liutkus et al., 2016), a
603 dataset of 100 publicly available music recordings in which each component instrument can be
604 downloaded separately. A single short portion of the vocal track from each of 22 different recordings
605 was extracted. The duration of these musical passages ranged from 9.5 to 24.6 seconds (mean 15.7,
606 SD 3.9). These 22 stimuli were then divided into three different sets that underwent three different
607 types of pitch morphing in Praat (Boersma & Weenink, 2021). For 6 of the stimuli, the pitch contour
608 was expanded by 30% on a time-point-by-time-point basis (Zatorre & Baum, 2012). This was done by
609 first extracting the fundamental frequency (F0) of each phrase using Praat (default settings, with a
610 time step of 0.01 seconds), then comparing the ratio between the F0 value at each time point and
611 the median F0 across all time points. The F0 of each time point was then adjusted to be equal to 1.3
612 times its original distance from the median F0, in semitones. (Previous perceptual research suggests
613 that the semitone scale is more relevant to the perception of pitch in speech than the Hz scale
614 [Nolan, 2003].) For 5 of the stimuli, F0 contours were contracted by 30% in a similar manner: the F0
615 of each time point was adjusted to be equal to 0.7 times its original distance from the median F0.
616 For 11 of the stimuli, the F0 values were set to be equivalent to their original values. However, to
617 ensure that any distortions due to the F0 manipulations were present across all stimuli, the
618 “unaltered” stimuli were constructed by first expanding the F0 contours by 30% and then
619 contracting them back to their original values.

620 Participants were told they would hear 22 melodies, some of which would be in tune, and
621 some of which would not be in tune. They were asked to indicate whether the melody they heard
622 was in tune (by pressing a button marked “in tune”) or out of tune (by pressing a button marked
623 “out of tune”). Task performance was summarized by calculating sensitivity (d'). Stimuli from this
624 test are available at <https://osf.io/hegxs/>.

625

626 3.2.5. BAT

627 This beat perception test used 22 stimuli drawn from Iversen & Patel (2008). These were
628 instrumental musical excerpts from a broad variety of musical genres, including rock, classical, and
629 jazz, onto which a sequence of 1 KHz 100 ms pure tones was superimposed. These tones were either
630 aligned with the beat or shifted away by 25% of the inter-beat-interval. Participants were asked to
631 indicate whether the tones were aligned with the beat or shifted away from the beat (by pressing
632 buttons marked “on the beat” or “off the beat”). As soon as participants selected their response, the
633 next stimulus was presented. A difference between this test and that reported in Iversen & Patel
634 (2008) is that only a short excerpt of each stimulus was presented (long enough to contain seven
635 beats, average duration 3.81 seconds), to keep the experiment relatively brief. Task performance
636 was summarized by calculating sensitivity (d'). Stimuli from this test are available at
637 <https://osf.io/hegxs/>.

638

639 3.2.6. Self-assessment of Music Perception Skills

640 Because the TAT is a novel measure, it is of interest to establish its validity by comparing
641 performance on this test to self-assessments of tonality perception ability. In addition, we were
642 interested in comparing the validity of the TAT to that of the BAT, an already established measure in
643 the field. Participants were asked to indicate on a scale from 1 to 7 the extent to which they agreed
644 with each of four statements, with 1 indicating “not at all” and 7 indicating “agree completely”.
645 These statements were: 1) “I can tell when people sing or play out of tune”, 2) “When I sing I have
646 no idea whether I’m in tune or not”, 3) “I can tell when people sing or play out of time with the
647 beat”, and 4) “When I clap to music I have no idea whether I’m on the beat or not”. A composite
648 measure of self-assessed tonality perception was calculated by averaging the response to question 1
649 with the inverted response to question 2. Similarly, a composite measure of self-assessed beat

650 perception was calculated by averaging the response to question 3 with the inverted response to
651 question 4.

652

653 3.2.7. Psychoacoustic Discrimination

654 Two tests were conducted to examine the precision of participants' perception of frequency
655 and amplitude rise time. Each test consisted of a single 2-down 1-up adaptive staircase design. In
656 each test, stimuli were drawn from a continuum of 101 stimuli, representing 101 different levels of
657 the target acoustic continuum (i.e. either frequency or rise time). In each trial, participants were
658 presented with 3 sounds, with either the first or the third being different from the other two, which
659 were identical. The two identical sounds always corresponded to level 1 of the target continuum.
660 Participants were asked to indicate which of the three sounds was different by pressing either a
661 button labelled "1" or a button labelled "3". Initially, the target sound level (in other words, the level
662 of the different sound) was set at 50. After every two correct responses, the target level decreased,
663 becoming more similar to the comparison level, while after every incorrect response, the target level
664 increased, becoming more dissimilar to the comparison level. The size of these level
665 increases/decreases changed across the block, becoming smaller after each subsequent "reversal"
666 or inflection point (in other words, two correct responses following a set of incorrect responses, or
667 one incorrect response following a set of at least two correct responses). These step sizes were 10,
668 5, 2, 1, 1, 1, and 1 before the first through seventh reversals, respectively. The block continued until
669 75 trials were presented or participants reached the seventh reversal, whichever came first.

670 The stimulus continua for the two tests were constructed as follows. For the frequency
671 discrimination test, the comparison stimulus was always presented at a fundamental frequency (F0)
672 of 330 Hz, while the target stimulus ranged from 330.3 to 360 Hz in 100 equal linear steps.
673 Frequency discrimination stimuli were 500-ms four-harmonic complex tones with equal amplitude
674 across harmonics with a 0.015 linear amplitude ramp at the beginning and end to avoid perception
675 of transient clicks. For the amplitude rise time test, the comparison stimulus was always presented

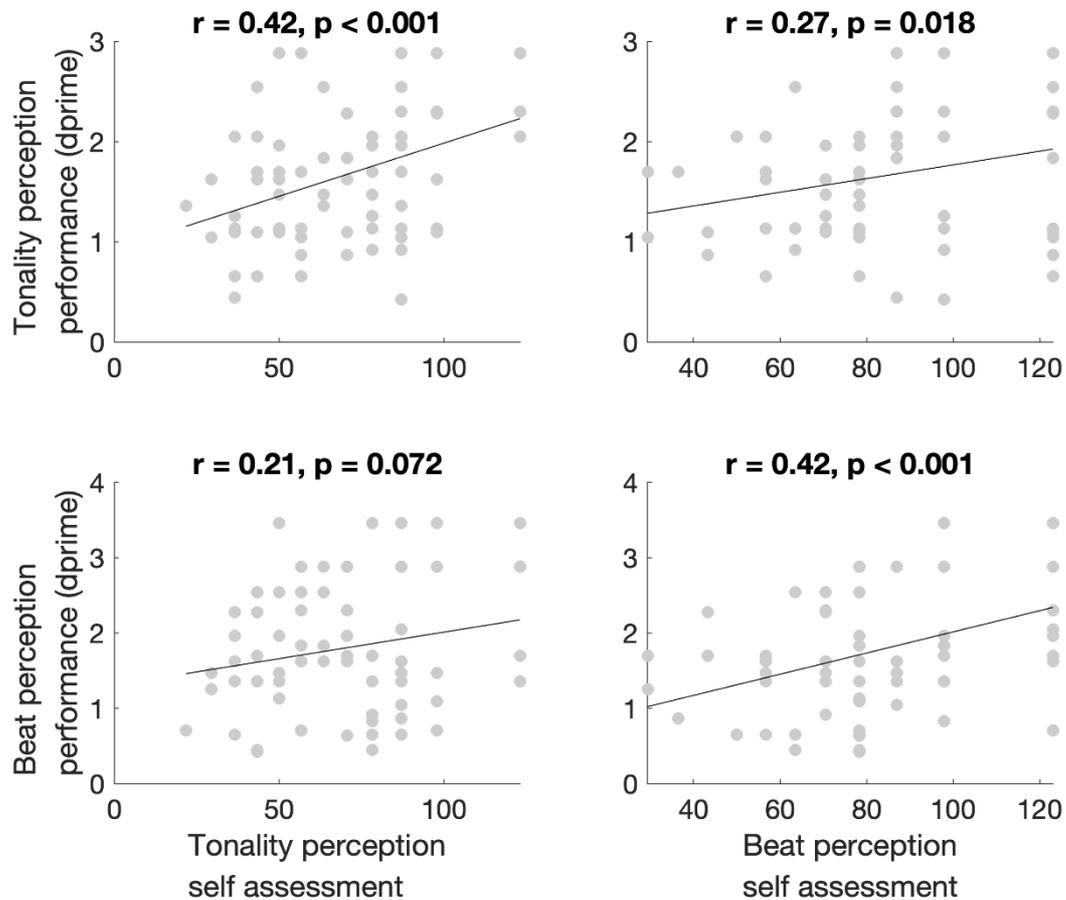
676 with a linear amplitude ramp at the beginning of the stimulus of duration 15 ms, while the target
677 stimulus ranged from a rise time of 17.8 to 300 ms in 100 equal linear steps. Rise time stimuli were
678 500-ms four-harmonic complex tones with equal amplitude across harmonics and an F0 of 330 Hz.
679 Thresholds for each test were calculated as the mean of the target stimulus levels at the second
680 through final reversals.

681 Any variables that were shown to be non-normally distributed according to a Jacque-Bera
682 test were transformed prior to analysis. A rau transformation was used for the attention-to-pitch
683 test and beat and tonality self-assessments, a log transformation was used for the frequency and
684 rise time discrimination thresholds, and years of musical training were converted to ranks.

685 3.3. Results

686 As in Experiment 1, across participants the song illusion was very reliably experienced, as
687 seen in the strong interaction between repetition and stimulus set ($F(7,525) = 85.3, p < 0.001$). There
688 was once again a main effect of repetition ($F(1,525) = 50.9, p < 0.001$) and of stimulus set ($F(1,75) =$
689 $182.0, p < 0.001$). Illusion strength and musical prior did not correlate ($r(74) = -0.12, p = 0.318$),
690 suggesting that these two measures index dissociable aspects of the speech-to-song illusion.

691 To test the validity of the TAT and BAT, performance on these measures was correlated with
692 self-ratings of tonality perception (the ability to detect whether one's own singing and others'
693 singing is in tune) as well as self-ratings of beat perception (the ability to detect whether one's own
694 singing/playing and others' singing/playing is on the beat). BAT performance was correlated with
695 beat perception self-rating ($r(74) = 0.42, p < 0.001$) but not tonality perception self-rating ($r(74) =$
696 $0.21, p = 0.072$). TAT performance was correlated with both tonality perception self-rating ($r(74) =$
697 $0.42, p < 0.001$) and beat perception self-rating ($r(74) = 0.27, p = 0.018$). Figure 4 displays the
698 relationship between TAT and BAT performance and tonality and beat perception self-rating. The
699 finding that performance on the TAT relates to self-assessment of both beat and tonality processing
700 suggests that while this measure may be a valid measure of musical aptitude, its specificity in
701 assessing tonality per se may not be optimal.



703

704

Figure 4. Upper left panel, relationship between tonality perception self-assessment and

705

performance on the TAT. Upper right panel, relationship between beat alignment perception self-

706

assessment and performance on the TAT. Lower left panel, relationship between tonality perception

707

self-assessment and performance on the BAT. Lower right panel, relationship between beat

708

perception self-assessment and performance on the BAT.

709

710

To investigate whether performance on the musical aptitude and auditory perception tests

711

was associated with musical training, Pearson correlations were used to investigate the relationship

712

between musical training and these measures. These correlations were corrected for multiple

713

comparisons using false discovery rate, Benjamini & Hochberg 1995. Years of musical training were

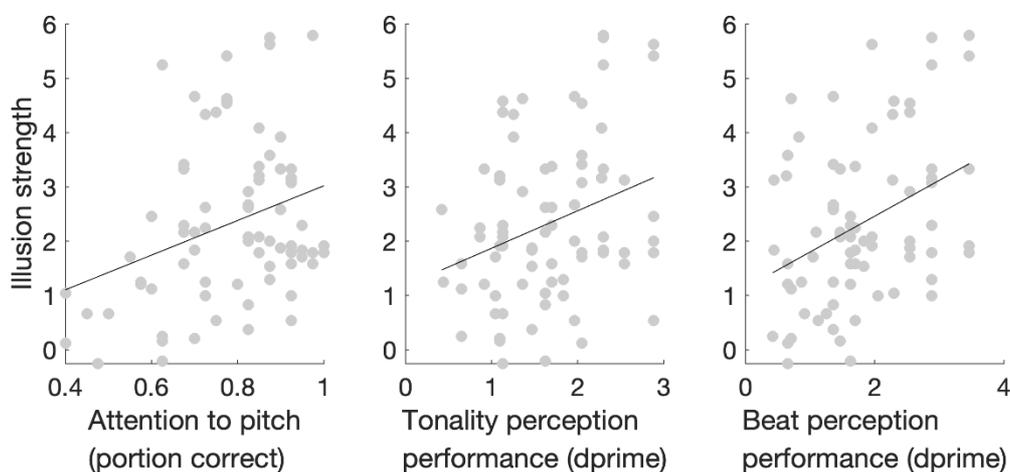
714

not significantly correlated with performance on any of the musical and auditory tests, including

715 attention to amplitude ($r(74) = -0.05$, $p(\text{corrected}) = 0.674$), attention to pitch ($r(74) = 0.25$,
716 $p(\text{corrected}) = 0.092$), TAT ($r(74) = 0.30$, $p(\text{corrected}) = 0.053$), BAT ($r(74) = 0.22$, $p(\text{corrected}) =$
717 0.116), frequency discrimination ($r(74) = -0.06$, $p(\text{corrected}) = 0.674$), and amplitude rise time
718 discrimination ($r(74) = -0.13$, $p(\text{corrected}) = 0.413$). These nonsignificant correlations underscore the
719 distinction between musical aptitude and musical training as measures of musical ability (Kragness
720 et al., 2020).

721 Pearson correlations were used to investigate the relationship between musical aptitude
722 and the strength of perception of the speech-to-song illusion, which was operationalised as the
723 difference in the average ratings of illusion and control stimuli after the eighth repetition (i.e.,
724 illusion strength in Experiment 1). These correlations were corrected for multiple comparisons using
725 false discovery rate (Benjamini & Hochberg 1995). Illusion strength was correlated with the ability to
726 direct attention to pitch within spoken phrases ($r(74) = 0.26$, $p(\text{corrected}) = 0.044$) but not with the
727 ability to direct attention to amplitude within spoken phrases ($r(74) = 0.09$, $p(\text{corrected}) = 0.454$).
728 Illusion strength was correlated with performance on the TAT ($r(74) = 0.30$, $p(\text{corrected}) = 0.027$)
729 and performance on the BAT ($r(74) = 0.38$, $p(\text{corrected}) = 0.005$). Psychophysical discrimination
730 thresholds were not correlated with illusion strength, including frequency discrimination ($r(74) =$
731 0.19 , $p(\text{corrected}) = 0.147$) and amplitude rise time discrimination ($r(74) = -0.14$, $p(\text{corrected}) =$
732 0.256). Figure 6 displays the relationship between illusion strength and musical aptitude.

733



734

735 **Figure 5.** Left panel, relationship between attention to pitch and final illusion-control rating
736 difference. Middle panel, relationship between TAT performance and final illusion-control rating
737 difference. Right panel, relationship between BAT performance and final illusion-control rating
738 difference.

739

740 Pearson correlations were used to investigate the relationship between musical aptitude
741 and musical prior, which was operationalised as mean rating across all stimuli over all eight
742 repetitions. These correlations were corrected for multiple comparisons using false discovery rate
743 (Benjamini & Hochberg 1995). Musical prior was not correlated with performance on any of the
744 auditory or musical tests, including attention to amplitude ($r(74) = -0.27$, $p(\text{corrected}) = 0.114$),
745 attention to pitch ($r(74) = 0.20$, $p(\text{corrected}) = 0.273$), TAT ($r(74) = 0.02$, $p(\text{corrected}) = 0.894$), BAT
746 ($r(74) = 0.11$, $p(\text{corrected}) = 0.495$), frequency discrimination ($r(74) = 0.02$, $p(\text{corrected}) = 0.894$), and
747 amplitude rise time discrimination ($r(74) = -0.16$, $p(\text{corrected}) = 0.332$).

748 To examine whether the predictors explain independent variance in illusion strength,
749 backwards linear regression was used, starting with a full model containing all (centered, scaled)
750 predictors and removing predictors based on AIC. (The Variance Inflation Factor for all predictors
751 was less than 1.6 in the full and reduced models, suggesting that multicollinearity was not a major
752 problem in this dataset.) The resulting model (see Table 2) explained 22% of variance in illusion
753 strength ($F(3,72) = 6.80$, $p < 0.001$). Predictors contained in the model included frequency
754 discrimination threshold ($\beta = 0.306$, $t = 1.964$, $p = 0.053$), BAT performance ($\beta = 0.449$, $t =$
755 2.764 , $p = 0.007$), and TAT performance ($\beta = 0.340$, $t = 2.083$, $p = 0.041$). See Figure 7 for a
756 depiction of predicted versus actual illusion strength values across participants.

757

758

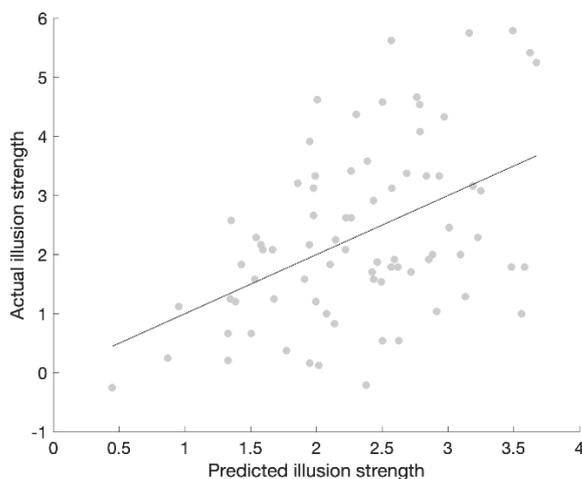
759

760

	Coefficient	Std. error	t value	p value
(Intercept)	2.327	0.154	15.149	< 0.001
Frequency discrimination	0.306	0.156	1.964	0.053
BAT	0.449	0.162	2.764	0.007
TAT	0.340	0.163	2.083	0.041

761 **Table 2.** Regression predicting cross-participant differences in illusion strength.

762



763

764 **Figure 6.** Scatterplot displaying predicted (x-axis) versus actual (y-axis) illusion strength values across
765 participants.

766

767 3.4. Discussion

768 We found that musical aptitude—including performance on the Tonality Alignment Test
769 (TAT) and Beat Alignment Test (BAT)—as well as the ability to direct attention to the pitch of speech
770 were linked to the strength of the speech-to-song illusion (the difference in musicality ratings
771 between strongly transforming versus weakly transforming stimuli) but not to musical prior (the
772 overall tendency to rate stimuli as musical). This result suggests that a key factor driving the song
773 illusion is perceptual sensitivity to the degree of musical structure present in speech. Our finding of a
774 link between individual differences in musical aptitude and individual differences in perception of
775 illusion strength confirms a key prediction of the hybrid account of the speech-to-song illusion, in

776 which satiation of lexical nodes frees up cognitive resources so that listeners can focus on extraction
777 of rhythmic and melodic information from stimuli. Our results suggest that not all listeners are
778 equally able to assess the musical structure of acoustic stimuli, and that this assessment skill may be
779 driven by sensitivity to the details of musical pitch and rhythmic patterns. Importantly, we find that
780 musical aptitude is not linked to individual differences in musical prior, suggesting that our results
781 are not driven by an overall tendency for certain listeners to perceive stimuli as musical, but truly
782 reflect differences in the ability to assess musical structure in non-musical stimuli.

783 The strongest predictor of the strength of the speech-to-song illusion was performance on
784 the BAT. The presence of an isochronous beat is one of the most reliable characteristics
785 distinguishing speech from music cross-culturally (Savage, Brown, Sakai, & Currie, 2015). Moreover,
786 the presence of a steady pulse has also been shown to be a useful feature when constructing
787 automatic classifiers of speech versus music (Scheirer & Slaney, 1997). Listeners may, therefore, be
788 broadly aware that song tends to contain an isochronous beat while speech does not and use the
789 presence or absence of a musical beat as relevant evidence when categorizing a stimulus as speech
790 versus song. Indeed, we have previously shown that within the stimulus set used in the current
791 study, stimuli with more isochronous beats (as extracted via a computational model of beat
792 perception) are perceived as more musical after repetition (Tierney et al., 2018b). Beat perception is
793 not limited to musicians, but rather is a broadly present skill in the general population: although
794 neural entrainment to musical beats is enhanced by musical training, it is clearly present in non-
795 musicians (Doelling & Poeppel, 2015). There is some evidence that beat perception emerges very
796 early in development: infants coordinate their movements with the musical tempo based on musical
797 pulse clarity (Zentner & Eerola, 2010) and can integrate auditory and somatosensory/proprioceptive
798 information when deciding which notes carry musical beats (Phillips-Silver & Trainor, 2005).
799 Nevertheless, although basic competence in beat perception is widespread, there are large
800 individual differences in this ability (Müllensiefen, Gingras, Musil, & Stewart, 2014; Tranchant,
801 Vuvan, & Peretz, 2016; Tranchant, Lagrois, Bellemare, Schultz, & Peretz, 2021), and impaired beat

802 perception may limit individuals' ability to detect musical structure in speech and other non-musical
803 stimuli.

804 Performance on the TAT was another significant predictor of the strength of the speech-to-
805 song illusion. The use of discrete scales is another characteristic that distinguishes speech from
806 music cross-culturally (Savage, Brown, Sakai, & Currie, 2015). Listeners may, as a result, take the
807 presence of tonality as evidence that a stimulus should be categorized as song rather than speech.
808 Tonality perception, like beat perception, is widespread in the general population: Western
809 European musicians and non-musicians show a preference for scale structure in melodic sequences
810 (Cross, Howell, & West, 1983). Our findings, however, suggest that not everyone is equally able to
811 detect the presence of scale structure, and that individual differences in this ability help determine
812 variability in the robustness of musicality perception in non-musical stimuli. The correlation between
813 tonality perception and speech-to-song perception (as well as with self-assessment of tonality
814 perception skills) demonstrates the validity of this novel measure, which may be useful in a variety
815 of future studies of tonality perception, including investigations of the relationship between musical
816 and language skills, studies of the effects of musical training, and research on the development of
817 musical abilities.

818 We find that the ability to direct attention to pitch in speech is linked to the robustness of
819 perception of the speech-to-song illusion. This relationship, however, does not reflect the influence
820 of general attentional skills, as the ability to direct attention to amplitude in speech was not linked
821 to perception of the illusion. This finding aligns with the proposal by Deutsch et al. (2011) that
822 perception of musical characteristics in speech is inhibited by default because pitch salience is down-
823 regulated during speech perception, but that repetition of speech disinhibits attention to pitch,
824 making possible the detection of musical characteristics in stimuli. Several findings support the
825 possibility that variation in other acoustic dimensions can interfere with pitch perception, possibly
826 by diverting attention away from pitch. For example, Warrier & Zatorre (2002), Allen & Oxenham
827 (2014), and Caruso & Balaban (2014) showed that the presence of variation in timbre can interfere

828 with pitch perception, while Russo, Vuvan, & Thompson (2019) found that vowel content can
829 interfere with relative pitch perception in speech stimuli. The idea that perceiving the speech-to-
830 song illusion can lead to increased pitch salience is supported by the finding that pitch-sensitive
831 cortical areas increase in activation when illusion stimuli are repeated (Tierney et al., 2013). This
832 increase in salience may be greater in individuals who can more readily direct attention to pitch,
833 leading to more robust perception of the illusion.

834 **4. General discussion**

835 Writing in 55 BCE, Cicero remarked in *De Oratore* that “even in speaking there may be a
836 concealed kind of music.” Over 2,000 years later, the discovery of the song illusion provided
837 compelling evidence for this claim (Deutsch, 2003; Deutsch et al., 2011). The current work on
838 individual differences in the experience of the illusion suggests that Cicero’s phrase should be
839 updated: “even in speaking there may be a concealed kind of music, *though only some can hear it.*”
840 Specifically, our results suggest that the ability to focus on and extract several different kinds of
841 auditory and musical cues is linked to the strength of the song illusion. These results align with prior
842 work on other auditory illusions, which has generally shown that the salience of (and ability to
843 perceive) different stimulus characteristics is linked to the influence these characteristics have on
844 the final percept. In the “Yanny/Laurel” illusion, for example, whether individuals perceive “Yanny”
845 versus “Laurel” is linked to their degree of prior exposure to lower versus higher auditory
846 frequencies in their environment (Gwilliams & Wallisch, 2020). Similarly, the magnitude of the
847 McGurk effect correlates with lipreading skill (Strand, Cooperman, Rowe, & Simenstad, 2014; Brown
848 et al., 2018), as well as the amount of time individuals spend fixating their gaze on a talker’s mouth
849 during audiovisual speech perception (Gurler, Doyle, Walker, Magnotti, & Beauchamp, 2015). The
850 song illusion, therefore, like other illusions, may involve the weighting of multiple potential cues to
851 which individuals are differentially sensitive.

852 Our findings support theoretical models of the song illusion which suggest that perception of
853 the illusion requires extraction of musical characteristics from the stimuli (Deutsch et al. 2011,

854 Tierney et al. 2018b). The need to extract musical characteristics from the stimuli could partly
855 explain the increase in musicality perception with repetition. Prior work has shown that perception
856 of pitch interval size is surprisingly poor after a single presentation of a melody, and that melody
857 repetition is necessary before participants achieve consistent interval perception (Deutsch, 1979).
858 The idea that the increase in song illusion strength with repetition is driven by the time it takes to
859 extract musical characteristics from the stimuli is supported by the finding that the increase in
860 musicality with repetition is similarly sized when speech stimuli and synthesized complex tone
861 stimuli with identical pitch contours are presented (Tierney et al., 2018a).

862 Our findings constrain theories of the song illusion in significant ways. The finding that
863 individual differences in perception of the illusion are linked to musical aptitude is not predicted by
864 models which suggest that the illusion is driven by deactivation of lexical nodes and continued
865 activation of syllable nodes (Castro et al., 2018). While our results indicate that Node Structure
866 Theory does not provide a *sufficient* explanation of the song illusion, our results do not imply that
867 deactivation of lexical nodes is *not* relevant to perception of the illusion. Our experiment was not
868 designed to test the theory that individual differences in the speed or robustness of lexical satiation
869 drive individual differences in perception of the illusion, an idea worth exploring in future research.

870 Here we find that musical aptitude predicts the strength of the song illusion, but musical
871 training does not. This may seem somewhat contradictory since the main goal of musical training is
872 to boost musical skills. Indeed, prior work has shown that musicians demonstrate more precise
873 mistuning perception (Hutchins, Roquet, & Peretz, 2012; Larrouy-Maestri, 2018) and stronger neural
874 entrainment to musical beats (Doelling & Poeppel, 2015; but see Hickey, Merseal, Patel, & Race,
875 2020). However, the lack of a relationship between musical training and illusion perception aligns
876 with other recent work finding that cognitive/perceptual abilities are more strongly linked to musical
877 aptitude than to degree of musical training (Swaminathan, Schellenberg, & Khalil, 2017;
878 Swaminathan & Schellenberg, 2019). One possibility is that strong individual differences in musical
879 aptitude exist prior to music training, and in fact help determine whether individuals begin (and stick

880 with) musical training (Kragness et al., 2020). Overall, our results suggest that individual differences
881 in the perception of musicality in speech are tied to musical skills that do not depend strongly on
882 formal musical training, possibly reflecting stable traits influenced by genetic differences between
883 individuals (Wesseldijk et al., 2021; Niarchou et al., 2021). It should be noted that we measured
884 musical training, rather than musical experience, and therefore it remains an open question whether
885 aspects of perception of the speech-to-song illusion are related to musical experience (which future
886 research could investigate using the Gold-MSI: Müllensiefen et al., 2014).

887 Although our finding of a relationship between musical aptitude and illusion strength
888 represents an initial step towards understanding the factors driving individual differences in
889 perception of the speech-to-song illusion, the results of our linear regression explained only 22% of
890 the variance in illusion strength. It is likely, therefore, that other major factors driving individual
891 differences in perception of the illusion remain to be uncovered. One possibility is that individuals
892 who tend to focus on the F0 rather than spectral information when assessing pitch contour may
893 more robustly perceive the illusion. Research on perception of the missing fundamental illusion has
894 shown that there exist large individual differences in the extent to which listeners rely on F0 versus
895 spectral information when judging the interval between two notes (Schneider et al., 2005; Ladd et
896 al., 2013). These divergent listening strategies may have consequences for the ability to extract
897 musical information from sound sequences. Although pitch contour can be as easily extracted from
898 inharmonic compared to harmonic sounds, the ability to extract F0 information facilitates judgment
899 of exact pitch intervals and tonality (McPherson & McDermott, 2018) as well as memory for pitch
900 (McPherson & McDermott, 2020). As a result, spectrally-biased listeners may have difficulty
901 detecting musical regularities latent in the speech-to-song illusion stimuli. Another possible factor
902 driving individual differences in perception of the illusion is the extent to which individuals find
903 music rewarding or absorbing: listeners with a greater hedonic or absorptive response to music may
904 be more keen to seek out musical characteristics in non-musical stimuli. This hypothesis could be
905 tested in future work using the Barcelona Musical Reward Questionnaire (Mas-Herrero, Marco-

906 Pallares, Lorenzo-Seva, Zatorre, & Rodriguez-Fornells, 2013) or the Absorption in Music Scale
907 (Sandstrom and Russo, 2013). Answers to specific questions within these scales may prove
908 especially useful for helping predict susceptibility to the song illusion, e.g., “I like to find patterns in
909 everyday sounds” from the Absorption in Music Scale. Individual differences in the vividness and
910 control of auditory imagery could help predict illusion strength as well (Halpern, 2015).

911 In conclusion, we find that there are strong, reliable individual differences in the tendency to
912 perceive certain repeated spoken phrases as sung. These differences seem to be largely independent
913 of demographic characteristics, including language and musical training, but are linked to individual
914 differences in specific perceptual and musical skills. This finding suggests that individual differences
915 in auditory abilities may strongly affect perceptual categorization, influencing not only categorization
916 judgments within a domain (such as speech perception; Jasmin, Dick, Holt, & Tierney, 2020) but even
917 perceptual categorization between domains (here, speech versus music). A deeper understanding of
918 these individual differences would help researchers predict whether a given listener will experience
919 a particular spoken phrase as sung when repeated. This in turn would provide a powerful tool for
920 exploring the cognitive and brain mechanisms underlying selective neural responses to speech and
921 music (Ogg, Moraczewski, Kuchinsky, & Slevc, 2019; Zuk, Teoh, & Lalor, 2020; Boebinger, Norman-
922 Haignere, McDermott, & Kanwisher, 2021), by using the same physical stimuli to elicit categorically
923 different perceptual experiences.

924

925 **5. Acknowledgements**

926 We would like to thank Peter Pfordresher, Jonathan De Souza, and three anonymous reviewers for
927 their helpful comments on previous versions of this manuscript.

928

929 **6. References**

930 Allen, E., & Oxenham, A. (2014). Symmetric interactions and interference between pitch and timbre.
931 *JASA*, *135*, 1371-1379.

932

933 Anwyl-Irvine, A., Massonnié, J., Flitton, A., Kirkham, N., & Evershed, J. (2020). Gorilla in our midst: an
934 online behavioral experiment builder. *Behavior Research Methods*, *52*, 388-407.

935

936 Benjamini, Y., & Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful
937 approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)*, *57*,
938 289-300.

939

940 Bigand, E., Delbé, C., Gerard, Y., & Tillmann, B. (2011). Categorization of extremely brief auditory
941 stimuli: domain-specific or domain-general processes? *PLoS ONE* *6*, e27024.

942

943 Boebinger, D., Norman-Haignere, S. V., McDermott, J. H., & Kanwisher, N. (2021). Music-selective
944 neural populations arise without musical training. *Journal of Neurophysiology*, *125*(6), 2237-2263.

945

946 Boersma, P., & Weenink, D. (2021). Praat: doing phonetics by computer [Computer program].
947 Version 6.1.42, retrieved 15 April 2021 from <http://www.praat.org/>

948

949 Brennan, D., & Stevens, C. (2002). Specialist musical training and the octave illusion: analytical
950 listening and veridical perception by pipe organists. *Acta Psychologica*, *109*, 301-314.

951

952 Brown, V., Hedayati, M., Zanger, A., Mayn, S., Ray, L., Dillman-Hasso, N., & Strand, J. (2018). What
953 accounts for individual differences in susceptibility to the McGurk effect? *PLoS ONE*, *13*, e0207160.

954

955 Caruso, V. C., & Balaban, E. (2014). Pitch and timbre interfere when both are parametrically varied.
956 *PLoS ONE*, *9*(1), e87065.

957

958 Castro, N., Mendoza, J., Tampke, E., & Vitevitch, M. (2018). An account of the speech-to-song illusion
959 using Node Structure Theory. *PLoS ONE*, *13*, e0198656.

960

961 Cicero, M. *On Oratory and Orators. With Notes Historical and Explanatory. A New Edition, Carefully*
962 *Revised and Corrected. In Two Volumes 1808.* Translation by William Guthrie.

963

964 Craig, J. (1979). The effect of musical training and cerebral asymmetries on perception of an auditory
965 illusion. *Cortex*, *15*, 671-677.

966

967 Cross, I., Howell, P., & West, R. (1983). Preferences for scale structure in melodic sequences. *Journal*
968 *of Experimental Psychology: Human Perception and Performance*, *9*, 444–460.

969

970 Danielsen, A., Nymoen, K., Anderson, E., Câmara, G., Langerød, M., Thompson, M., & London, J.
971 (2019). Where is the beat in that note? Effects of attack, duration, and frequency on the perceived
972 timing of musical and quasi-musical sounds. *Journal of Experimental Psychology: Human Perception*
973 *and Performance*, *45*, 402-418.

974

975 Davidson, B., Power, R., & Michie, P. (1987). The effects of familiarity and previous training on
976 perception of an ambiguous musical figure. *Perception and Psychophysics*, *41*, 601-608.

977

978 Deutsch, D. (1975). Musical illusions. *Scientific American*, 233, 92-105.

979

980 Deutsch, D. (1995). *Musical Illusions and Paradoxes*. Philomel Records.

981

982 Deutsch, D. (2003). *Phantom Words and Other Curiosities*. Philomel Records.

983

984 Deutsch, D. (1979). Octave generalization and the consolidation of melodic information. *Canadian*

985 *Journal of Psychology*, 33, 201-205.

986

987 Deutsch, D., Henthorn, T., & Lapidis, R. (2011). Illusory transformation from speech to song. *JASA*,

988 129, 2245-2252.

989

990 Ding, N., Patel, A. D., Chen, L., Butler, H., Luo, C., & Poeppel, D. (2017). Temporal modulations in

991 speech and music. *Neuroscience & Biobehavioral Reviews*, 81, 181-187.

992

993 Doelling, K., & Poeppel, D. (2015). Cortical entrainment to music and its modulation by expertise.

994 *PNAS*, 112, E6233-E6242.

995

996 Falk, S., Rathcke, T., & Dalla Bella, S. (2014). When speech sounds like music. *Journal of Experimental*

997 *Psychology: Human Perception and Performance*, 40, 1491-1506.

998

999 Graber, E., Simchy-Gross, R., & Margulis, E. (2017). Musical and linguistic listening modes in the

1000 speech-to-song illusion bias timing perception and absolute pitch memory. *Journal of the Acoustical*

1001 *Society of America*, 142, 3593-3602.

1002

1003 Gurler, D., Doyle, N., Walker, E., Magnotti, J., & Beauchamp, M. (2015). A link between individual
1004 differences in multisensory speech perception and eye movements. *Attention, Perception, and*
1005 *Psychophysics*, *77*, 1333-1341.

1006

1007 Gwilliams, L., & Wallisch, P. (2020). Immediate ambiguity resolution in speech perception based on
1008 prior acoustic experience. *PsyArXiv*. doi:10.31234/osf.io/yptfr

1009

1010 Halpern, A. (2015). Differences in auditory imagery self-report predict neural and behavioral
1011 outcomes. *Psychomusicology: Music, Mind, and Brain*, *25*, 37-47.

1012

1013 Hickey, P., Merseal, H., Patel, A. D., & Race, E. (2020). Memory in time: Neural tracking of low-
1014 frequency rhythm dynamically modulates memory formation. *Neuroimage*, *213*, 116693

1015

1016 Hutchins, S., Roquet, C., & Peretz, I. (2012). The vocal generosity effect: how bad can your singing
1017 be? *Music Perception*, *30*, 147-159.

1018

1019 Irwin, J., Whalen, D., & Fowler, C. (2006). A sex difference in visual influence on heard speech.
1020 *Perception & Psychophysics*, *68*, 582-592.

1021

1022 Iversen, J. R., & Patel, A. D. (2008). The Beat Alignment Test (BAT): Surveying beat processing abilities
1023 in the general population. In K. Miyazaki, M. Adachi, Y Hiraga, Y Nakajima, & M. Tsuzaki (Eds.),
1024 *Proceedings of the 10th International Conference on Music Perception & Cognition (ICMPC10)* (CD-
1025 ROM; pp. 465–468). Adelaide, Australia: Causal Productions.

1026

1027 Jaisin, K., Suphanchaimat, R., Candia, M., & Warren, J. (2016). The speech-to-song illusion is reduced
1028 in speakers of tonal (vs. non-tonal) languages. *Frontiers in Psychology*, *7*, 662.

1029

1030 Jarque, C., & Bera, A. (1980). Efficient tests for normality, homoscedasticity and serial independence
1031 of regression residuals. *Economics Letters*, *6*, 255-259.

1032

1033 Jasmin, K., Dick, F., Holt, L., & Tierney, A. (2020). Tailored perception: individuals' speech and music
1034 perception strategies fit their perceptual abilities." *Journal of Experimental Psychology: General*, *149*,
1035 914-934.

1036

1037 Jasmin, K., Sun, H., & Tierney, A. (2021). Effects of language experience on domain-general
1038 perceptual strategies. *Cognition*, *206*, 104481.

1039

1040 Kachlicka, M., Saito, K., & Tierney, A. (2019). Successful second language learning is tied to robust
1041 domain-general auditory processing and stable neural representation of sound. *Brain and Language*,
1042 *192*, 15-24.

1043

1044 Kawahara, H., & Irino, T. (2005). Underlying principles of a high-quality speech
1045 manipulation system STRAIGHT and its application to speech segregation. In P.
1046 Divenyi (Ed.), *Speech separation by humans and machines* (pp. 167–180). Boston,
1047 MA: Kluwer Academic Publishers.

1048

1049 Kragness, H., Swaminathan, S., Cirelli, L., & Schellenberg, G. (2020). Individual differences in musical
1050 ability are stable over time in childhood. *Developmental Science*. doi:10.1111/desc.13081

1051

1052 Krizman, J., Bonacina, S., & Kraus, N. (2019). Sex differences in subcortical auditory processing
1053 emerge across development. *Hearing Research*, *380*, 166-174.

1054

1055 Ladd, R., Turnbull, R., Browne, C., Caldwell-Harris, C., Ganushchak, L., Swoboda, K., Woodfield, V., &
1056 Dediu, D. (2013). Patterns of individual differences in the perception of missing-fundamental tones.
1057 *Journal of Experimental Psychology: Human Perception and Performance*, *39*, 1386-1397.
1058
1059 Larrouy-Maestri, P. (2018). I know it when I hear it: on listeners' perception of mistuning. *Music &*
1060 *Science*, *1*, 2059204318784582.
1061
1062 Larrouy-Maestri, P., Harrison, P., & Müllensiefen, D. (2019). The mistuning perception test: a new
1063 measurement instrument. *Behavior Research Methods*, *51*, 663-675.
1064
1065 Liutkus, A., Stöter, F., Rafii, Z., Kitamura, D., Rivet, B., Ito, N., Ono, N., & Fontecave, J. (2016). The
1066 2016 signal separation evaluation campaign. In *International Conference on Latent Variable Analysis*
1067 *and Signal Separation*, pp. 323-332. Springer, Cham.
1068
1069 MacKay, D., Wulf, G., Yin, C., & Abrams, L. (1993). Relations between word perception and
1070 production: new theory and data on the verbal transformation effect. *Journal of Memory and*
1071 *Language*, *32*, 624-646.
1072
1073 Margulis, E., Simchy-Gross, R., & Black, J. (2015). Pronunciation difficulty, temporal regularity, and
1074 the speech-to-song illusion. *Frontiers in Psychology*, *6*, 48.
1075
1076 Mas-Herrero, E., Marco-Pallares, J., Lorenzo-Seva, U., Zatorre, R., & Rodriguez-Fornells, A. (2013).
1077 Individual differences in music reward experiences. *Music Perception*, *31*, 118-138.
1078
1079 McPherson, M., & McDermott, J. (2018). Diversity in pitch perception revealed by task dependence.
1080 *Nature Human Behaviour*, *2*, 52-66.

1081

1082 McPherson, M., & McDermott, J. (2020). Time-dependent discrimination advantages for harmonic
1083 sounds suggest efficient coding for memory. *PNAS*, *117*, 32169-32180.

1084

1085 Müllensiefen, D., Gingras, B., Musil, J., & Stewart, L. (2014). The musicality of non-musicians: an
1086 index for assessing musical sophistication in the general population. *PLoS one*, *9*(2), e89642.

1087

1088 Mullin, H., Norkey, E., Kodwani, A., Vitevitch, M., & Castro, N. (2021). Does age affect perception of
1089 the Speech-to-Song illusion? *PLoS one*, *16*, e0250042.

1090

1091 Niarchou, M., Gustavson, D., Sathirapongsasuti, J., Anglada-Tort, M., Eising, E., Bell, E., McArthur,
1092 E., Straub, P., The 23andMe Research Team, McAuley, J., Capra, J., Ullén, F., Creanza, N., Mosing,
1093 M., Hinds, D., Davis, L., Jacoby, N., & Gordon, R. (2021). Unravelling the genetic architecture of
1094 musical rhythm: a large-scale genome-wide association study of beat synchronization.

1095 *bioRxiv*, 836197. doi:10.1101/8361

1096

1097 Nolan, F. (2003). Intonational equivalence: an experimental evaluation of pitch scales. In *Proceedings*
1098 *of the 15th International Congress of Phonetic Sciences, Barcelona* (Vol. 39).

1099

1100 Norman-Haignere, S., Kanwisher, N. G., & McDermott, J. H. (2015). Distinct cortical pathways for
1101 music and speech revealed by hypothesis-free voxel decomposition. *Neuron*, *88*(6), 1281-1296.

1102

1103 Ogg, M., Slevc, R., & Idsardi, W. (2017). The time course of sound category identification: insights
1104 from acoustic features. *Journal of the Acoustical Society of America*, *142*, 3459-3473.

1105

1106 Ogg, M., Moraczewski, D., Kuchinsky, S., & Slevc, R. (2019) Separable neural representations of
1107 sound sources: speaker identity and musical timbre. *NeuroImage*, *191*, 116-126.
1108
1109 Ogg, M., Carlson, T., & Slevc, R. (2020) The rapid emergence of auditory object representations in
1110 cortex reflect central acoustic attributes. *Journal of Cognitive Neuroscience*, *32*, 111-123.
1111
1112 Patel, A. (2008). *Music, language, and the brain*. Oxford University Press.
1113
1114 Patterson, R., Uppenkamp, S., Johnsrude, I., & Griffiths, T. (2002). The processing of temporal pitch
1115 and melody information in auditory cortex. *Neuron*, *36*, 767-776.
1116
1117 Penagos, H., Melcher, J., & Oxenham, A. (2004). A neural representation of pitch salience in
1118 nonprimary human auditory cortex revealed with functional magnetic resonance imaging. *Journal of*
1119 *Neuroscience*, *24*, 6810-6815.
1120
1121 Pfordresher, P., Mantell, J., & Pruitt, T. (2021). Effects of intention in the imitation of sung and
1122 spoken pitch. *Psychological Research*. doi:10.1007/s00426-021-01527-0
1123
1124 Phillips-Silver, J., & Trainor, L. J. (2005). Feeling the beat: Movement influences infant rhythm
1125 perception. *Science*, *308*, 1430–1430
1126
1127 Pilotti, M., & Khurshid, A. (2004). Semantic satiation effect in young and older adults. *Perceptual and*
1128 *Motor Skills*, *98*, 999-1016.
1129
1130 Pilotti, M., Simcox, T., Baldy, J., & Schauss, F. (2011). Are verbal transformation sensitive to age
1131 differences and stimulus properties? *The American Journal of Psychology*, *124*, 89-97.

1132

1133 Pressnitzer, D., Graves, J., Chambers, C., de Gardelle, V., & Egré, P. (2018). Auditory perception:

1134 Laurel and Yanny together at last. *Current Biology*, 28, R737-R759.

1135

1136 Prince, J. (2014). Pitch structure, but not selective attention, affects accent weightings in metrical

1137 grouping. *Journal of Experimental Psychology: Human Perception and Performance*, 40, 2073-2090.

1138

1139 Rathcke, T., Falk, S., & Dalla Bella, S. (in press). Music to your ears: sentence sonority and listener

1140 background modulate the "speech-to-song illusion". *Music Perception*.

1141

1142 Rowland, J., Kasdan, A., & Poeppel, D. (2019). There is music in repetition: looped segments of

1143 speech and nonspeech induce the perception of music in a time-dependent manner. *Psychonomic*

1144 *Bulletin & Review*, 26, 583-590.

1145

1146 Russo, F., Vuvar, D., & Thompson, W. (2019). Vowel content influences relative pitch perception in

1147 vocal melodies. *Music Perception*, 37, 57-65.

1148

1149 Sandstrom, G. M., & Russo, F. A. (2013). Absorption in music: Development of a scale to identify

1150 individuals with strong emotional responses to music. *Psychology of Music*, 41(2), 216-228.

1151

1152 Savage, P., Brown, S., Sakai, E., & Currie, T. (2015). Statistical universals reveal the structures and

1153 functions of human music. *PNAS*, 112, 8987-8992.

1154

1155 Scheirer, E., & Slaney, M. (1997). Construction and evaluation of a robust multifeature speech/music

1156 discriminator. In *1997 IEEE International Conference in Acoustics, Speech, and Signal Processing* (pp.

1157 1331– 1334).

1158

1159 Schneider, P., Sluming, V., Roberts, N., Scherg, M., Goebel, R., Specht, H., Dosch, H., Bleeck, S.,
1160 Stippich, C., & Rupp, A. (2005). Structural and functional asymmetry of lateral Heschl's gyrus reflects
1161 pitch perception preference. *Nature Neuroscience*, *8*, 1241-1247.

1162

1163 Simchy-Gross, R., & Margulis, E. (2018). The sound-to-music illusion: repetition can musicalize
1164 nonspeech sounds. *Music and Science*, *1*, 1-6.

1165

1166 Smith, L. (1984). Semantic satiation affects category membership decision time but not lexical
1167 priming. *Memory and Cognition*, *12*, 483-488.

1168

1169 Strand, J., Cooperman, A., Rowe, J., & Simenstad, A. (2014). Individual differences in susceptibility to
1170 the McGurk Effect: links with lipreading and detecting audiovisual incongruity. *JSLHR*, *57*, 2322-2331.

1171

1172 Swaminathan, S., Schellenberg, G., & Khalil, S. (2017). Revisiting the association between music
1173 lessons and intelligence: training effects or music aptitude? *Intelligence*, *62*, 119-124.

1174

1175 Swaminathan, S., & Schellenberg, G. (2019). Musical ability, music training, and language ability in
1176 childhood. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *46*, 2340-2348.

1177

1178 Swaminathan, S., Kragness, H., & Schellenberg, G. (2021). The Musical Ear Test: Norms and
1179 correlates from a large sample of Canadian undergraduates. *Behavior Research Methods*.

1180 doi:10.3758/s13428-020-01528-8

1181

1182 Tierney, A., Dick, F., Deutsch, D., & Sereno, M. (2013). Speech versus song: multiple pitch-sensitive
1183 areas revealed by a naturally occurring musical illusion. *Cerebral Cortex*, *23*, 240-254.

1184

1185 Tierney, A., Patel, A., & Breen, M. (2018a). Repetition enhances the musicality of speech and tone
1186 stimuli to similar degrees. *Music Perception*, 35, 573-578.

1187

1188 Tierney, A., Patel, A., & Breen, M. (2018b). Acoustic foundations of the speech-to-song illusion.
1189 *Journal of Experimental Psychology: General*, 6, 888-904.

1190

1191 Toscano, J., & McMurray, B. (2010). Cue integration with categories: weighting acoustic cues in
1192 speech using unsupervised learning and distributional statistics. *Cognitive Science*, 34, 434-464.

1193

1194 Tranchant, P., Vuvan, D. T., & Peretz, I. (2016). Keeping the beat: A large sample study of bouncing
1195 and clapping to music. *PloS one*, 11(7), e0160178.

1196

1197 Tranchant, P., Lagrois, M. É., Bellemare, A., Schultz, B. G., & Peretz, I. (2021). Co-occurrence of
1198 Deficits in Beat Perception and Synchronization Supports Implication of Motor System in Beat
1199 Perception. *Music & Science*, 4, 2059204321991713.

1200

1201 Vanden Bosch der Nederlanden, C., Hannon, E., & Snyder, J. (2015a). Finding the music of speech:
1202 musical knowledge influences pitch processing in speech. *Cognition*, 143, 135-140.

1203

1204 Vanden Bosch der Nederlanden, C., Hannon, E., & Snyder, J. (2015b). Everyday musical training is
1205 sufficient to perceive the speech-to-song illusion. *Journal of Experimental Psychology: General*, 2,
1206 e43-e49.

1207

1208 Vitevitch, M., Ng, J., Hatley, E., & Castro, N. (2020). Phonological but not semantic influences on the
1209 speech-to-song illusion. *Quarterly Journal of Experimental Psychology*. doi: 10.1177/1747021820969

1210

1211 Warren, R., & Gregory, R. (1958). An auditory analogue of the visual reversible figure. *The American*
1212 *Journal of Psychology*, *71*, 612-613.

1213

1214 Warren, R., & Warren, R. (1966). A comparison of speech perception in childhood, maturity, and old
1215 age by means of the verbal transformation effect. *Journal of Verbal Learning and Verbal Behavior*, *5*,
1216 142-146.

1217

1218 Warrier, C., & Zatorre, R. (2002). Influence of tonal context and timbral variation on perception of
1219 pitch. *Perception & Psychophysics*, *64*, 198-207.

1220

1221 Wesseldijk, L., Gordon, R., Mosing, M., & Ullén, F. (2021). Music and verbal ability—a twin study of
1222 genetic and environmental associations. *Psychology of Aesthetics, Creativity, and the Arts*.
1223 doi:10.1037/aca0000401

1224

1225 Wesseldijk, L., Mosing, M., & Ullén, F. (2021). Why is an early start of training related to musical skills
1226 in adulthood? A genetically informative study. *Psychological Science*, *32*, 3-13.

1227

1228 Zatorre, R., & Baum, S. (2012). Musical melody and speech intonation: singing a different tune? *PLoS*
1229 *Biology*, *10*, e1001372.

1230

1231 Zentner, M., & Eerola, T. (2010). Rhythmic engagement with music in infancy. *PNAS*, *107*, 5768-5773.

1232

1233 Zhang, J., Susino, M., McPherson, G., & Schubert, E. (2020). The definition of a musician in music
1234 psychology: a literature review and the six-year rule. *Psychology of Music*, *48*, 389-409.

1235

1236 Zuk, N. J., Teoh, E. S., & Lalor, E. C. (2020). EEG-based classification of natural sounds reveals
1237 specialized responses to speech and music. *NeuroImage*, 210, 116558.

1238