

The following publication Zhang, G., Shao, J., Zhang, C., & Wang, L. (2022). The Perception of Lexical Tone and Intonation in Whispered Speech by Mandarin-Speaking Congenital Amusics. *Journal of Speech, Language, and Hearing Research*, 65(4), 1331-1348 is available at https://dx.doi.org/10.1044/2021_JSLHR-21-00345.
Journal of Speech, Language, and Hearing Research is available at <https://pubs.asha.org/toc/jslhr/65/4>.

1 **The Perception of Lexical Tone and Intonation in Whispered Speech by** 2 **Mandarin-speaking Congenital Amusics**

3 Gaoyuan Zhang¹, Jing Shao², Caicai Zhang^{3*}, Lan Wang⁴

4
5 ¹Department of Chinese Language and Literature, Peking University

6 ²Department of English Language and Literature, Hong Kong Baptist University, Hong Kong,
7 SAR, China

8 ³Research Centre for Language, Cognition, and Neuroscience, Department of Chinese and
9 Bilingual Studies, The Hong Kong Polytechnic University, Hong Kong SAR, China

10 ⁴Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China

11

12 *Corresponding author:

13 Caicai Zhang: Room EF741, Department of Chinese and Bilingual Studies, The Hong
14 Kong Polytechnic University, 11 Yuk Choi Rd, Hung Hom, Hong Kong SAR, China.

15 Tel: (+852) 34008465. Email address: caicai.zhang@polyu.edu.hk.

16

17 Running head: Whispered speech perception in amusia

18

19 This work was supported by grants from the National Natural Science Foundation of China
20 (NSFC: 11504400), the Research Grants Council of Hong Kong (ECS: 25603916), the PolyU
21 Start-up Fund for New Recruits, and Shenzhen Fundamental Research Program
22 (JCYJ20160429184226930; JCYJ20170413161611534).

23

24 **Abstract**

25 **Purpose:** A fundamental feature of human speech is variation, including the manner of
26 phonation, as exemplified in the case of whispered speech. In the current study, we employed
27 whispered speech to examine an unresolved issue about congenital amusia, a neurodevelopmental
28 disorder of musical pitch processing, which also affects speech pitch processing like lexical tone
29 and intonation perception. The controversy concerns whether amusia is a pitch-processing
30 disorder or can affect speech processing beyond pitch.

31 **Method:** We examined lexical tone and intonation recognition in 19 Mandarin-speaking amusics
32 and 19 matched controls in phonated and whispered speech, where fundamental frequency (F0)
33 information is either present or absent.

34 **Results:** The results revealed that the performance of congenital amusics was inferior to that of
35 controls in lexical tone identification in both phonated and whispered speech. These impairments
36 were also detected in identifying intonation (statements/questions) in phonated and whispered
37 modes. Across the experiments, regression models revealed that F0 and non-F0 (duration,
38 intensity and formant frequency) acoustic cues predicted tone and intonation recognition in
39 phonated speech, whereas non-F0 cues predicted tone and intonation recognition in whispered
40 speech. There were significant differences between amusics and controls in the use of both F0
41 and non-F0 cues.

42 **Conclusions:** The results provided the first evidence that the impairments of amusics in lexical
43 tone and intonation identification prevail into whispered speech, and support the hypothesis that
44 the deficits of amusia extend beyond pitch processing.

45

46 **Keywords:** congenital amusia, lexical tone perception, intonation perception, whispered speech,
47 Mandarin Chinese.

48

49 **Introduction**

50 **Deficits of congenital amusia**

51 Congenital amusia (amusia hereafter) is an innate neurodevelopmental disorder that affects fine-
52 grained pitch (pitch is the perceptual correlate of F0) processing throughout the lifetime (Hyde &
53 Peretz, 2003; Peretz et al., 2002, 2003, 2007). It is believed that individuals with amusia (amusics
54 hereafter) have difficulties in detecting out-of-tune melodies and singing in tune although they
55 have normal hearing, intelligence, sufficient exposure to music and no brain injury (Ayotte et al.,
56 2002; Mignault Goulet et al., 2012; Peretz et al., 2002). It has been reported that amusics are not
57 only impaired in musical pitch processing, but also in the processing of several dimensions in
58 speech that rely on pitch, such as the perception of lexical tone, intonation (statement/question),
59 and emotional prosody (Cheung et al., 2021; Huang et al., 2015a; Liu et al., 2010; Nan et al.,
60 2010; Patel et al., 2008; Thompson et al., 2012) These results indicate that the pitch-processing
61 deficit in amusia is not restricted to music, but transfers to the language domain (Vuvan et al.,
62 2015).

63 However, an unresolved issue is whether amusics' inferior performance is restricted to pitch
64 processing or not, and relatedly whether amusics are impaired in non-pitch processing in speech.
65 Amusia, in its most common form, is thought to reflect impaired pitch processing, and therefore
66 sometimes also called tone-deafness (Cousineau et al., 2015). Previous research has mainly
67 focused on pitch processing (Albouy et al., 2013, 2015; Ayotte et al., 2002; Foxtan, 2004; Peretz
68 et al., 2002, 2005), revealing that amusics have reduced sensitivity to fine-grained pitch
69 differences (Foxtan, 2004; Liu et al., 2012), and inferior pitch memory (Gosselin et al., 2009;
70 Tillmann et al., 2009; Williamson & Stewart, 2010). However, it is unclear whether amusics are
71 deficient at processing other, non-pitch auditory cues, and if yes what cues are affected. A few
72 studies have reported that amusics might have inferior durational, frequency, or amplitude
73 processing abilities beyond pitch processing (Jones, Zalewski, et al., 2009; Lehmann et al., 2015;

74 Peretz & Vuvar, 2017; Phillips-Silver et al., 2011; Whiteford & Oxenham, 2017). Jones,
75 Zalewski et al. (2009) reported that compared with the musically intact control group, the amusic
76 group had poor performance in pitch, duration pattern discrimination and auditory gap detection
77 tasks. Phillips-Silver et al. (2011) proposed a new form of congenial amusia, namely beat-
78 deafness, based a group of individuals who had deficiencies in perceiving and producing the
79 musical beat. Peretz & Vuvar, (2017) divided individuals into pitch-based amusics and time-
80 based amusics (i.e., beat-deaf amusics) based on the criteria of performing below the cutoff score
81 (2 SD below the mean) on two pitch-related tests (Scale and Off-key) and above the cutoff in the
82 Off-beat test for the former, and performing above the cutoff score on the pitch-related tests and
83 below the cutoff in the Off-beat test for the latter. The authors reported a prevalence rate of 1.5%
84 for pitch-based amusics and 3.1% for time-based amusics. An important difference between these
85 two subtypes of amusia is that time-based amusics are often associated with other developmental
86 disorders such as dyscalculia and dyslexia, whereas the pitch-based form of amusia is believed to
87 emerge in isolation, and relatively free of language problems (Peretz & Vuvar, 2017). But two
88 recent studies reported association between pitch processing deficits and dyslexia (Couvignou et
89 al., 2019; Couvignou & Kolinsky, 2021). Whiteford & Oxenham, (2017) reported that amusics
90 were impaired in both low frequency (500Hz and 2000Hz) detection and high-frequency (8000Hz)
91 detection that is beyond the frequency range of musical pitch, indicating that amusia is not a
92 deficit selective to the musical attributes of pitch. They also assessed the participants' sensitivity
93 to frequency modulation detection and amplitude modulation detection in a one-interval yes/no
94 task and a standard two-alternative forced choice task. They found that the ability of amusics to
95 detect frequency modulation and amplitude modulation was significantly poorer than that of
96 controls in both tasks. Although the amplitude detection task did not involve any pitch-related
97 changes, the amusics were still at a disadvantage, which questioned the general assumption that
98 amusia is a pitch-specific deficit.

99 Regarding speech processing, some evidence indicates that amusics' impairments extend
100 beyond pitch processing, affecting phonological awareness, segmental processing, or speech
101 comprehension (Jiang et al., 2012; Jones, Lucker, et al., 2009; Liu et al., 2015; Sun et al., 2017;
102 Zhang et al., 2017). It has been found that amusia may affect phonological processing such as
103 phonological awareness (Jones, Lucker, et al., 2009; Sun et al., 2017), and phonological
104 representations of lexical tones as revealed by categorical perception studies (Huang et al., 2015a;
105 Jiang et al., 2012; Zhang et al., 2017). The results of Zhang et al. (2017) showed that Cantonese-
106 speaking amusics had impairment in the discrimination of vowels in addition to pure tones and
107 lexical tones, but their ability of voice-onset time (VOT) processing remained largely intact. Liu
108 et al. (2015) examined the speech intelligibility of Mandarin-speaking amusics in perceiving
109 Mandarin sentences with natural and flattened F0 curves in quiet and in noise. The authors found
110 that Mandarin-speaking amusics had deficits in speech intelligibility for both natural-F0 and flat-
111 F0 sentences regardless of noise, and their deficit in speech intelligibility was not associated with
112 their pitch perception deficit. These results led the authors to argue that segmental processing
113 might be impaired in amusics, independent of pitch processing. However, a problem with the
114 previous studies is that these additional speech-processing deficits in amusia are likely to stem
115 from their low-level auditory pitch deficit. For instance, a phonological deficit in categorical
116 perception of lexical tones might originate from the low-level pitch-processing deficit in Chinese
117 speakers with amusia from birth. The transfer of a low-level auditory deficit to a higher-level
118 phonological deficit has been reported in other types of developmental disorders such as
119 developmental dyslexia (Goswami et al., 2010, 2011). Even in the case of Liu et al. (2015),
120 where pitch was flattened and neutralized, some pitch cues were still present. To investigate
121 whether amusics are impaired in other areas of speech processing that do not involve pitch, it is
122 thus necessary to employ speech conditions where pitch is absent, like in the case of whispered
123 speech.

124 **Perception of lexical tones in whispered speech**

125 Compared with phonated speech, the most obvious characteristic of whispered speech is that the
126 dominant perceptual cue – fundamental frequency (F0 hereafter) is absent. Previous studies have
127 shown that compared with phonated tones, listeners have difficulty in identifying lexical tones in
128 whispered speech, but the accuracy is above chance level, meaning that lexical tones can still be
129 recognized to some extent (Gao, 2002; Jensen, 1958; Jiao et al., 2015; Jiao & Xu, 2019; Yang et
130 al., 2005). For instance, Gao (2002) examined tone recognition in monosyllables in whispered
131 Mandarin speech, and found that the rank of recognition accuracy for the four tones in the order
132 of increasing difficulty is T3 (low dipping tone), T4 (high falling tone), T1 (high level tone) and
133 T2 (high rising tone). Jiao et al. (2015) showed that the identification rates for phonated tones in
134 Mandarin were: T1 (98.9%), T2 (98.9%), T4 (94.7%) and T3 (94.2%); in contrast, the rates for
135 whispered tones were: T3 (85%), T4 (66.9%), T2 (35.8%) and T1 (21.7%).

136 Most research have found that there in fact is tone perception in whispered speech rather
137 than contextually determined interpretations (Heeren, 2015). Therefore, there must be non-F0
138 cues in the speech signal that can compensate for the absence of F0 in whispered tone recognition
139 (Jiao & Xu, 2019). Some researchers proposed some secondary cues, such as duration (Gao,
140 2002; Li & Guo, 2012; Li & Rong, 2012; Liu & Samuel, 2004; Yang et al., 2005), intensity
141 contour (Gao, 2002; Li & Guo, 2012), and formant frequency (Eklund & Traunmüller, 1997;
142 Higashikawa et al., 1996; Kallail & Emanuel, 1984; Li & Xu, 2005; Matsuda & Kasuya, 1999),
143 that were exaggerated by the speakers to the needs of the listeners. For example, the duration of
144 tones in whispered speech was much longer than that in phonated speech(Li & Guo, 2012; Li &
145 Rong, 2012; Yang et al., 2005), and the duration differences across the four lexical tones were
146 exaggerated in whispered Mandarin(Li & Guo, 2012). With regard to the intensity contour,
147 compared with phonated speech, Gao (2002) found that speakers exaggerated the intensity
148 contour when they were whispering, and concluded that intensity contour played a crucial role in
149 whispered tone perception especially for Tone 3 and Tone 4. Higashikawa et al. (1996) did the
150 formant analysis of the vowel /a/ with low, mid and high pitch in whispered and phonated

151 Japanese and found that the formants were higher overall in whispered speech than in phonated
152 speech. However, there is no consensus on the question of whether non-F0 acoustic are
153 exaggerated or not. Indeed, a few papers reported that non-F0 acoustic cues were not enhanced in
154 whispered utterances. For example, Chang & Yao (2007) found that there were similar
155 differences of duration and average intensity across the four tones in whispered and normal
156 speech. Jiao & Xu (2019) examined duration, intensity, spectral tilt and formants of phonated
157 tones and whispered tones respectively, and showed that there was no special articulatory
158 manoeuver enhancement in whispered tones. Even if the non-F0 acoustic cues were not
159 exaggerated in whispered utterances, they might carry certain acoustic distinctions to support tone
160 identification.

161 To summarize, the majority of previous studies have focused on examining the enhancement
162 of acoustic cues in whispered speech. Although there is some controversy regarding whether non-
163 F0 acoustic cues are enhanced or not, such cues are believed to compensate for the lack of F0 and
164 support tone recognition in whispered speech to some extent. However, few studies have directly
165 examined the relationship between tone recognition performance and the distribution of non-F0
166 acoustic cues across the four tones. Furthermore, it remains unclear whether non-F0 acoustic cues
167 are employed differentially by individuals with and without amusia among native Chinese
168 speakers.

169 **Perception of intonation in whispered speech**

170 It is generally agreed that statements are characterized by a falling intonation, whereas questions
171 are associated with a rising F0 contour, which conveys the meaning of non-finality and inquiry
172 (Ohala, 1983). Studies on different languages have shown that questions are marked by acoustic
173 characteristics such as a final increase in F0, higher F0 level, wider F0 range, longer syllable
174 duration on the final position, and higher intensity level (Gussenhoven & Chen, 2000; Hirschberg
175 & Ward, 1992; Ho, 1977; Ma et al., 2006).

176 As in the case of whispered tone perception, previous studies have shown that compared
177 with phonated speech, listeners have difficulty in identifying boundary tone and intonation in
178 whispered speech, but the accuracy is above chance level, meaning that boundary tone and
179 intonation can still be recognized to some extent (Heeren & Heuven, 2009; Jiao & Xu, 2019).
180 Some studies suggested that there is a compensatory strategy beneficial for perceiving boundary
181 tone and intonation in whispered speech (Heeren, 2015; Heeren & Heuven, 2009). But it is not
182 entirely clear what non-F0 acoustic cues facilitated intonation perception in whispered speech,
183 and whether they are enhanced in the whispered mode (Jiao & Xu, 2019). Heeren & Heuven
184 (2009) found that statements and questions could still be identified well above chance in
185 whispered Dutch, and formant frequency and intensity differences that correlated with high and
186 low boundary tones might be possible perceptual cues for linguistic intonation. Jiao & Xu (2019)
187 investigated the production and perception of questions and statements in whispered Chinese,
188 which showed that compared to whispered statements, whispered questions had flattened spectral
189 slope. Yet they speculated that this acoustic cue did not assist the identification of questions since
190 whispered questions were recognized much less accurately than statements. They concluded that
191 there was no evidence of effective production enhancement for intonation identification in
192 whispered Mandarin.

193 Again, a limitation of previous studies on whispered intonation is that few studies have
194 directly examined the relationship between intonation perception performance and the
195 distribution of non-F0 acoustic cues. Instead, the potential contribution of a certain acoustic cue
196 to perception was inferred indirectly by examining whether there is exaggeration of that cue in
197 production or whether it differs significantly between whispered statements and questions. While
198 these analyses are certainly informative, it is necessary to examine whether the statement-
199 question distance in one or a set of non-F0 acoustic cues can explain the perceptual performance
200 of Mandarin listeners, and whether there are individual differences between listeners with high
201 and low musical aptitude (i.e., musically intact controls and congenital amusics).

202 **The current study**

203 As mentioned above, whispered speech provides an ideal case to examine whether amusics have
204 impairment in other aspects of the linguistic domain other than pitch processing. In addition, it
205 remains unexplored to what extent the distribution of non-F0 acoustic cues in whispered speech
206 can predict listeners' recognition performance of lexical tones and speech intonation respectively.
207 Relatedly, no studies have examined before whether and if so how acoustic cues are differentially
208 employed by listeners with amusia and intact musical abilities in the recognition of phonated and
209 whispered tones and intonation.

210 To address the aforementioned questions, the current study compared the performance of
211 Mandarin-speaking amusics and controls in recognizing lexical tones and speech intonation
212 (statements/questions) in phonated and whispered speech through a series of identification tests.
213 We then examined to what extent the participants' musical pitch, rhythm and memory abilities, as
214 assessed by the Montreal Battery of Evaluation of Amusia (MBEA) (Peretz et al., 2003; Vuvan et
215 al., 2018), can predict their tone and intonation recognition performance in phonated and
216 whispered speech respectively using regression analyses. Regarding the contribution of acoustic
217 cues to tone and intonation recognition in phonated and whispered speech, the following set of
218 acoustic cues were selected and measured based on previous studies(Jiao et al., 2015; Jiao & Xu,
219 2019; Lima et al., 2016; Liu et al., 2012): phonated tone – F0, duration, intensity and formant
220 frequency; whispered tone – duration, intensity and formant frequency; phonated intonation – F0,
221 duration and intensity; whispered intonation – duration and intensity. We then examined which
222 set of acoustic cues can predict the recognition performance of lexical tones and intonation in
223 phonated and whispered speech in amusics and controls respectively using regression analyses.

224 There are two possible predictions. According to the hypothesis that amusia is primarily a
225 pitch-processing disorder and relatively free of language problems(Ayotte et al., 2002; Hyde &
226 Peretz, 2003; Peretz & Vuvan, 2017), the amusics would perform comparably to the controls in
227 the recognition of lexical tones and intonation in whispered speech, where pitch is absent, and yet

228 inferiorly in the case of phonated speech. Furthermore, amusics are less likely to rely on pitch
229 cues and may tend to employ other acoustic cues on which they have no or less severe auditory
230 deficits for lexical tone and intonation recognition in phonated and whispered speech. On the
231 other hand, given the hypothesis that the amusics are impaired in speech processing beyond pitch
232 (Jiang et al., 2012; Jones, Lucker, et al., 2009; Liu et al., 2015; Sun et al., 2017; Zhang et al.,
233 2017), we expect amusics to perform worse than controls even in whispered speech, and the
234 presence of non-F0 acoustic cues may not sufficiently enable them to compensate for their
235 deficient pitch processing. Furthermore, fewer (F0 or non-F0) acoustic cues may explain the
236 recognition performance in amusics.

237

238 **Materials and Methods**

239 **Participants**

240 Nineteen Mandarin-speaking amusics and 19 matched musically intact controls were recruited in
241 this study. All participants were native speakers of Mandarin Chinese from northern China, right-
242 handed and reported no previous history of speech, hearing, neurological or psychiatric
243 impairments. No participants had any formal musical training. All the participants were selected
244 by the MBEA, which consists of six subtests: scale, contour, interval, rhythm, meter and memory
245 (Peretz et al., 2003). Among the six subtests, scale, contour and interval are pitch-based tests that
246 concern melodic organization, rhythm and meter are duration-based tests that concern temporal
247 organization, and memory concerns incidental memory of previous heard melodies. To separate
248 amusics from controls, the pitch composite scores (sum of the number of correct trials in the three
249 pitch-based subtests) were calculated, and those participants who scored at or below 65 were
250 classified as amusics (Liu et al., 2012, Vuvar et al., 2018). Results of independent-samples t-tests
251 confirmed that the amusics performed significantly worse than the controls on the pitch
252 composite scores and the global score (percent of correct responses averaged across the six

253 subtests), as well as on the percent of correct responses of each subtest ($ps < .001$). The
254 demographic characteristics of the participants are shown in Table 1. The experimental
255 procedures were approved by the Human Subjects Ethics committee of the Shenzhen Institutes of
256 Advanced Technology, Chinese Academy of Science. Informed written consent was obtained
257 from the participants in compliance with the experiment protocols.

258 **Stimuli**

259 The speech stimuli were recorded by a 27-year-old female Mandarin speaker and a 28-year-old
260 male Mandarin speaker, who are native speakers of Mandarin Chinese from northern China. All
261 the word (tone) and sentence (intonation) stimuli are provided in supplementary Table 1. We
262 recruited two speakers to record the speech stimuli in order to avoid the influence of
263 idiosyncrasies of individual speakers, for the reason that the stimuli produced by some speakers
264 might be easier to identify than those by other speakers according to a previous study (Gao,
265 2002). The two speakers were asked to produce the isolated words and sentences both in
266 whispered and phonated mode as naturally as possible. The sound recording was done in a
267 soundproof room using Praat (Boersma & Weenink, 2001), with 44.1 kHz sampling rate. We
268 used ProsodyPro (Xu, 2013) to extract the acoustic parameters (F0, duration and intensity) of the
269 stimuli. The measurements of F1 and F2 in both phonated and whispered speech followed the
270 method described in (Sharifzadeh et al., 2012). The F1 and F2 were measured from the average of
271 the steady portion where the formants were relatively clear and steady by a trained phonetician.

272 To assess tone identification in Mandarin Chinese, 36 words with nine base (C)V syllables
273 (/ta/, /ti/, /tu/, /pa/, /pi/, /tʃu/, /a/, /i/, /u/) contrasting the four lexical tones were selected as the
274 stimuli (supplementary Table 1). Each word was produced three times by each speaker under
275 phonated and whispered conditions. Overall, there were a total of 432 tokens (9 syllables \times 4
276 tones \times 2 phonation modes \times 3 repetitions \times 2 speakers). From the three repetitions, one token
277 with accurate and clear pronunciation was selected and used as the stimuli in the ensuing

278 perception tasks, totalling 144 tokens. Figure 1 shows the time-normalized F0 contours of the
279 four Mandarin tones from the two speakers, averaged across all the syllables which contained the
280 same tones in the phonated stimulus set. We measured F0 mean, F0 SD, duration, mean intensity,
281 F1 and F2 of the phonated stimuli produced by the two speakers. A series of one-way ANOVAs
282 with the factor of *lexical tone* were conducted on each acoustic cue, with the *p*-value corrected for
283 multiple comparisons ($.05/6 = .008$). There were significant differences between the four
284 phonated tones in F0 mean, F0 SD, duration, and intensity ($ps \leq .05$), but no significant effect in
285 F1 or F2 (see supplementary Table 2 for the acoustic cues of the four tones and statistical results).
286 The post hoc analyses found that every tone was significantly different from each other in F0
287 mean ($T1 > T4 > T2 > T3$, $ps < .001$) and F0 SD ($T4 > T2 > T3 > T1$, $ps \leq .003$); for duration, the
288 results fell into the pattern of $T3 > T1 \approx T2 > T4$ ($ps \leq 0.003$); for intensity, T3 was significantly
289 lower than T1 and T4 ($ps \leq .03$). We also measured duration, mean intensity, F1 and F2 of the
290 whispered stimuli produced by the two speakers. The results of one-way ANOVA with the factor
291 of *lexical tone* with correction for multiple comparisons ($.05/4 = .0125$) revealed that there was a
292 significant difference between the four whispered tones in duration ($p < .001$), but not in intensity,
293 F1 or F2 (see supplementary Table 3 for the acoustic cues of the four whispered tones and
294 statistical results). The post hoc analyses found that the duration of T3 was significantly longer
295 than the other three tones, while the duration of T4 was significantly shorter than the other three
296 tones ($ps \leq .001$).

297 To assess intonation identification in Mandarin Chinese, 25 statement-question pairs sharing
298 the same words were constructed as the stimuli (see supplementary Table 1 for the sentences).
299 Five sentence lengths were included (4, 5, 6, 7 and 10 syllables), with five sentences for each
300 sentence length. For each length, four out of the five sentences contained words with identical
301 tones on every position (e.g., the sentence ‘张薇开车’ consisted of only T1) and the last sentence
302 contained words with varied tones (e.g., the sentence ‘李刚讲课’ consisted of varied tones),

303 which ensured that all four tones appeared on every position for each sentence length to avoid any
304 potential influence of lexical tones on the intonation F0 patterns (Yan, 2016). Each sentence was
305 produced twice by each speaker under the phonated and whispered conditions. Overall there were
306 a total of 400 tokens (5 sentences \times 2 intonations \times 5 length types \times 2 phonation modes \times 2
307 repetitions \times 2 speakers). From the two repetitions, one token with accurate and clear
308 pronunciation was selected and used as stimuli in the ensuing perception tasks, totalling 200
309 tokens. Figure 2 displays the real-time F0 contours of one pair of statement and question
310 produced by the male speaker. As can be seen, the differences between statements and questions
311 were not only present in the F0, but also in the total sentence duration. For phonated speech, we
312 tested the acoustic characteristics (F0, duration and intensity) of the whole sentences and their
313 final syllables produced by the two speakers which followed the method in Liu et al. (2012) and
314 Lima et al. (2016). Paired-samples *t*-tests with the factor of *intonation* (corrected *p*-value at .05/9
315 = .006) indicated that significant intonation differences were detected on all acoustic cues of the
316 whole sentences, and on the F0 mean and intensity of the final syllables (*p*s < .001) (see
317 supplementary Table 4 for the acoustic cues of statements and questions and statistical results).
318 For whispered speech, we measured the acoustic cues (mean intensity and duration) of the whole
319 sentences and their final syllables produced by the two speakers, which echoes acoustic
320 characteristics in the phonated mode. Paired-samples *t*-tests with the factor of *intonation*
321 (corrected *p*-value at .05/4 = .0125) indicated that statements had significantly longer overall
322 duration than questions (*p* < .001), but significantly shorter duration on the final syllable (*p* =
323 .001) (see supplementary Table 5 for the acoustic cues of statements and questions and statistical
324 results).
325

326 **Procedure**

327 The study included two tasks: lexical tone identification and intonation identification, both of
328 which were implemented using E-prime 2.0. In both tasks, the two phonation modes (phonated
329 and whispered speech) were presented in separate blocks. The stimuli from the two speakers were
330 also presented in two separate sub-blocks in order to avoid the effect of talker variation. The
331 order of these two tasks (the lexical tone identification task and the intonation identification task)
332 was counterbalanced across the participants. Furthermore, half of the participants completed the
333 phonated speech block first, and the other half completed the whispered speech block first. The
334 presentation order of the two sub-blocks (two speakers) within each task was also
335 counterbalanced across the participants.

336 In the lexical tone identification task, the stimuli were presented three times, resulting in a
337 total of 216 trials ($9 \text{ syllables} \times 4 \text{ tones} \times 3 \text{ repetitions} \times 2 \text{ speakers}$) in each of the two phonation
338 modes. Within each sub-block, all the trials were presented randomly. In each trial, a fixation
339 occurred at first for 500ms, followed by the presentation of a spoken stimulus via the headphones.
340 The participants were asked to identify the tone of the stimulus by pressing buttons 1-4 referring
341 to the four lexical tones in Mandarin on a computer keyboard. The experiment only proceeded to
342 the next trial when a response was received. Participants were instructed to respond as quickly as
343 possible. There were practices before each task to familiarize the participants with the
344 experimental procedure. The practice contained the three syllables /a, i, u/ with the four tones in
345 both phonation modes presented only once in random order. Half of the participants practiced on
346 the stimuli produced by the female speaker in the phonation mode and the stimuli produced by
347 the male speaker in the whispered mode; this was reversed in the other half of the participants.

348 As for the intonation identification task, the stimuli were presented twice, resulting in a total
349 of 200 trials ($5 \text{ sentences} \times 2 \text{ intonations} \times 5 \text{ length types} \times 2 \text{ repetition} \times 2 \text{ speakers}$) in each of
350 the two phonation modes. The buttons “Q” and “S” which refer to questions and statements

351 respectively were response buttons on a computer keyboard. The other aspects were the same as
352 those in the tone identification task. There were also practices before each task to familiarize the
353 participants with the experimental procedure. The practice contained ten trials comprising the 10-
354 syllable sentences with the two intonation patterns in both phonation modes produced by the two
355 speakers, which were randomly presented once.

356

357 **Data analysis**

358 For the tone identification task, accuracy was recorded and analysed. Accuracy was the
359 percentage of trials correctly identified for each tone per subject. For the intonation identification
360 task, performance was scored as the sensitivity index d' (Macmillan & Creelman, 2005). We used
361 the sensitivity index d' to analyse the intonation identification data for the reason that Jiao & Xu
362 (2019) indicated that statement is likely to be treated as a default choice when identification was
363 challenging and as a result the signal detection method allows us to avoid any response bias
364 (Irwin et al., 1992). The d' was computed as the z-score of the hit rate minus that of the false
365 alarm rate for each phonation mode per subject. Specifically, the hit rate was the rate of “question”
366 responses to the “question” test items, while the false alarm rate was the rate of “question”
367 responses to the “statement” test items. *Group* \times *lexical tone* \times *phonation mode* repeated-
368 measures ANOVA was conducted on the accuracy of the tone identification task. *Group* \times
369 *phonation mode* repeated-measures ANOVA was conducted on the d' of the intonation
370 identification task.

371 Two sets of regression analyses were conducted to examine (1) to what extent the
372 participants’ musical pitch, duration and memory abilities can explain their recognition
373 performance in phonated and whispered speech respectively, and (2) what acoustic cues can
374 explain the listeners’ recognition performance in phonated and whispered speech respectively,
375 and whether different cues were employed by amusics and controls. For the first set of analyses,

376 multiple linear regression models were constructed on the average accuracy of tone identification
377 (averaged across four tones) for the phonated and whispered mode separately, collapsing amusics
378 and controls, with melodic organization, temporal organization and melodic memory as three
379 predictors. Similar regression analysis was conducted on the d' score of intonation identification
380 for the phonated and whispered mode separately. To keep the set of predictors small and to avoid
381 collinearity, we combined the six MBEA subtests into three sets: melodic organization (which is
382 the average accuracy of the three subtests: scale, contour and interval), temporal organization
383 (which is the average accuracy of the two subtests: rhyme and meter), and melodic memory
384 (Peretz et al., 2003). Prior to the regression models, we conducted bivariate Pearson correlations
385 (two-tailed) to estimate the degree of association between the three sets of musical abilities and
386 tone/intonation identification performance, and between the three sets of musical abilities
387 themselves. Only musical abilities that showed significant correlation with tone/intonation
388 identification performance were then included in the linear regression models in a stepwise
389 manner to determine their relative contribution¹.

390 For the second set of analysis, we employed logistic regression to examine the relationship
391 between the acoustic cues for tone and intonation types and the two groups' responses separately
392 for the phonated and whispered mode. Multinomial logistic regression analyses were conducted
393 for tones and binominal logistic regression analyses were conducted for intonations. For the
394 phonated mode, the following acoustic cues were included: F0 mean, F0 SD, duration, intensity,
395 F1 and F2. For the whispered mode, the following acoustic cues were included: duration,
396 intensity, F1 and F2. Accuracy instead of d' was used as the dependent variable in the regression
397 analyses on intonation perception because the difference of statements and questions was
398 collapsed in the d' score. For each regression model, we ensured that the VIF value was below 5
399 (Zhang & Dong, 2004) to avoid collinearity. In the phonated mode, we excluded the F0 mean and
400 intensity of the final syllables since they were highly correlated with the F0 mean and intensity of

401 the whole sentences respectively ($r > 0.8$), in order to reduce the VIF to be below 5. The details
402 of the correlation analyses are presented in supplementary Table 6.

403

404 **Results**

405 **Lexical tone identification task**

406 Figure 3 shows the tone identification accuracy under phonated and whispered conditions for the
407 two groups of participants. There was a significant main effect of *group* ($F(1, 36) = 8.79, p =$
408 $.005, \eta_{\text{partial}}^2 = 0.20$), where the score of the control group ($M = 0.81, SD = 0.016$) was
409 significantly higher than that of the amusic group ($M = 0.75, SD = 0.016$). The *group* factor did
410 not interact with the other two factors (*lexical tone* and *phonation mode*). Significant main effects
411 of *phonation mode* ($F(1, 36) = 308.68, p < .001, \eta_{\text{partial}}^2 = 0.90$), *lexical tone* ($F(2.49, 89.93) =$
412 $58.09, p < .001, \eta_{\text{partial}}^2 = 0.62$), and a significant interaction between *lexical tone* and
413 *phonation mode* ($F(2.63, 94.57) = 67.98, p < .001, \eta_{\text{partial}}^2 = 0.65$) were also detected. To
414 explore the two-way interaction, we first conducted independent-samples t-tests to examine the
415 effect of *phonation mode* within each lexical tone. A significant effect of *phonation mode* was
416 observed in all lexical tones ($ps \leq .01$), where the accuracy of the phonated mode was always
417 higher than that of the whispered mode. Then one-way ANOVAs with the factor of *lexical tone*
418 within each phonated mode were conducted, revealing a significant effect of *lexical tone* in both
419 phonation modes ($ps \leq .02$). The post hoc analyses revealed that the rank of identification
420 accuracy of the four lexical tones were different in the two phonation modes. For the phonated
421 mode, the accuracy rank of the four tones was $T4 \approx T1 \approx T3 > T2$, where the accuracy of T2 was
422 significantly lower than the other three tones ($ps < .05$). Nonetheless, the accuracy rank of the
423 four tones was $T3 > T4 > T1 \approx T2$ in the whispered mode. The accuracy for T3 was significantly
424 higher than the other three tones ($ps < .001$) and the accuracy of T4 was also significantly higher
425 than T2 and T1 ($ps < .001$) in the whispered mode.

426 Confusion matrixes across the tones were constructed for each phonation mode. The details
427 of confusion matrixes are presented in supplementary Table 7 and 8. All the tones were
428 recognized above chance level. In the phonated mode, the confusion pattern of controls was
429 roughly similar to that of amusics, that is, T1 and T2 were to some extent confused with each
430 other, and T3 was more often confused with T2, whereas T4 exhibited no clear confusion bias. In
431 the whispered mode, the confusion pattern differed from that of the phonated mode, but controls
432 and amusics exhibited roughly similar confusion patterns except for T1. For both controls and
433 amusics, they were likely to confuse T2 with T3, and to a less extent also confuse T3 with T2,
434 and T4 was more often confused with T1.

435

436 **Intonation identification task**

437 Figure 4 shows the d' of intonation identification under phonated and whispered conditions for
438 the two groups. There was a significant main effect of *group* ($F(1, 36) = 13.11, p = .001,$
439 $\eta_{partial}^2 = 0.27$), where the d' score of the control group ($M = 1.91, SD = 1.1$) was significantly
440 higher than that of the amusic group ($M = 1.48, SD = 1.05$). There was a significant main effect of
441 *phonation mode* ($F(1, 36) = 855.84, p < .001, \eta_{partial}^2 = 0.96$), where the d' score of the
442 phonated mode ($M = 2.67, SD = 0.53$) was significantly higher than that of the whispered mode
443 ($M = 0.72, SD = 0.39$). The interaction between *group* and *phonation mode* was not significant.

444 As for the confusion matrix in the phonated mode (see supplementary Table 9), the controls
445 showed comparable identification accuracy for statements and questions, whereas there was a
446 clear decline in the identification accuracy for questions in amusics. In whispered utterances (see
447 supplementary Table 9), statements were identified with over 80% accuracy in both groups,
448 whereas questions were recognized at around chance-level in the control group and below

449 chance-level in the amusic group. These observations further confirm that statement is likely to
450 be treated as a default choice by listeners (Jiao & Xu, 2019).

451

452 **Regression analyses**

453 Two sets of regression analyses were conducted. The first linear analysis aimed to test whether
454 the participants' musical pitch (melodic organization), rhythm (temporal organization), and
455 melodic memory abilities can account for their performance in the two phonation modes. The
456 bivariate correlations showed that all three musical abilities were significantly correlated with the
457 participants' performance on the four identification tests ($ps \leq .02$), and thus were all included in
458 the regression models (see Table 2 and supplementary Figure 1). The stepwise regression models
459 were significant in tone and intonation identification for both phonation modes ($ps \leq .001$) (see
460 Table 3). In the phonated mode, temporal organization was a significant predictor for tone
461 identification, and melodic memory was a significant predictor for intonation identification (ps
462 $< .001$), with increased scores of the two predictors contributing to higher accuracy of tone and
463 intonation identification. The pattern was different in the whispered mode. Melodic memory was
464 a significant predictor for both whispered tone and intonation identification ($ps \leq .001$), with
465 increased scores of melodic memory contributing to higher accuracy of whispered tone and
466 intonation identification.

467 The second analysis concerns whether the acoustic cues can account for the participants'
468 responses, and whether amusics and controls employed these acoustic cues in a different manner
469 in their perception. We conducted multinomial logistical regression for tone identification and
470 binominal regression for intonation identification. Tables 4-7 show the main findings. The figures
471 are displayed in supplementary materials (supplementary Figure 2-5). For phonated tone (Table 4)
472 and intonation (Table 5), the models were significant in both cases ($ps < .001$). The Nagelkerke
473 R^2 of the estimate of tone and intonation identification was 0.91 and 0.81 respectively. Note that

474 in the tone model, the second tone in each pair was used as the baseline for the contrast (e.g., T1
475 was the baseline in the T2 vs. T1 pair), and controls were used as the baseline for comparison
476 with amusics. We found that F0 mean, F0 SD, duration, intensity, F1 and F2 significantly
477 predicted the identification of almost all tone pairs ($ps \leq .003$). Importantly, non-F0 cues were
478 also significant predictors, which indicated that F0 is not the only acoustic cue that can
479 differentiate the four tones. The group factor was significant only in the tone pairs including T1
480 ($ps \leq .04$). There were multiple significant interactions between group and acoustic cues, but the
481 specific patterns varied across the tone pairs. For instance, the interaction between group and F0
482 mean was significant in all tone pairs except for T4-T1 ($ps \leq .003$). Note that the F0 mean of T2
483 (vs. T1) and T3 (vs. T1 and T2) is *lower* than the baseline tone in each of these contrasts (see
484 supplementary Table 2). The *positive* coefficients in these contrasts suggested that an increase in
485 F0 mean was *more likely* to lead to the identification of T2 and T3 (i.e., the *wrong* tones) in these
486 contrasts by the amusics compared to the controls, which implied that the amusics may have used
487 F0 mean less efficiently in the identification of these tones. Similarly, the F0 mean of T4 (vs. T2
488 and T3) is *higher* than the baseline tone in each of these contrasts (see supplementary Table 2).
489 The *negative* coefficients in these contrasts suggested that an increase in F0 mean is *less likely* to
490 lead to the identification of T4 (i.e., the *correct* tone) in these contrasts by the amusics, which
491 also implied worse usage of F0 mean in these tone contrasts by the amusics compared to the
492 controls. For easy reference, cases where a certain acoustic cue led to worse or better
493 performance in amusics were marked differentially in Table 4. Overall, the results showed that
494 amusics employed F0 mean, F0 SD, duration and intensity worse than controls, but employed F1
495 and F2 better than controls in various tone contrasts.

496 In the phonated intonation model, we found that sentence F0 mean, sentence F0 SD,
497 sentence F0 direction, sentence intensity, final syllable F0 SD, final syllable duration and group
498 significantly predicted intonation identification accuracy ($ps \leq .01$). Again, several non-F0 cues
499 were significant predictors in the model, which implied that non-F0 acoustic cues are likely to

500 contribute to phonated intonation identification. The interactions between group and sentence F0
501 mean, sentence F0 SD, sentence F0 direction, sentence intensity and final syllable duration were
502 also significant ($ps < .001$). Note that statement was used as the baseline for contrast with
503 question, and controls were used as the baseline for contrast with amusics. The value of question
504 on each of the aforementioned acoustic cue was *greater* than that of statement (see
505 Supplementary Table 4). The negative coefficients in these significant interactions indicate that as
506 the acoustic value increases, it is less likely for the amusics to choose questions (i.e., the correct
507 intonation) over statements compared to the controls.

508 As for the whispered mode, the models were significant for tone (Table 6) and intonation
509 (Table 7) ($ps < .001$). The Nagelkerke R^2 of the estimate of tone and intonation identification was
510 0.38 and 0.09 respectively. In the whispered tone model, we found that duration, intensity, F1 and
511 F2 significantly predicted the identification of almost all tone pairs ($ps \leq .02$). The interactions
512 between group and duration, F1 and F2 were significant in several tone pairs ($ps \leq .04$). Overall,
513 the amusics employed duration worse than the controls in the T2-T1, T4-T1, T3-T2, T4-T3 pairs;
514 they also employed F1 worse in the T3-T2 pair, and F2 worse in the T4-T1 pair. In contrast, the
515 amusics employed F2 better than the controls in the T2-T1 and T3-T1 pairs.

516 In the whispered intonation model, we found that sentence duration and final syllable
517 duration significantly predicted intonation identification ($ps < .001$). The interaction between
518 group and sentence duration was also significant in the model ($p = .02$). Since the sentence
519 duration was longer in statements than in questions, the significant interaction (with positive
520 coefficient) indicates that when the sentence duration increases, it is more likely for the amusics
521 to identify the intonation pattern as question (i.e., the wrong intonation) compared to the controls.

522

523 **Discussion**

524 The present study examined the identification of lexical tone and intonation by Mandarin-
525 speaking amusics and controls in the phonated and whispered modes. In the lexical tone
526 identification task, the results showed that Mandarin-speaking amusics demonstrated an overall
527 lower accuracy compared with the controls regardless of the phonated and whispered modes.
528 Likewise, the amusics performed inferiorly with respect to the controls in terms of the d' scores
529 in the intonation identification task in both modes. These results indicated that the impairment of
530 amusics extends to tone and intonation identification in whispered speech, where the F0 is absent,
531 implying that amusics are likely to be impaired in other aspects of speech processing other than
532 pitch. In the text below, we discussed the results of the current study in relation to the two
533 research questions raised in the introduction: (1) whether amusia is a pitch-processing disorder or
534 whether it affects other aspects of the linguistic domain beyond pitch processing; (2) what
535 acoustic cues other than the F0 can predict the listeners' recognition performance of lexical tones
536 and speech intonation in the phonated and whispered mode respectively, and whether different
537 acoustic cues were employed by amusics and controls.

538 **Impairment of amusics in tone and intonation identification in the phonated and** 539 **whispered modes**

540 In the phonated mode, the accuracy of the amusic group was significantly lower than the control
541 group, for both lexical tone identification and intonation identification. The inferior performance
542 of amusics in tone identification is in line with the results of several studies on tonal language
543 speakers with amusia (Nan et al., 2010; Shao et al., 2019, 2016; Shao & Zhang, 2020; Wang &
544 Peng, 2014; Zhang et al., 2018), which have suggested that amusia is a domain-general disorder
545 rather than being restricted to the musical domain (Douglas & Bilkey, 2007; Patel et al., 2008;
546 Thompson, 2007; Vuvan et al., 2015). Furthermore, we found that the accuracy of most tones was
547 above 90% except for T2 in the amusic group, which echoed with the results in (Nan et al., 2010)
548 that T2 was the most difficult tone to identify for Mandarin-speaking amusics, especially for

549 those with lexical tone agnosia, who were markedly impaired in lexical tone perception (i.e., with
550 scores below 3 SD of the controls' scores in tone identification and discrimination). A plausible
551 explanation is that similar acoustical characteristics shared by T2 and T3 may exacerbate the
552 confusion (Nan et al., 2010), which is evidenced by the high confusion rate of T2 with T3 (and to
553 some extent with T1) in the confusion matrix in the current study (see supplementary Table 7).
554 On the other hand, this result differed from Liu et al. (2012), who did not find differences
555 between the Mandarin-speaking amusics and controls in the tone identification task, but the
556 discrepancy can be attributed to task differences. Since Liu and colleagues asked the participants
557 to recognize the lexical tone stimuli as words using Chinese characters, instead of as tonal
558 categories, it is likely that their task involved less abstract phonological processing, and as a
559 result did not reveal the group difference.

560 With regard to intonation identification in the phonated mode, compared with controls, the
561 significantly lower accuracy of the amusic group is consistent with the results in Jiang et al.
562 (2010) and Liu et al. (2010), but not with Liu et al. (2012). Liu et al. (2012) have found that
563 Mandarin-speaking amusics performed as well as controls on intonation identification in natural
564 speech. Material differences may explain the somewhat different results of these two studies. The
565 average pitch range of statements and questions was 10.49 and 9.59 semitones respectively in the
566 current study (statement-question difference: 0.9 semitones), whereas that of statements and
567 questions was 11.35 and 7.75 semitones in Liu et al. (2012) (statement-question difference: 3.6
568 semitones). The average pitch excursion of the final syllable of statements and questions was 6.88
569 and 5.83 semitones respectively in the current study (statement-question difference: 1.05
570 semitones), whereas that of statements and questions was 6.82 and 3.59 semitones in Liu et al.
571 (2012) (statement-question difference: 3.23 semitones). It is clear from the comparison above that
572 the phonated stimuli contained smaller differences between statements and questions in the
573 current study than those in Liu et al. (2012). Larger differences may enable amusics to distinguish
574 and identify statements and questions, therefore showing comparable performance with controls.

575 In support of this argument, previous studies have revealed that amusics' average threshold for
576 discriminating pitch direction is around two semitones (Foxton, 2004; Liu et al., 2010). Thus
577 naturally produced smaller pitch contrasts in the current study may be more sensitive in revealing
578 the amusics' intonation processing deficit.

579 In the whispered mode, amusics again had worse performance than controls in both lexical
580 tone and intonation identification tasks, which is consistent with the previous findings that the
581 amusics' impairments extend beyond pitch processing, affecting phonological awareness,
582 segmental processing or speech comprehension (Jiang et al., 2012; Jones, Lucker, et al., 2009;
583 Liu et al., 2015; Sun et al., 2017; Zhang et al., 2017). It also complements the previous findings
584 by further providing evidence that amusics are impaired in lexical tone and intonation
585 identification even when the F0 is absent. Altogether, there is convergent evidence for the notion
586 that the deficits of amusia exist outside of pitch processing. That being said, these results do not
587 necessarily negate the hypothesis that amusia is a pitch-processing disorder, because amusics are
588 indeed impaired in lexical tone and intonation identification in the phonated mode. We will return
589 to the discussion of the deficits of amusics in non-pitch processing in speech after discussing the
590 contribution of acoustic cues to lexical tone and intonation identification in phonated and
591 whispered speech first below.

592 **Contribution of acoustic cues to phonated and whispered tone and intonation** 593 **identification**

594 It remains controversial which non-F0 acoustic cues facilitate tone and intonation identification
595 where the F0 is absent. Previous studies have probed this question from the perspective of the
596 enhancement of other acoustic cues (e.g., duration, intensity and formant frequency) in whispered
597 speech compared to phonated speech, but no consensus has been reached on this issue.
598 Furthermore, few previous studies have directly examined the relationship between acoustic cues
599 and the participants' identification performance.

600 The current study filled this gap and generated new results from the regression analyses. We
601 found that duration, intensity, F1 and F2 were significant predictors for whispered tone
602 identification. These results corroborated with the proposals that duration (Jiao & Xu, 2019; Yang
603 et al., 2005), mean intensity (Jiao & Xu, 2019; Li & Guo, 2012), and formant frequency (Heeren
604 & Heuven, 2009; Li & Xu, 2005) facilitate whispered tone recognition. Intriguingly, duration,
605 intensity, F1 and F2 were also significant predictors in *phonated* tone identification, which
606 indicates that F0 is not the only acoustic cue that contributes to tone identification, although it is a
607 dominant one, as evidenced by the drop in tone identification accuracy from the phonated mode
608 to the whispered mode. As for whispered intonation identification, the results of regression
609 analyses indicated that sentence duration and final syllable duration were significant acoustic
610 predictors.

611 Crucially, we found significant interactions between group and several acoustic cues in the
612 regression models on tone and intonation identification in phonated and whispered speech, which
613 indicates that there were relative weaknesses in the usage of various acoustic cues by the amusics.
614 In phonated tone identification, the amusics not only employed F0 cues (F0 mean and SD) less
615 efficiently than the controls, but also exhibited worse usage of duration and intensity in almost all
616 tone contrasts. On the other hand, the amusics appeared to have employed F1 and F2 cues better
617 than the controls in almost all tone contrasts. But it is worth noting that the acoustic distinction in
618 F1 and F2 among the four tones was relatively small and not significant (whereas there were
619 significant tone differences in F0 mean, F0 SD, duration and intensity; see supplementary Table
620 2), which implies that F1 and F2 may not be the most optimal cues to employ in phonated tone
621 identification. In whispered tone identification, where the F0 was absent, non-F0 acoustic cues
622 including duration, intensity, F1 and F2 presumably played a greater role. Here, the amusics not
623 only employed duration cues less efficiently compared to the controls, but also demonstrated
624 worse performance in the usage of F1 and F2 in several tone pairs; they only used F2 more
625 efficiently than the controls in the T2-T1 and T3-T1 pairs. Likewise, in phonated intonation

626 identification, the amusics not only employed F0 cues less efficiently compared to the controls,
627 but also exhibited a weakness in the usage of sentence intensity and final syllable duration cues.
628 Where the F0 cues were absent, duration cues (sentence and final syllable duration) significantly
629 predicted intonation identification, and the amusics continued to show inferior performance in the
630 usage of sentence duration relative to the controls. Taken together, these observations appeared to
631 suggest that in speech signals (phonated or whispered) with rich acoustic redundancies where
632 multiple acoustic cues index a functional contrast (e.g., statement vs. question or the four tones),
633 the amusics may not employ the most optimal acoustic cues for the contrast or use them less
634 efficiently compared to the controls.

635 How to explain the worse performance of amusics in whispered tone and intonation
636 identification in the current study? As the regression analyses revealed weaknesses in the
637 amusics' usage of duration and to some extent formant frequency cues in whispered speech, a
638 most straightforward explanation is that the amusic participants recruited in this study may have
639 impaired duration or formant frequency processing. This explanation is compatible with the
640 various findings that amusics have inferior durational, frequency, or intensity processing abilities
641 beyond pitch processing (Jones et al., 2009; Lehmann et al., 2015; Peretz & Vuvan, 2017;
642 Phillips-Silver et al., 2011; Whiteford & Oxenham, 2017). Future studies should directly examine
643 fine-grained duration and formant frequency processing (e.g., using threshold tasks) together with
644 whispered speech perception in amusics in a single study, so as to further reveal which sub-
645 domain of acoustic processing best explains the performance of amusics in whispered speech
646 perception.

647 An alternative explanation is that the phonological representations of lexical tone and
648 intonation are impaired in amusics. Several previous studies have indicated that Chinese speakers
649 with amusia are impaired in the phonological representation of lexical tones (Huang et al., 2015a;
650 Jiang et al., 2012; Zhang et al., 2017) . For instance, Jiang et al. (2012) examined the performance
651 of amusics and controls in the categorical perception of lexical tone, and found that amusics

652 performed less categorically, exhibiting less between-category benefit than the controls, which
653 suggested that there was a deficit of higher-level phonological processing of lexical tones of the
654 amusic group. According to this view, regardless of whether amusics are deficient in earlier
655 auditory processing of pitch and non-pitch cues, when the acoustic cues are mapped onto the
656 phonological representation of lexical tones in the categorization process, an impairment in
657 amusics is detected, even in the case of whispered speech.

658 Finally, although it is not our primary interest, the finding that temporal organization
659 significantly predicted phonated tone identification, and melodic memory significantly predicted
660 whispered tone, phonated intonation and whispered intonation identification requires an
661 explanation. It is unexpected that melodic organization, which is related to pitch processing, was
662 not a significant predictor of phonated tone or intonation identification, where F0 was a dominant
663 cue. It may be because melodic organization, temporal organization and melodic memory are
664 highly correlated with each other (see Table 2), and temporal organization or melodic memory
665 may be able to explain more unique variances than melodic organization in these models.
666 Another explanation is that the three pitch-based MBEA subtests do not purely assess pitch
667 processing, but also involve musical knowledge. In either case, the finding that temporal
668 organization was a significant predictor of phonated tone identification reinforces the view above
669 that acoustic duration cues contributed to phonated tone identification. The contribution of
670 melodic memory to intonation identification (phonated or whispered) may be because the
671 intonation tasks used long sentence materials, where memory capacities of the melodic patterns
672 are crucial for the identification performance. However, it is not entirely clear why melodic
673 memory also significantly predicted whispered tone identification. That being said, all these
674 results must be replicated in future studies for more rigorous interpretation.

675 To conclude, we found that Mandarin-speaking amusics showed degraded performance of
676 lexical tone and intonation identification in both phonated and whispered modes compared to
677 musically intact listeners. The results indicated that although only around 7% of amusics self-

678 report that they have difficulties in understanding other people’s speech in daily life (Liu et al.,
679 2015), their deficits affected phonated and whispered lexical tone and intonation processing in the
680 laboratory. The results of the current study are consistent with the hypothesis that the impairment
681 of amusia is domain general, rather than limited to the musical domain. Moreover, our findings
682 indicate that the impairment is not confined to pitch processing, but extend to other aspects
683 beyond pitch processing (Jones et al., 2009; Jones, Lucker, et al., 2009; Lehmann et al., 2015; Liu
684 et al., 2015; Whiteford & Oxenham, 2017; Zhang et al., 2017). It is likely that amusia is a
685 syndromic disorder frequently accompanied by deficiencies of other kinds (Jones et al., 2009;
686 Jones, Lucker, et al., 2009). This study is the first to examine whispered speech perception in
687 amusics, and revealed that amusics have impairments in other aspects of the linguistic domain,
688 which sheds further light on the nature of the deficits underlying amusia. These findings also have
689 real-world implications for the diagnosis and treatment of amusia. However, there are some
690 remaining issues to be addressed in future studies. First, future studies with a large sample of
691 amusics should separate them into pitch- and time-based forms of amusia (Peretz & Vuvan, 2017)
692 and further examine if there are subgroup differences in phonated and whispered tone and
693 intonation perception. It should be noted that the MBEA temporal organization subtests are
694 complex tasks that assess more than duration processing. It is recommended that future studies
695 use tasks that probe into duration processing (e.g., duration threshold tasks) to examine whether
696 time-based amusics truly have duration (and intensity) processing deficits and how these
697 problems contribute to their perception in phonated and whispered speech. Second, and related to
698 the first point, future large-scale studies may separate Mandarin-speaking amusics into lexical
699 tone agnostics and those without severe tone perception deficits (Nan et al., 2010; Huang et al.,
700 2015a; Huang et al., 2015b; Nan et al., 2016), and examine if there are subgroup differences in
701 phonated and whispered tone and intonation perception. Future studies should also investigate
702 whether the finding of the current study generalizes to amusic individuals in other tonal languages
703 (e.g., Cantonese) or non-tonal language (e.g., English), and with different tasks (e.g.,

704 discrimination task). It will also be of interest to use event-related potentials (ERPs) to probe
705 passive and active processing of whispered speech in amusics and examine if there are any
706 processing differences between the two listening conditions (Moreau et al., 2003; Zhang & Shao,
707 2018).

708

709 **Acknowledgements**

710 This work was supported by grants from the National Natural Science Foundation of China
711 (NSFC: 11504400; <http://www.nsf.gov.cn/>) and the Research Grants Council of Hong Kong
712 (ECS: 25603916; <https://www.ugc.edu.hk/eng/rgc/>) to CZ. This work was also supported by the
713 grant of the National Key R & D Program of China (2020YFC2004100) to LW. We thank Mr.
714 Yulin Wen for help with data collection. We thank Dr. Xiaocong Chen for help with statistical
715 analysis.

716

717 **Footnote**

718 ¹We conducted simultaneous linear regression and reported the results in the supplementary
719 materials (see supplementary Table10). Nonetheless, the results of simultaneous regression
720 generated some puzzling results. Although the correlations between the scores of melodic
721 organization and identification performance are positive (see Table 2), the standardized
722 coefficients for melodic organization turned out to be negative in the regression models (although
723 not significant). It might be because the correlation among the three MBEA composites are very
724 high ($ps < .001$). That is, the multicollinearity of the independent variables is strong. As a result,
725 we entered the predictors into the regression models in a stepwise manner. The significant
726 predictors in the stepwise regression models were almost identical to those in the simultaneous
727 model, except that melodic memory was a significant predictor in whispered intonation

728 identification in the stepwise model (but not in the simultaneous model). Forward and backward
729 models generated the same results as the stepwise regression models.

730

731 **References**

- 732 Albouy, P., Mattout, J., Bouet, R., Maby, E., Sanchez, G., Aguera, P.-E., Daligault, S., Delpuech,
733 C., Bertrand, O., Caclin, A., & Tillmann, B. (2013). Impaired pitch perception and
734 memory in congenital amusia: The deficit starts in the auditory cortex. *Brain, 136*(5),
735 1639–1661. <https://doi.org/10.1093/brain/awt082>
- 736 Albouy, P., Mattout, J., Sanchez, G., Tillmann, B., & Caclin, A. (2015). Altered retrieval of
737 melodic information in congenital amusia: Insights from dynamic causal modeling of
738 MEG data. *Frontiers in Human Neuroscience, 9*.
739 <https://doi.org/10.3389/fnhum.2015.00020>
- 740 Ayotte, J., Peretz, I., & Hyde, K. (2002). Congenital amusia: A group study of adults afflicted
741 with a music-specific disorder. *Brain, 125*(2), 238–251.
742 <https://doi.org/10.1093/brain/awf028>
- 743 Boersma, P., & Weenink, D. (2001). *Praat: Doing Phonetics by Computer*.
744 <https://uvafon.hum.uva.nl/praat/>
- 745 Chang, C., & Yao, Y. (2007, August 6). *Tone Production in whispered Mandarin*. Proceedings of
746 the 16th international congress of phonetics sciences, Saarbrücken, Germany.
- 747 Cheung, Y. L., Zhang, C., & Zhang, Y. (2021). Emotion processing in congenital amusia: The
748 deficits do not generalize to written emotion words. *Clinical Linguistics & Phonetics,*
749 *35*(2), 101–116. <https://doi.org/10.1080/02699206.2020.1719209>
- 750 Cousineau, M., Oxenham, A. J., & Peretz, I. (2015). Congenital amusia: A cognitive disorder
751 limited to resolved harmonics and with no peripheral basis. *Neuropsychologia, 66*, 293–
752 301. <https://doi.org/10.1016/j.neuropsychologia.2014.11.031>

753 Couvignou, M., Peretz, I., & Ramus, F. (2019). Comorbidity and cognitive overlap between
754 developmental dyslexia and congenital amusia. *Cognitive Neuropsychology*, 1–17.
755 <https://doi.org/10.1080/02643294.2019.1578205>

756 Couvignou, M., & Kolinsky, R. (2021). Comorbidity and cognitive overlap between
757 developmental dyslexia and congenital amusia in children. *Neuropsychologia*, 107811.
758 <https://doi.org/https://doi.org/10.1016/j.neuropsychologia.2021.107811>

759 Douglas, K. M., & Bilkey, D. K. (2007). Amusia is associated with deficits in spatial processing.
760 *Nature Neuroscience*, 10(7), 915–921. <https://doi.org/10.1038/nn1925>

761 Eklund, I., & Traunmüller, H. (1997). Comparative study of male and female whispered and
762 phonated versions of the long vowels of Swedish. *Phonetica*, 54(1), 1–21.
763 <https://doi.org/10.1159/000262207>

764 Foxton, J. M. (2004). Characterization of deficits in pitch perception underlying “tone deafness.”
765 *Brain*, 127(4), 801–810. <https://doi.org/10.1093/brain/awh105>

766 Gao, M. (2002). *Tones in whispered Chinese: Articulatory features and perceptual cues* [Master
767 Thesis]. University of Victoria.

768 Gussenhoven, C., & Chen, A. (2000). *Universal and language-specific effects in the perception of*
769 *question intonation*. 91–94.

770 Heeren, W. F. L. (2015). Vocalic correlates of pitch in whispered versus normal speech. *The*
771 *Journal of the Acoustical Society of America*, 138(6), 3800–3810.
772 <https://doi.org/10.1121/1.4937762>

773 Heeren, W., & Heuven, V. J. V. (2009). *Perception and production of boundary tones in*
774 *whispered Dutch*. 2411–2414.

775 Higashikawa, M., Nakai, K., Sakakura, A., & Takahashi, H. (1996). Perceived pitch of whispered
776 vowels-relationship with formant frequencies: A preliminary study. *Journal of Voice*,
777 10(2), 155–158. [https://doi.org/10.1016/S0892-1997\(96\)80042-7](https://doi.org/10.1016/S0892-1997(96)80042-7)

778 Hirschberg, J., & Ward, G. (1992). The influence of pitch range, duration, amplitude and spectral
779 features on the interpretation of the rise-fall-rise intonation contour in English. *Journal of*
780 *Phonetics*, 20(2), 241–251. [https://doi.org/10.1016/S0095-4470\(19\)30625-4](https://doi.org/10.1016/S0095-4470(19)30625-4)

781 Ho, A. T. (1977). Intonation Variation in a Mandarin Sentence for Three Expressions:
782 Interrogative, Exclamatory and Declarative. *Phonetica*, 34(6), 446–457.
783 <https://doi.org/10.1159/000259916>

784 Huang, W.-T., Liu, C., Dong, Q., & Nan, Y. (2015a). Categorical perception of lexical tones in
785 mandarin-speaking congenital amusics. *Frontiers in Psychology*, 6.
786 <https://doi.org/10.3389/fpsyg.2015.00829>

787 Huang, W.-T., Nan, Y., Dong, Q., & Liu, C. (2015b). Just-noticeable difference of tone pitch
788 contour change for Mandarin congenital amusics. *The Journal of the Acoustical Society*
789 *of America*, 138(1), EL99–EL104. <https://doi.org/10.1121/1.4923268>

790 Hyde, K. L., & Peretz, I. (2003). “Out-of-pitch” but still “in-time.” *Annals of the New York*
791 *Academy of Sciences*, 999(1), 173–176. <https://doi.org/10.1196/annals.1284.023>

792 Irwin, R. J., Hautus, M. J., & Stillman, J. A. (1992). Use of the receiver operating characteristic in
793 the study of taste perception. *Journal of Sensory Studies*, 7(4), 291–314.
794 <https://doi.org/10.1111/j.1745-459X.1992.tb00196.x>

795 Jensen, M. K. (1958). Recognition of word tones in whispered speech. *WORD*, 14(2–3), 187–196.
796 <https://doi.org/10.1080/00437956.1958.11659663>

797 Jiang, C., Hamm, J. P., Lim, V. K., Kirk, I. J., Chen, X., & Yang, Y. (2012). Amusia Results in
798 Abnormal Brain Activity following Inappropriate Intonation during Speech
799 Comprehension. *PLoS ONE*, 7(7), e41411. <https://doi.org/10.1371/journal.pone.0041411>

800 Jiang, C., Hamm, J. P., Lim, V. K., Kirk, I. J., & Yang, Y. (2010). Processing melodic contour
801 and speech intonation in congenital amusics with Mandarin Chinese. *Neuropsychologia*,
802 48(9), 2630–2639. <https://doi.org/10.1016/j.neuropsychologia.2010.05.009>

803 Jiao, L., Ma, Q., Wang, T., & Xu, Y. (2015). *Perceptual Cues of Whispered Tones: Are They*
804 *Really Special?* 5.

805 Jiao, L., & Xu, Y. (2019). Whispered Mandarin has no production-enhanced cues for tone and
806 intonation. *Lingua*, 218, 24–37. <https://doi.org/10.1016/j.lingua.2018.01.004>

807 Jones, J. L., Lucker, J., Zalewski, C., Brewer, C., & Drayna, D. (2009). Phonological processing
808 in adults with deficits in musical pitch recognition. *Journal of Communication Disorders*,
809 42(3), 226–234. <https://doi.org/10.1016/j.jcomdis.2009.01.001>

810 Jones, J. L., Zalewski, C., Brewer, C., Lucker, J., & Drayna, D. (2009). Widespread Auditory
811 Deficits in Tune Deafness. *Ear & Hearing*, 30(1), 63–72.
812 <https://doi.org/10.1097/AUD.0b013e31818ff95e>

813 Kallail, K. J., & Emanuel, F. W. (1984). Formant-Frequency Differences Between Isolated
814 Whispered and Phonated Vowel Samples Produced by Adult Female Subjects. *Journal of*
815 *Speech, Language, and Hearing Research*, 27(2), 245–251.
816 <https://doi.org/10.1044/jshr.2702.251>

817 Lehmann, A., Skoe, E., Moreau, P., Peretz, I., & Kraus, N. (2015). Impairments in musical
818 abilities reflected in the auditory brainstem: Evidence from congenital amusia. *European*
819 *Journal of Neuroscience*, 42(1), 1644–1650. <https://doi.org/10.1111/ejn.12931>

820 Li, B., & Guo, Y. (2012, May 26). *Mandarin Tone Contrast in Whisper*. Third International
821 Symposium on Tonal Aspects of Languages, Nanjing, China.

822 Li, B., & Rong, R. (2012). Tones in whispered Mandarin. *2012 8th International Symposium on*
823 *Chinese Spoken Language Processing*, 422–425.
824 <https://doi.org/10.1109/ISCSLP.2012.6423539>

825 Li, X., & Xu, B. (2005). Formant comparison between whispered and voiced vowels in
826 Mandarin. *Acta Acustica United with Acustica*, 91(6), 1079–1085.

827 Lima, C. F., Brancatisano, O., Fancourt, A., Müllensiefen, D., Scott, S. K., Warren, J. D., &
828 Stewart, L. (2016). Impaired socio-emotional processing in a developmental music
829 disorder. *Scientific Reports*, 6(1), 34911. <https://doi.org/10.1038/srep34911>

830 Liu, F., Jiang, C., Thompson, W. F., Xu, Y., Yang, Y., & Stewart, L. (2012). The Mechanism of
831 Speech Processing in Congenital Amusia: Evidence from Mandarin Speakers. *PLoS*
832 *ONE*, 7(2), e30374. <https://doi.org/10.1371/journal.pone.0030374>

833 Liu, F., Jiang, C., Wang, B., Xu, Y., & Patel, A. D. (2015). A music perception disorder
834 (congenital amusia) influences speech comprehension. *Neuropsychologia*, 66, 111–118.
835 <https://doi.org/10.1016/j.neuropsychologia.2014.11.001>

836 Liu, F., Patel, A. D., Fourcin, A., & Stewart, L. (2010). Intonation processing in congenital
837 amusia: Discrimination, identification and imitation. *Brain*, 133(6), 1682–1693.
838 <https://doi.org/10.1093/brain/awq089>

839 Liu, S., & Samuel, A. G. (2004). Perception of Mandarin Lexical Tones when F0 Information is
840 Neutralized. *Language and Speech*, 47(2), 109–138.
841 <https://doi.org/10.1177/00238309040470020101>

842 Ma, J. K.-Y., Ciocca, V., & Whitehill, T. L. (2006). Effect of intonation on Cantonese lexical
843 tones. *The Journal of the Acoustical Society of America*, 120(6), 3978–3987.
844 <https://doi.org/10.1121/1.2363927>

845 Macmillan, N. A., & Creelman, C. D. (2005). *Detection theory: A user's guide* (2nd ed).
846 Lawrence Erlbaum Associates.

847 Matsuda, M., & Kasuya, H. (1999). *Acoustic nature of the whisper. 1*, 147–140.

848 Mignault Goulet, G., Moreau, P., Robitaille, N., & Peretz, I. (2012). Congenital amusia persists in
849 the developing brain after daily music listening. *PLoS ONE*, 7(5), e36860.
850 <https://doi.org/10.1371/journal.pone.0036860>

851 Moreau, P., Jolicœur, P., & Peretz, I. (2013). Pitch discrimination without awareness in
852 congenital amusia: Evidence from event-related potentials. *Brain and Cognition*, *81*(3),
853 337–344. <https://doi.org/10.1016/j.bandc.2013.01.004>

854 Nan, Y., Sun, Y., & Peretz, I. (2010). Congenital amusia in speakers of a tone language:
855 Association with lexical tone agnosia. *Brain*, *133*(9), 2635–2642.
856 <https://doi.org/10.1093/brain/awq178>

857 Nan, Y., Huang, W., Wang, W., Liu, C., & Dong, Q. (2016). Subgroup differences in the lexical
858 tone mismatch negativity (MMN) among Mandarin speakers with congenital amusia.
859 *Biological Psychology*, *113*, 59–67. <https://doi.org/10.1016/j.biopsycho.2015.11.010>

860 Ohala, J. J. (1983). Cross-Language Use of Pitch: An Ethological View. *Phonetica*, *40*(1), 1–18.
861 <https://doi.org/10.1159/000261678>

862 Patel, A. D., Wong, M., Foxton, J., Lochy, A., & Peretz, I. (2008). Speech intonation perception
863 deficits in musical tone deafness (congenital amusia). *Music Perception*, *25*(4), 357–368.
864 <https://doi.org/10.1525/mp.2008.25.4.357>

865 Peretz, I., Ayotte, J., Zatorre, R. J., Mehler, J., Ahad, P., Penhune, V. B., & Jutras, B. (2002).
866 Congenital amusia: A disorder of fine-grained pitch discrimination. *Neuron*, *33*(2), 185–
867 191. [https://doi.org/10.1016/S0896-6273\(01\)00580-3](https://doi.org/10.1016/S0896-6273(01)00580-3)

868 Peretz, I., Brattico, E., & Tervaniemi, M. (2005). Abnormal electrical brain responses to pitch in
869 congenital amusia. *Annals of Neurology*, *58*(3), 478–482.
870 <https://doi.org/10.1002/ana.20606>

871 Peretz, I., Champod, A. S., & Hyde, K. (2003). Varieties of musical disorders: The montreal
872 battery of evaluation of amusia. *Annals of the New York Academy of Sciences*, *999*(1),
873 58–75. <https://doi.org/10.1196/annals.1284.006>

874 Peretz, I., Cummings, S., & Dubé, M.-P. (2007). The genetics of congenital amusia (tone
875 deafness): A family-aggregation study. *The American Journal of Human Genetics*, *81*(3),
876 582–588. <https://doi.org/10.1086/521337>

877 Peretz, I., & Vuvan, D. T. (2017). Prevalence of congenital amusia. *European Journal of Human*
878 *Genetics*, 25(5), 625–630. <https://doi.org/10.1038/ejhg.2017.15>

879 Phillips-Silver, J., Toiviainen, P., Gosselin, N., Piché, O., Nozaradan, S., Palmer, C., & Peretz, I.
880 (2011b). Born to dance but beat deaf: A new form of congenital amusia.
881 *Neuropsychologia*, 49(5), 961–969.
882 <https://doi.org/10.1016/j.neuropsychologia.2011.02.002>

883 Shao, J., Lau, R. Y. M., Tang, P. O. C., & Zhang, C. (2019). The Effects of Acoustic Variation on
884 the Perception of Lexical Tone in Cantonese-Speaking Congenital Amusics. *Journal of*
885 *Speech, Language, and Hearing Research*, 62(1), 190–205.
886 https://doi.org/10.1044/2018_JSLHR-H-17-0483

887 Shao, J., & Zhang, C. (2020). Dichotic Perception of Lexical Tones in Cantonese-Speaking
888 Congenital Amusics. *Frontiers in Psychology*, 11, 1411.
889 <https://doi.org/10.3389/fpsyg.2020.01411>

890 Shao, J., Zhang, C., Peng, G., Yang, Y., & Wang, W. S.-Y. (2016). *Effect of Noise on Lexical*
891 *Tone Perception in Cantonese-Speaking Amusics*. 272–276.
892 <https://doi.org/10.21437/Interspeech.2016-891>

893 Sharifzadeh, H. R., McLoughlin, I. V., & Russell, M. J. (2012). A Comprehensive Vowel Space
894 for Whispered Speech. *Journal of Voice*, 26(2), e49–e56.
895 <https://doi.org/10.1016/j.jvoice.2010.12.002>

896 Sun, Y., Lu, X., Ho, H. T., & Thompson, W. F. (2017). Pitch discrimination associated with
897 phonological awareness: Evidence from congenital amusia. *Scientific Reports*, 7(1),
898 44285. <https://doi.org/10.1038/srep44285>

899 Thompson, W. F. (2007). Exploring variants of amusia: Tone deafness, rhythm impairment, and
900 intonation insensitivity. *Proceedings of the Inaugural International Conference on Music*
901 *Communication Science*, 159–163.

902 <https://researchers.mq.edu.au/en/publications/exploring-variants-of-amusia-tone->
903 [deafness-rhythm-impairment-and-](https://researchers.mq.edu.au/en/publications/exploring-variants-of-amusia-tone-)
904 Thompson, W. F., Marin, M. M., & Stewart, L. (2012). Reduced sensitivity to emotional prosody
905 in congenital amusia rekindles the musical protolanguage hypothesis. *Proceedings of the*
906 *National Academy of Sciences*, *109*(46), 19027–19032.
907 <https://doi.org/10.1073/pnas.1210344109>
908 Vuvan, D. T., Nunes-Silva, M., & Peretz, I. (2015). Meta-analytic evidence for the non-
909 modularity of pitch processing in congenital amusia. *Cortex*, *69*, 186–200.
910 <https://doi.org/10.1016/j.cortex.2015.05.002>
911 Vuvan, D. T., Paquette, S., Mignault Goulet, G., Royal, I., Felezeu, M., & Peretz, I. (2018). The
912 Montreal Protocol for Identification of Amusia. *Behavior Research Methods*, *50*(2), 662–
913 672. <https://doi.org/10.3758/s13428-017-0892-8>
914 Wang, X., & Peng, G. (2014). Phonological processing in Mandarin speakers with congenital
915 amusia. *The Journal of the Acoustical Society of America*, *136*(6), 3360–3370.
916 <https://doi.org/10.1121/1.4900559>
917 Whiteford, K. L., & Oxenham, A. J. (2017b). Auditory deficits in amusia extend beyond poor
918 pitch perception. *Neuropsychologia*, *99*, 213–224.
919 <https://doi.org/10.1016/j.neuropsychologia.2017.03.018>
920 Xu, Y. (2013). *ProsodyPro—A tool for large-scale systematic prosody analysis*. 7–10.
921 Yan, J. (2016). *Experimental study on foreign accent of Chinese learners*. Beijing language and
922 culture university.
923 Yang, L., Li, Y., & Xu, B. (2005). The establishment of a Chinese whisper database and
924 perceptual experiment (in Chinese). *Journal of NanJing University (Natural Sciences)*,
925 *41*(3), 311–317.

- 926 Zhang, C., Shao, J., & Huang, X. (2017). Deficits of congenital amusia beyond pitch: Evidence
927 from impaired categorical perception of vowels in Cantonese-speaking congenital
928 amusics. *PLOS ONE*, 12(8), e0183151. <https://doi.org/10.1371/journal.pone.0183151>
- 929 Zhang, C., & Shao, J. (2018). Normal pre-attentive and impaired attentive processing of lexical
930 tones in Cantonese-speaking congenital amusics. *Scientific Reports*, 8(1), 8420.
931 <https://doi.org/10.1038/s41598-018-26368-7>
- 932 Zhang, G., Shao, J., Huang, X., Wang, L., & Zhang, C. (2018). Unequal Impairment of Native
933 and Non-native Tone Perception in Cantonese-speaking Congenital Amusics. *9th*
934 *International Conference on Speech Prosody 2018*, 562–566.
935 <https://doi.org/10.21437/SpeechProsody.2018-114>
- 936 Zhang, W., & Dong, W. (2004). *Advanced tutorial of SPSS statistical analysis (in Chinese)*.
937 Higer Education Press.
- 938
- 939

940 Table 1. Demographic characteristics of participants. The results of independent-samples t-
 941 tests comparing the amusics and controls in age and the scores of MBEA test are also reported
 942 here. n.s. = not significant. The *p*-value was corrected for multiple comparisons on the MBEA
 943 tests (.05/8 = .006).

	Amusics	Controls	<i>t</i> -value	<i>p</i> -value	Cohen's <i>d</i>
Male/Female (total)	9/10 (19)	9/10 (19)	/	/	/
Mean Age (range)	24.37 (20-30)	24.42 (20-31)	-0.07	n.s.	/
MBEA					
Scale (SD)	55.16% (14.24)	85.90% (11.16)	-7.41	<i>p</i> < 0.001	2.40
Contour (SD)	58.56% (15.50)	94.23% (4.76)	-9.59	<i>p</i> < 0.001	3.11
Interval (SD)	58.57% (7.85)	93.02% (3.78)	-17.24	<i>p</i> < 0.001	5.59
Rhythm (SD)	61.13% (13.75)	93.54% (6.88)	-9.20	<i>p</i> < 0.001	2.98
Meter (SD)	50.53% (10.44)	84.39% (12.07)	-9.25	<i>p</i> < 0.001	3.00
Memory (SD)	71.06% (16.12)	96.49% (3.92)	-6.68	<i>p</i> < 0.001	2.17
Pitch composite score (SD)	51.69 (8.16)	81.94 (4.62)	-14.01	<i>p</i> < 0.001	4.56
Global (SD)	58.84% (7.32)	91.26% (4.65)	-16.30	<i>p</i> < 0.001	5.29

944

945

946 Table 2. Results of bivariate correlations between the dependent MBEA scores and
 947 identification performance, and between the three MBEA scores. * $p < 0.05$, ** $p < 0.01$, *** p
 948 < 0.001 .

	Phonated tone	Whispered tone	Phonated intonation	Whispered intonation	Melodic organization	Temporal organization
Melodic organization	0.50**	0.39*	0.57***	0.44**		
Temporal organization	0.61***	0.47**	0.57***	0.45**	0.91***	
Melodic memory	0.47**	0.54***	0.65***	0.50**	0.82***	0.78***

949

950

951 Table 3. Results of stepwise linear regression models with the MBEA scores as predictors on
 952 tone and intonation identification. Note: The values represent standardized regression
 953 coefficients for the predictors retained in the model. Empty cells indicate that the predictor was
 954 not retained in the model. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Phonation mode	MBEA predictors			Adjusted R ² of the model
	Melodic organization	Temporal organization	Melodic memory	
Tone				
Phonated mode		0.61***		0.35***
Whispered mode			0.54***	0.27***
Intonation				
Phonated mode			0.65***	0.41***
Whispered mode			0.50**	0.23**

955

956 Table 4. Results of multinomial logistic regression models with the acoustic cues as predictors
 957 on phonated tone identification. Note: The tone following versus was used as the baseline for
 958 the contrast (e.g., T1 was the baseline in the T2 vs. T1 pair). Controls were used as the baseline
 959 for comparison with amusics. The values represent regression coefficients (*B* (odds ratio)) for
 960 the predictors. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. For significant interactions between group
 961 and acoustic cues, cases where a certain acoustic cue led to worse performance in amusics than
 962 in controls were marked in bold. In contrast, cases where a certain acoustic cue led to better
 963 performance in amusics than in controls were marked with #.

Tones	Main effects of predictors						
	F0 mean	F0 SD	Duration	Intensity	F1	F2	Group
T2 vs. T1	-0.03*** (0.97)	0.35*** (1.41)	-0.01*** (0.99)	-0.28 *** (0.75)	0.002*** (1.002)	-0.001** (0.99)	-15.57** (1.74e-7)
T3 vs. T1	-0.11 *** (0.90)	0.43*** (1.54)	0.02*** (1.02)	-0.57*** (0.57)	0.005*** (1.005)	0.0003 (1.0003)	-15.84* (1.32e-7)
T4 vs. T1	-0.001 (0.99)	0.39*** (1.48)	-0.04*** (0.96)	0.23*** (1.26)	-0.005*** (0.99)	-0.002*** (0.99)	-12.10* (5.78e-6)
T3 vs. T2	-0.08*** (0.93)	0.09** (1.09)	0.04*** (1.04)	-0.28*** (0.75)	0.003*** (1.003)	0.001*** (1.001)	-0.28 (0.76)
T4 vs. T2	0.03*** (1.03)	0.04*** (1.04)	-0.03*** (0.97)	0.52*** (1.68)	-0.007*** (0.99)	-0.002*** (0.99)	3.50 (33.22)
T4 vs. T3	0.11*** (1.12)	-0.05 (0.96)	-0.07*** (0.94)	0.80*** (2.22)	-0.01*** (0.99)	-0.002*** (0.99)	3.78 (43.89)
Interactions between group and acoustic cues							

T2 vs. T1	0.01** (1.01)	-0.14*** (0.87)	0.008*** (1.008)	0.16** (1.17)	-0.001 (0.99)	0.0004*# (1.0004)	/
T3 vs. T1	0.06*** (1.06)	-0.19*** (0.83)	-0.01*** (0.99)	0.26** (1.30)	-0.002**# (0.99)	-0.0002 (0.99)	/
T4 vs. T1	-0.005 (0.99)	-0.12*** (0.88)	0.02*** (1.02)	0.06 (1.07)	0.001 (1.001)	0.001***# (1.001)	/
T3 vs. T2	0.05*** (1.05)	-0.05 (0.95)	-0.02*** (0.98)	0.10 (1.11)	-0.001*# (0.99)	-0.001***# (0.99)	/
T4 vs. T2	-0.02** (0.99)	0.02 (1.02)	0.009*** (1.01)	-0.10 (0.91)	0.002**# (1.002)	0.0005***# (1.0005)	/
T4 vs. T3	-0.06*** (0.94)	0.07*# (1.07)	0.03*** (1.03)	-0.20* (0.82)	0.003***# (1.003)	0.001***# (1.001)	/

964

965

966 Table 5. Results of binominal regression models with the acoustic cues as predictors on
 967 phonated intonation identification. Note: The statement was used as the baseline for contrast
 968 with the question. Controls were used as the baseline for comparison with amusics. The values
 969 represent regression coefficients (*B* (odds ratio)) for the predictors. ‘F0 dir’ means F0 direction.
 970 * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. For significant interactions between group and acoustic
 971 cues, cases where a certain acoustic cue led to worse performance in amusics than in controls
 972 were marked in bold. In contrast, cases where a certain acoustic cue led to better performance in
 973 amusics than in controls were marked with #.

Main effects of predictors							
Sentence					Final syllable		Group
F0 mean	F0 SD	F0 dir	Duration	Intensity	F0 SD	Duration	
0.02	0.11	2.17	-0.0002	1.19	0.02	0.02	43.99
***	***	***		***	*	***	***
(1.02)	(1.11)	(8.75)	(0.99)	(3.29)	(1.02)	(1.02)	(1.27e+19)
Interactions between group and acoustic cues							
-0.007	-0.05	-1.07	-0.0002	-0.55	-0.02	-0.01	/
***	***	***		***		***	
(0.99)	(0.95)	(0.34)	(0.99)	(0.58)	(0.98)	(0.99)	

974

975

976 Table 6. Results of multinomial logistic regression models with the acoustic cues as predictors
 977 on whispered tone identification. Note: The tone following versus was used as the baseline for
 978 the contrast (e.g., T1 was the baseline in the T2 vs. T1 pair). Controls were used as the baseline
 979 for comparison with amusics. The values represent regression coefficients (*B* (odds ratio)) for
 980 the predictors. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$. For significant interactions between group
 981 and acoustic cues, cases where a certain acoustic cue led to worse performance in amusics than
 982 in controls were marked in bold. In contrast, cases where a certain acoustic cue led to better
 983 performance in amusics than in controls were marked with #.

Tones	Main effects of predictors				
	Duration	Intensity	F1	F2	Group
T2 vs. T1	-0.002*** (0.99)	0.004 (1.004)	0.001*** (1.001)	0.001*** (1.001)	0.78 (2.18)
T3 vs. T1	0.008*** (1.008)	0.04** (1.04)	-0.002*** (0.99)	-0.0001 (0.99)	1.17 (3.22)
T4 vs. T1	-0.01*** (0.99)	-0.08*** (0.92)	0.002*** (1.002)	0.001*** (1.001)	-0.99 (0.42)
T3 vs. T2	0.01*** (1.01)	0.04* (1.04)	-0.003*** (0.99)	-0.001*** (0.99)	0.39 (1.48)
T4 vs. T2	-0.01*** (0.99)	-0.09*** (0.92)	0.001** (1.001)	0.001*** (1.001)	-1.77 (0.17)
T4 vs. T3	-0.02*** (0.98)	-0.12*** (0.89)	0.004*** (1.004)	0.001*** (1.001)	-2.16 (0.12)
Interaction between the group and acoustic cues					
T2 vs. T1	0.002* (1.002)	-0.009 (0.99)	-0.0003 (0.99)	-0.0005***# (0.99)	/

T3 vs. T1	-0.001 (0.99)	-0.009 (0.99)	0.001 (1.001)	-0.0003*# (0.99)	/
T4 vs. T1	0.003*** (1.003)	0.005 (1.008)	-0.0003 (0.99)	-0.0005** (0.99)	/
T3 vs. T2	-0.003*** (0.99)	0.0002 (1.0002)	0.001* (1.001)	0.0002 (1.0002)	/
T4 vs. T2	0.002 (1.002)	0.01 (1.01)	0.00006 (1.00006)	0.000008 (1.000008)	/
T4 vs. T3	0.004*** (1.004)	0.01 (1.01)	-0.001 (0.99)	-0.0002 (0.99)	/

984

985

986 Table 7. Results of binominal logistic regression models with the acoustic cues as predictors on
 987 whispered intonation identification. Note: The statement was used as the baseline for contrast
 988 with the question. Controls were used as the baseline for comparison with amusics. The values
 989 represent regression coefficients (*B* (odds ratio)) for the predictors. * $p < 0.05$, ** $p < 0.01$, ***
 990 $p < 0.001$. For significant interactions between group and acoustic cues, cases where a certain
 991 acoustic cue led to worse performance in amusics than in controls were marked in bold. In
 992 contrast, cases where a certain acoustic cue led to better performance in amusics than in controls
 993 were marked with #.

Main effects of predictors				
Sentence		Final syllable		Group
Duration	Intensity	Duration	Intensity	
-0.001*** (0.99)	0.03 (1.03)	0.01*** (1.01)	0.02 (1.02)	-0.62 (0.54)
Interaction between the group and acoustic cues				
0.0003* (1.0003)	0.03 (1.03)	-0.002 (0.99)	-0.02 (0.98)	/

994

995

996 Figure 1. The F0 contours of the four Mandarin tones produced in the phonated mode by two
997 native speakers.

998

999 Figure 2. Real-time F0 contours of a statement-question pair produced by the male speaker. This
1000 sentence is ‘高兵喝鸡汤./?’ /kau55 piəŋ55 xə55 tɕi55 tʰaŋ55/ ‘Gao Bing drinks chicken soup./?’,
1001 in which all the syllables carried T1.

1002

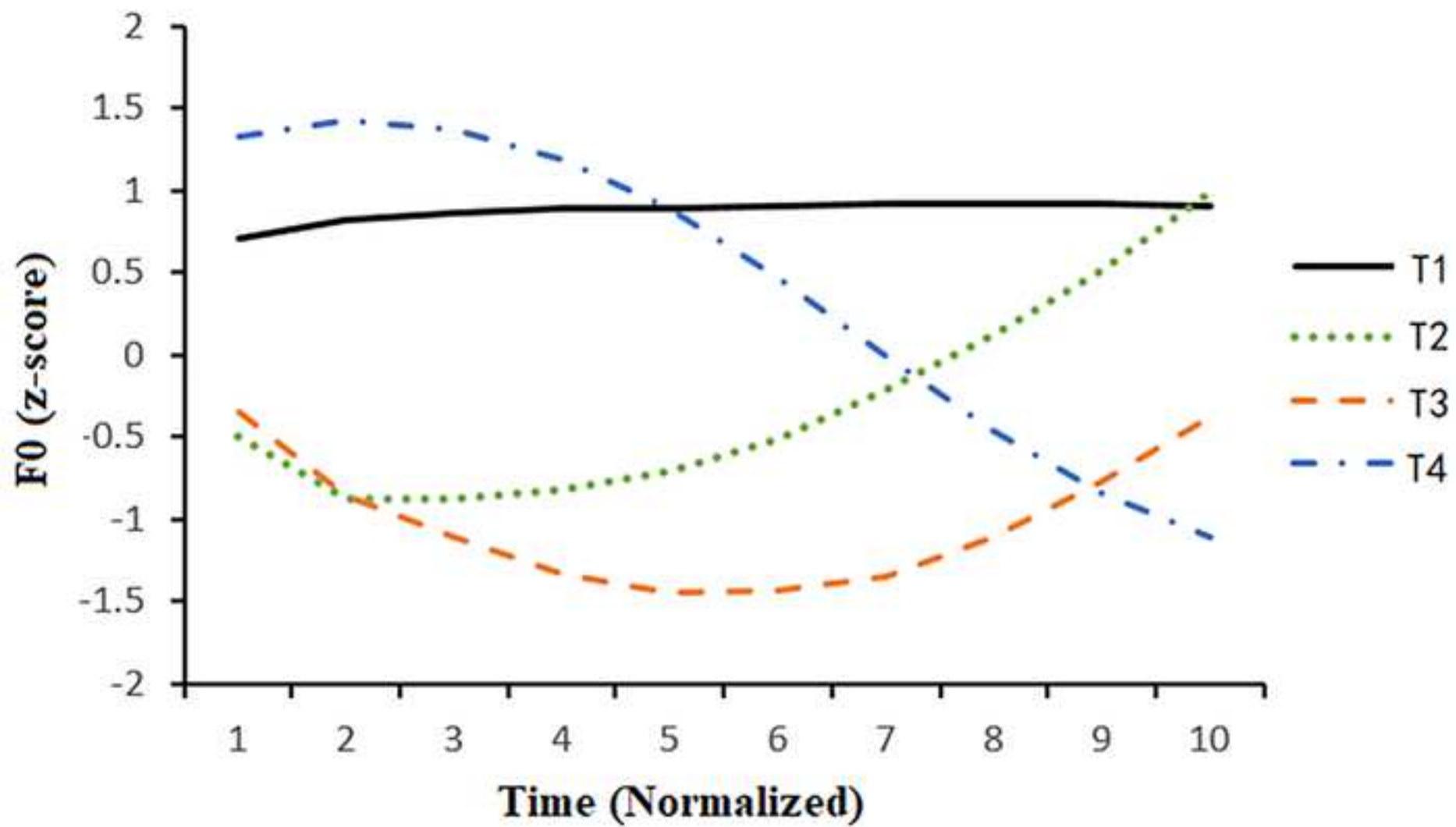
1003 Figure 3. The tone identification accuracy in the phonated and whispered mode in the amusics
1004 and controls.

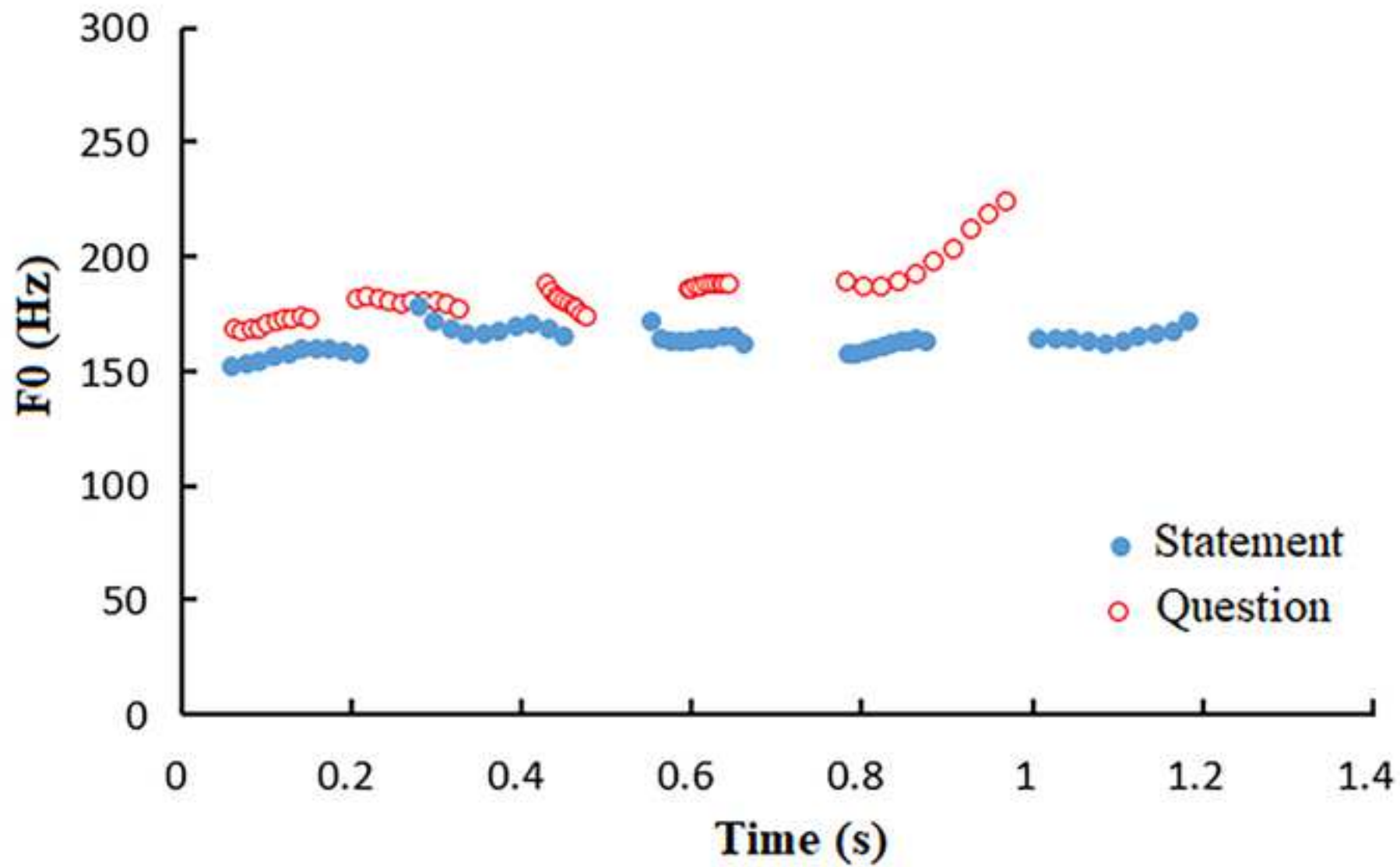
1005

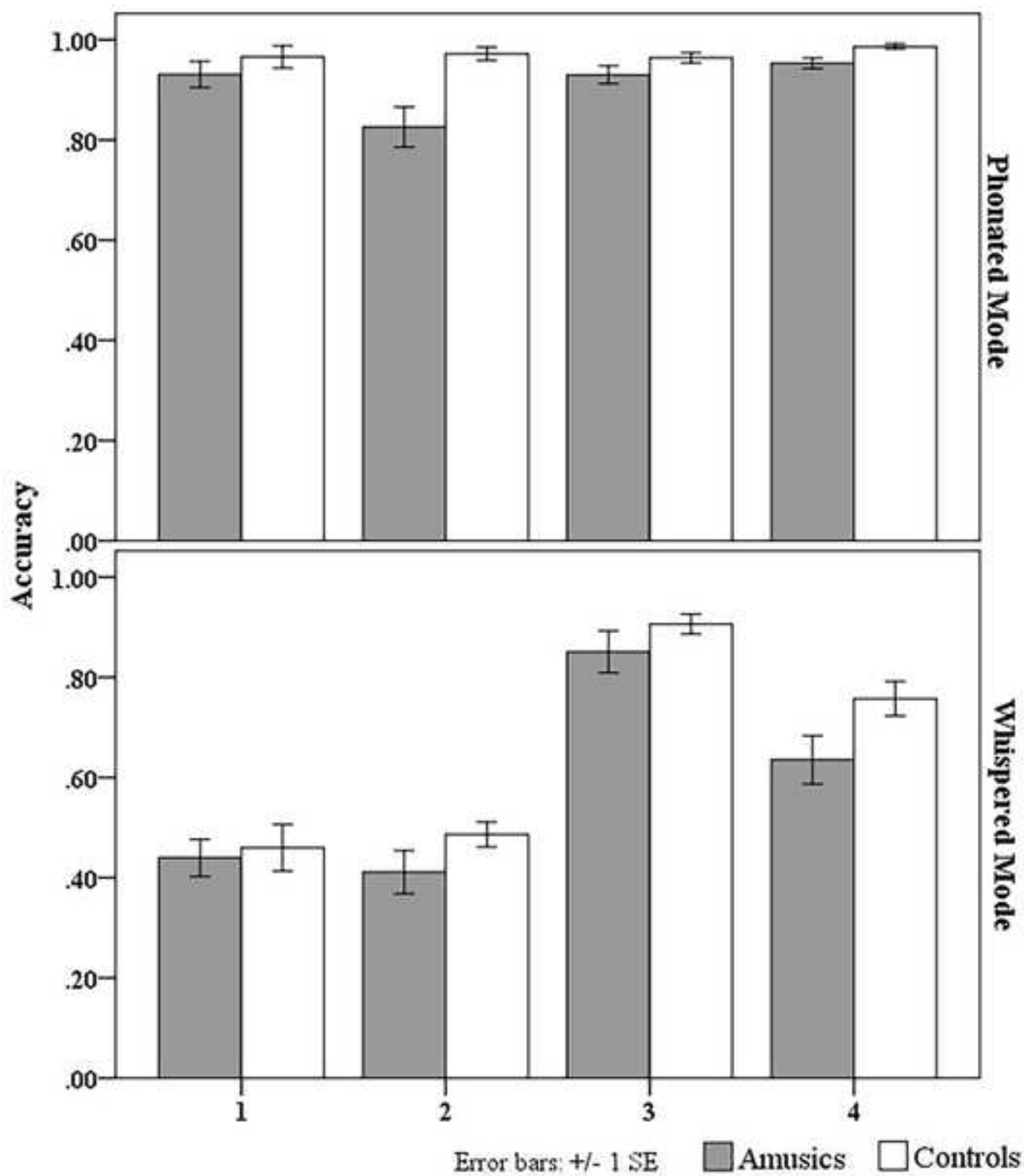
1006 Figure 4. The d’ of intonation identification in the phonated and whispered mode in the amusics and
1007 controls.

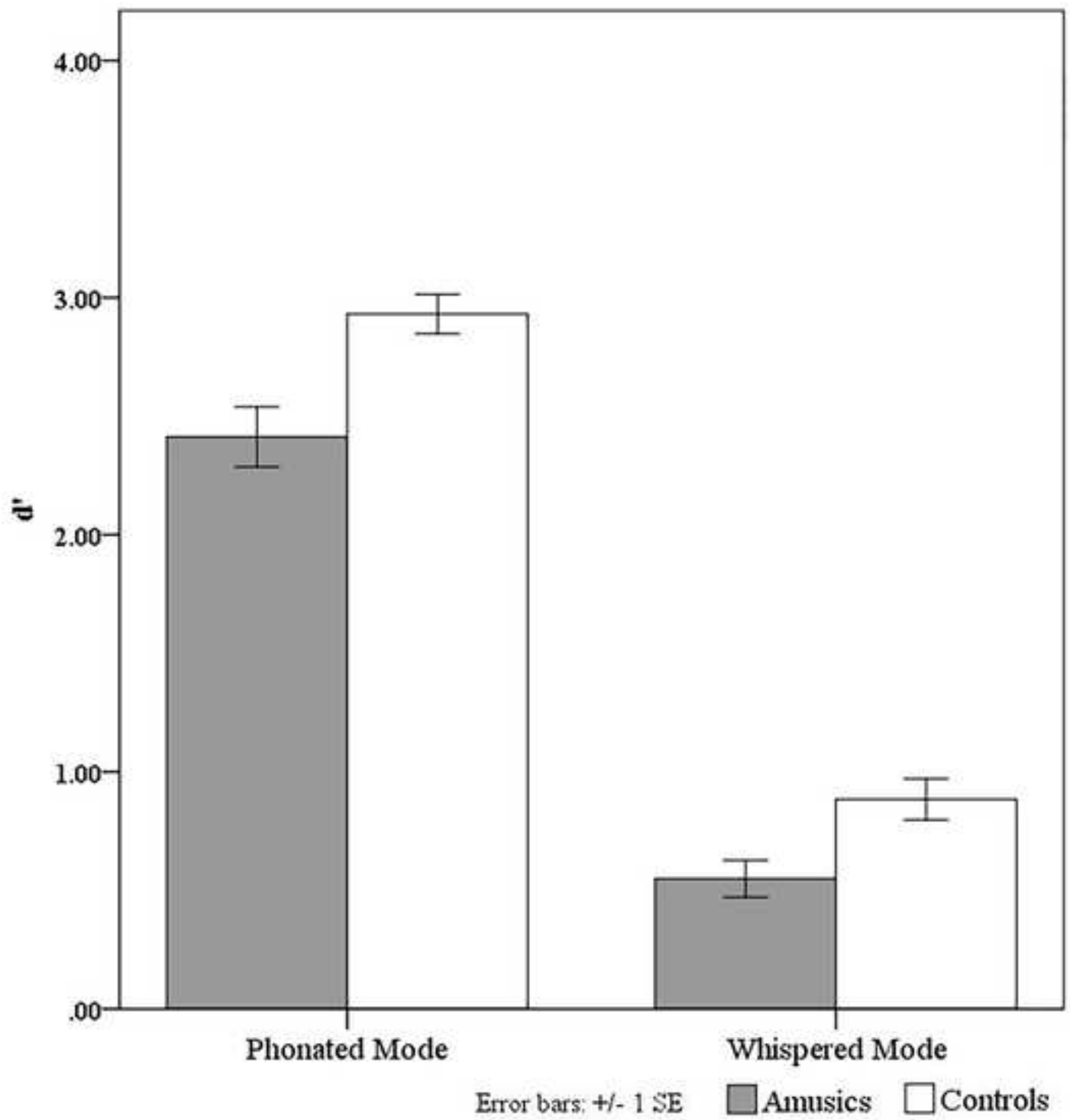
1008

1009 The supplementary information file includes a total of ten tables and five figures. The tables
1010 showed the word and sentence stimuli, acoustic cues of the four tones and two intonation patterns
1011 in the phonated and whispered mode, confusion matrices and other supportive statistical analysis
1012 results. The figures displayed the relationships between the participants' tone and intonation
1013 identification performance (in both phonation modes) and their MBEA scores and acoustic cues.









Supplementary Information

Table 1. A list of syllables and sentences for production and perception experiments.

Syllables									
Tone	/ta/	/ti/	/tu/	/pa/	/pi/	/tʃu/	/a/	/i/	/u/
1	搭 (build)	低 (low)	督 (supervise)	八 (eight)	逼 (force)	猪 (pig)	啊 (oh)	一 (one)	屋 (house)
2	答 (answer)	敌 (enemy)	毒 (poison)	拔 (pull)	鼻 (nose)	竹 (bamboo)	啊 (eh)	遗 (pity)	无 (nothing)
3	打 (fight)	底 (bottom)	赌 (bet)	把 (handle)	笔 (pen)	主 (lord)	啊 (what)	已 (already)	五 (five)
4	大 (big)	地 (land)	肚 (belly)	爸 (father)	币 (coin)	祝 (wish)	啊 (ah)	易 (easy)	误 (error)
Sentences									
4 syllables	张薇开车。/? /tʃaŋ55 uei55 kʰai55 tʃʰə55/ (Zhang Wei drives the car./?)				5 syllables	高兵喝鸡汤。/? /kau55 piŋ55 xə55 tɕi55 tʰaŋ55/ (Gao Bing drinks chicken soup./?)			
	王梅划船。/? /uaŋ35 mei35 xua35 tʃʰuan35/ (Wang Mei boats./?)					罗婷学轮滑。/? /luo35 tʰiŋ35 eyɛ35 luən35 xua35/ (Luo Ting learns skating./?)			
	李敏点火。/? /li214 miən214 dien214 huo214/ (Li Min makes a fire./?)					李伟买雨伞。/? /li214 uei214 mai214 y214 san214/ (Li Wei buys an umbrella./?)			
	叶亮睡觉。/? /iɛ51 liɑŋ51 ʃuei51 tɕiau51/ (Ye Liang sleeps./?)					赵志看电视。/? /tʃəu51 tʃʰi51 kʰan51 tien51 ʃʌ51/ (Zhao Zhi watches TV./?)			
	李刚讲课。/? /li214 kaŋ55 tɕiaŋ214 kʰə51/ (Li Gang gives a lesson./?)					李刚交水费。/? /li214 kaŋ55 tɕiau55 ʃuei214 fei51/ (Li Gang pays water fee./?)			

6 syllables	张薇担心肖英。/? /tʂɑŋ55 uei55 tan55 ɕiən55 ɕiɑu iəŋ55/ (Zhang Wei worries about Xiao Ying./?)	7 syllables	高兵今天喝鸡汤。/? /kau55 piəŋ55 tɕiən55 tʰien55 xə55 tɕi55 tʰɑŋ55/ (Gao Bing drinks chicken soup today./?)
	王梅怀疑刘宁。/? /uaŋ35 mei35 xuai35 i35 liou35 niəŋ35/ (Wang Mei suspects Liu Ning.)		罗婷明年学轮滑。/? /luo35 tʰiəŋ35 miəŋ35 niən35 ɕyɛ35 luən35 xua35/ (Luo Ting will learn skating next year./?)
	李敏反感刘雨。/? /li214 miən214 fan214 kan214 liou35 y214/ (Li Min is disgusted with Liu Yu)		李伟五点买雨伞。/? /li214 uei214 u214 tien214 mai214 y214 san214/ (Li Wei buys an umbrella at 5 o'clock./?)
	叶亮害怕赵丽。/? /ie51 liaŋ51 xai51 pʰa51 tʂɑu51 li51/ (Ye Liang is afraid of Zhang Li./?)		赵志半夜看电视。/? /tʂɑu51 tʂʅ51 pan51 ie51 kʰan51 tien51 ʂʅ51/ (Zhao Zhi watches TV at midnight./?)
	李刚讨厌吕梦。/? /li214 kaŋ55 tʰau214 ien51 ly214 mən̄51/ (Li Gang hates Lü Meng./?)		李刚九号交水费。/? /li214 kaŋ55 tɕiəu214 xau51 tɕiəu55 ʂuei214 fei51/ (Li Gang pays water fee on ninth./?)
10 syllables	张薇担心肖英开车发晕。/? /tʂɑŋ55 uei55 tan55 ɕiən55 ɕiɑu55 iəŋ55 kʰai55 tʂʰə55 fa55 yən55/ (Zhang Wei worries about Xiao Ying having a carsickness./?)		
	王梅怀疑刘宁划船着迷。/? /uaŋ35 mei35 xuai35 i35 liou35 niəŋ35 xua35 tʂʰuan35 tʂɑu35 mi35/ (Wang Mei suspects Liu Ning of indulging in boating./?)		
	李敏反感刘雨点火取暖。/? /li214 miən214 fan214 kan214 liou35 y214 tien214 xuo214 tɕʰy214 nuan214/ (Li Min is disgusted with Liu Yu making a fire for warmth./?)		
	叶亮害怕赵丽睡觉做梦。/? /ie51 liaŋ51 xai51 pʰa51 tʂɑu51 li51 ʂuei51 tɕiəu51 tsuo51 mən̄51/ (Ye Liang is afraid of Zhang Li dreaming when sleeping./?)		

李刚讨厌吕梦讲课紧张。 /? /li214 kaŋ55 t ^h au214 ien51 ly214 məŋ51 tɕiaŋ214 k ^h ə51 tɕiən214 tɕaŋ55/ (Li Gang hates Lü Meng to be nervous when teaching./?)

Table 2. Acoustic characteristics of the four Mandarin tones produced in the phonated mode and results of one-way ANOVAs conducted to compare the four tones on each acoustic cue (the *p*-value was corrected for multiple comparisons: .05/6 = .008).

Tone	F0 _{mean} (Hz)	F0 _{SD}	Duration (ms)	Intensity (dB)	F1 (Hz)	F2 (Hz)
T1 (SD)	194.11 (8.75)	3.53 (2.39)	516.95 (64.93)	72.06 (3.46)	586.05 (360.61)	1655.66 (842.74)
T2 (SD)	154.83 (11.06)	23.03 (4.01)	470.17 (54.68)	71.44 (2.70)	577.33 (369.32)	1664.7 (828.89)
T3 (SD)	128.04 (6.02)	15.59 (3.21)	662.97 (60.79)	68.01 (1.85)	574.02 (375.87)	1621.73 (875.03)
T4 (SD)	173.68 (8.2)	34.58 (5.93)	355.29 (68.84)	72.81 (2.98)	613 (354.73)	1677.88 (834.75)
<i>p</i> value	<i>p</i> < .001	<i>p</i> < .001	<i>p</i> < .001	<i>p</i> = .005	n.s.	n.s.
$\eta^2_{partial}$	0.90	0.89	0.78	0.33	0.02	0.01

Table 3. Acoustic characteristics of the four Mandarin tones produced in the whispered mode and results of one-way ANOVAs conducted to compare the four tones on each acoustic cue (the *p*-value was corrected for multiple comparisons: $.05/4 = .0125$).

Tone	Duration (ms)	Intensity (dB)	F1 (Hz)	F2 (Hz)
T1 (SD)	515.94 (42.05)	54.02 (5.35)	874.46 (259.63)	1867.19 (621.35)
T2 (SD)	472.93 (44.75)	54.08 (5.23)	833.65 (290.29)	1826.28 (632.03)
T3 (SD)	629.10 (63.24)	51.74 (5.48)	809.94 (290.36)	1789.99(666.22)
T4 (SD)	365.81 (58.34)	55.81 (4.17)	838.73 (276.04)	1872.67 (602.65)
<i>p</i> value	$p < .001$	n.s.	n.s.	n.s.
$\eta^2_{partial}$	0.78	0.08	0.08	0.03

Table 4. Acoustic characteristics of statements and questions in the phonated mode and results of t-tests conducted to compare statement and question on each acoustic cue (the *p*-value was corrected for multiple comparisons: $.05/9 = .005$).

Sentence					
Intonation	F0 mean (Hz)	F0 SD	F0 direction	Duration (ms)	Intensity (dB)
Statement (SD)	163.28 (28.16)	22.34 (10.81)	-0.42 (0.40)	1745.91 (541.99)	69.46 (1.41)
Question (SD)	213.36 (26.77)	29.27 (14.96)	0.22 (0.60)	1463.45 (446.62)	72.41 (1.12)
<i>p</i> value	<i>p</i> < .001	<i>p</i> < .001	<i>p</i> < .001	<i>p</i> < .001	<i>p</i> < .001
Cohen's <i>d</i>	1.82	0.53	1.26	0.57	2.32
Final syllable					
Intonation	F0 mean (Hz)	F0 SD	Duration (ms)	Intensity (dB)	
Statement (SD)	152.32 (36.02)	19.33 (12.31)	292.69 (47.06)	67.05 (2.16)	
Question (SD)	222.48 (46.96)	24.12 (11.48)	313.13 (38)	72.51 (2.19)	
<i>p</i> value	<i>p</i> < .001	n.s.	n.s. (<i>p</i> = .012)	<i>p</i> < .001	
Cohen's <i>d</i>	1.68	0.4	0.48	2.51	

Table 5. Acoustic characteristics of statements and questions in the whispered mode and results of t-tests conducted to compare statement and question on each acoustic cue (the *p*-value was corrected for multiple comparisons: $.05/4 = .0125$).

Intonation	Sentence		Final Syllable	
	Duration (ms)	Mean Intensity (dB)	Duration (ms)	Mean Intensity (dB)
Statement (SD)	1645.67 (482.22)	53.46 (2.49)	305.11 (61.13)	52.29 (5.17)
Question (SD)	1453.95 (470.21)	53.25 (2.09)	330.95 (42.82)	53.54 (4.00)
<i>p</i> value	$p < .001$	n.s.	$p = .001$	n.s.
Cohen's <i>d</i>	0.40	0.10	0.49	0.27

Table 6. Results of correlations among the nine acoustic cues of the phonated sentences. Note: ‘S’ means sentence (e.g., ‘S F0 mean’ means the mean F0 of the whole sentences). ‘Fs’ means final syllable (e.g., ‘FS F0 mean’ means the mean F0 of the final syllable).

Pearson Correlation/ items	S F0 mean	S F0 SD	S F0 dir	S duration	S intensity	Fs F0 mean	Fs F0 SD	Fs duration	Fs intensity
S F0 mean	1	0.52	0.21	0.06	-0.2	0.84	0.25	0.18	-0.002
S F0 SD		1	- 0.03	0.008	-0.23	0.32	0.61	0.18	-0.03
S F0 dir			1	-0.09	0.44	0.49	- 0.24	0.08	0.44
S duration				1	-0.29	0.04	- 0.04	-0.05	0.34
S intensity					1	- 0.002	- 0.19	-0.19	0.83
Fs F0 mean						1	0.02	0.007	0.17
Fs F0 SD							1	0.15	0.14
Fs duration								1	0.12
Fs intensity									1

Table 7. Confusion matrix of tone identification in the phonated mode. For each tonal category of the stimuli, the target tone response was in bold, and the tone response receiving the highest confusion rate was italicized.

Group	Heard	T1 (%)	T2 (%)	T3 (%)	T4 (%)
	Original				
Controls	T1	96.59	<i>2.92</i>	0.19	0.29
	T2	<i>2.14</i>	97.17	0.29	0.39
	T3	0.39	<i>2.34</i>	96.39	0.68
	T4	0.29	<i>0.78</i>	0.29	98.64
Amusics	T1	93.08	<i>4.87</i>	1.36	0.58
	T2	6.63	82.55	<i>9.16</i>	1.56
	T3	1.75	<i>4.39</i>	92.98	0.88
	T4	1.56	<i>1.85</i>	1.27	95.32

Table 8. Confusion matrix of tone identification in the whispered mode. For each tonal category of the stimuli, the target tone response was in bold, and the tone response receiving the highest confusion rate was italicized.

Group	Heard Original	T1 (%)	T2 (%)	T3 (%)	T4 (%)
	Controls	T1	46.00	12.38	16.96
T2		12.86	48.64	29.53	8.97
T3		1.27	<i>7.50</i>	90.64	0.58
T4		<i>14.23</i>	4.58	5.46	75.73
Amusics	T1	43.96	<i>21.05</i>	18.42	16.57
	T2	17.64	41.13	<i>32.94</i>	8.28
	T3	4.58	8.28	85.09	2.05
	T4	<i>18.23</i>	9.36	8.87	63.55

Table 9. Confusion matrix of intonation identification.

Phonated mode	Heard	Question	Statement	Whispered mode	Question	Statement
	Original	(%)	(%)		(%)	(%)
Controls	Question	96.89	3.11	Controls	51.21	48.79
	Statement	1.95	98.05		12.42	87.58
Amusics	Question	89.63	10.37	Amusics	41.26	58.74
	Statement	3.47	96.53		16.58	83.42

Table 10. Results of simultaneous linear regression models with the MBEA scores as predictors on tone and intonation identification. Note: The values represent standardized regression coefficients for the predictors retained in the model. The values in parentheses present VIF. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Phonation mode	MBEA predictors			R ² of the model
	Melodic organization	Temporal organization	Melodic memory	
Tone				
Phonated mode	-0.38 (7.22)	0.89* (6.2)	0.09 (3.06)	0.62**
Whispered mode	-0.61	0.55	0.60*	0.59**
Intonation				
Phonated mode	-0.07	0.21	0.55*	0.66***
Whispered mode	-0.08	0.21	0.40	0.51*

Figure 1. The relationship with the MBEA scores and tone and intonation identification in both phonation modes.

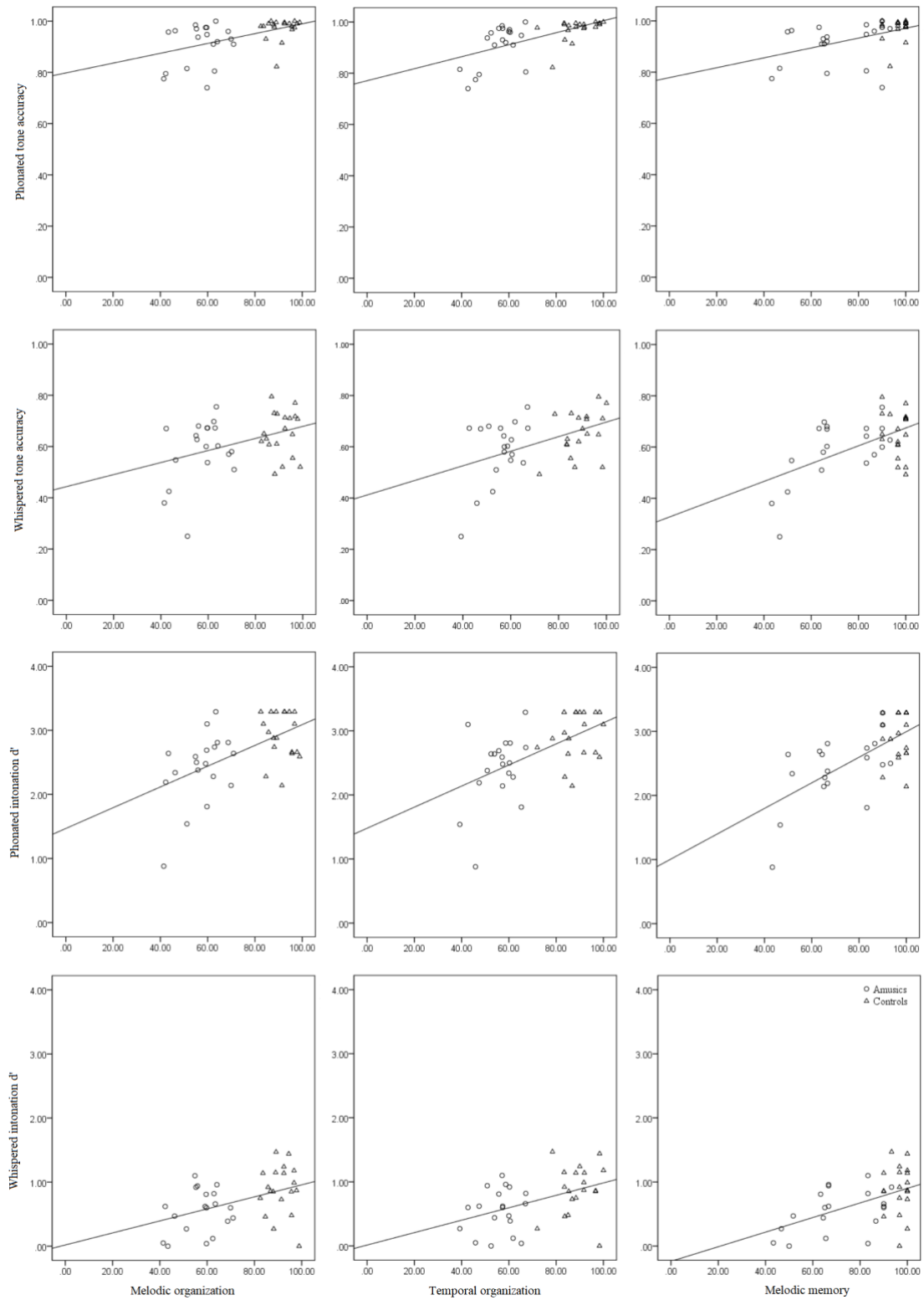
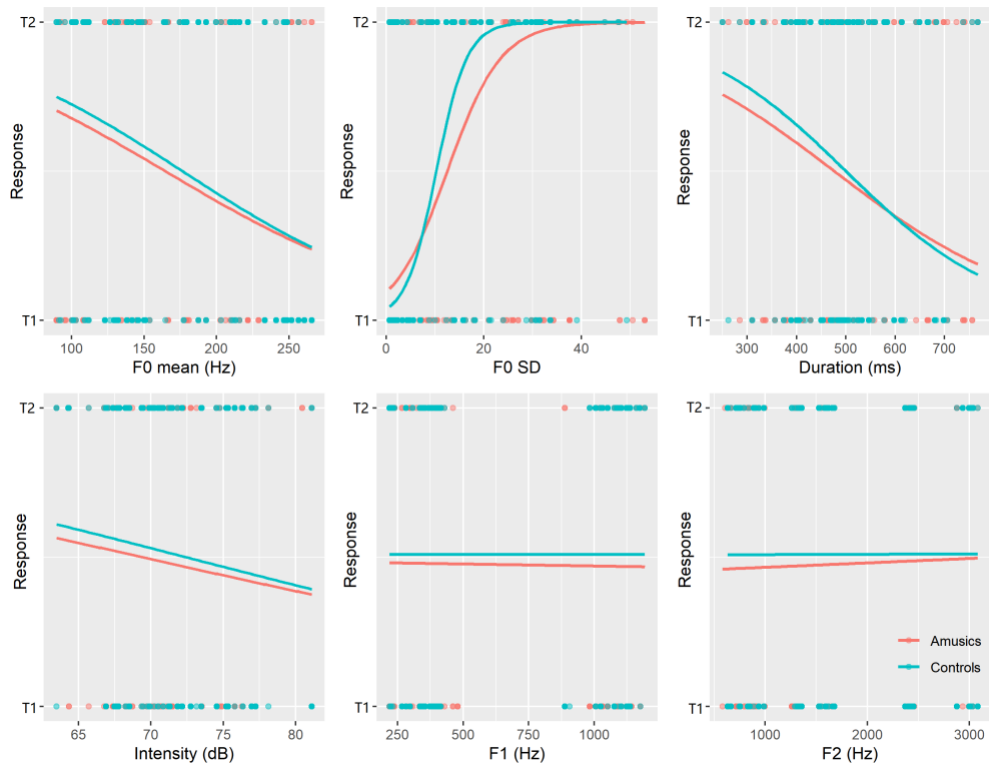
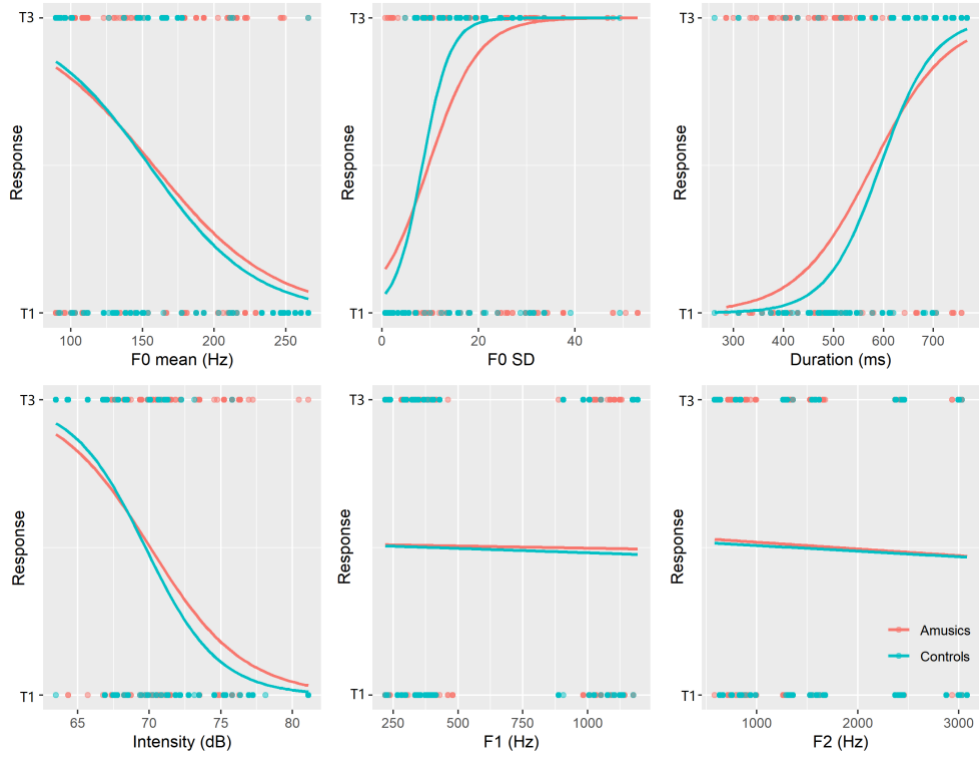


Figure 2. The relationship between the acoustic cues and phonated tone identification.

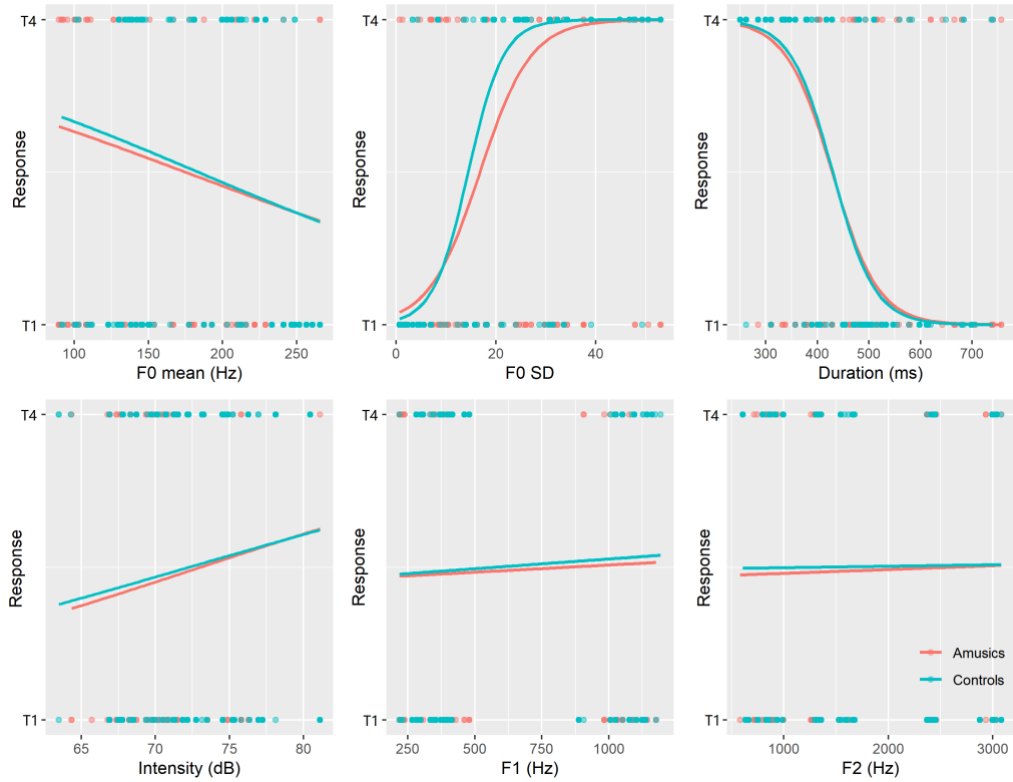
(a) The relationship between the acoustic cues and phonated T2 and T1 identification.



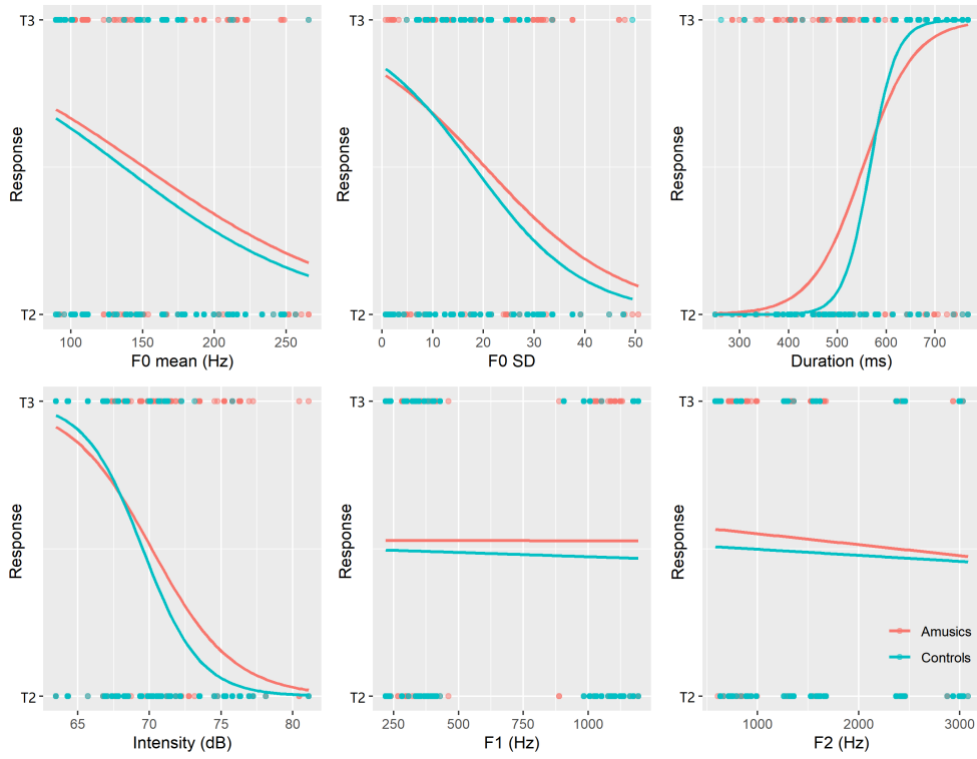
(b) The relationship between the acoustic cues and phonated T3 and T1 identification.



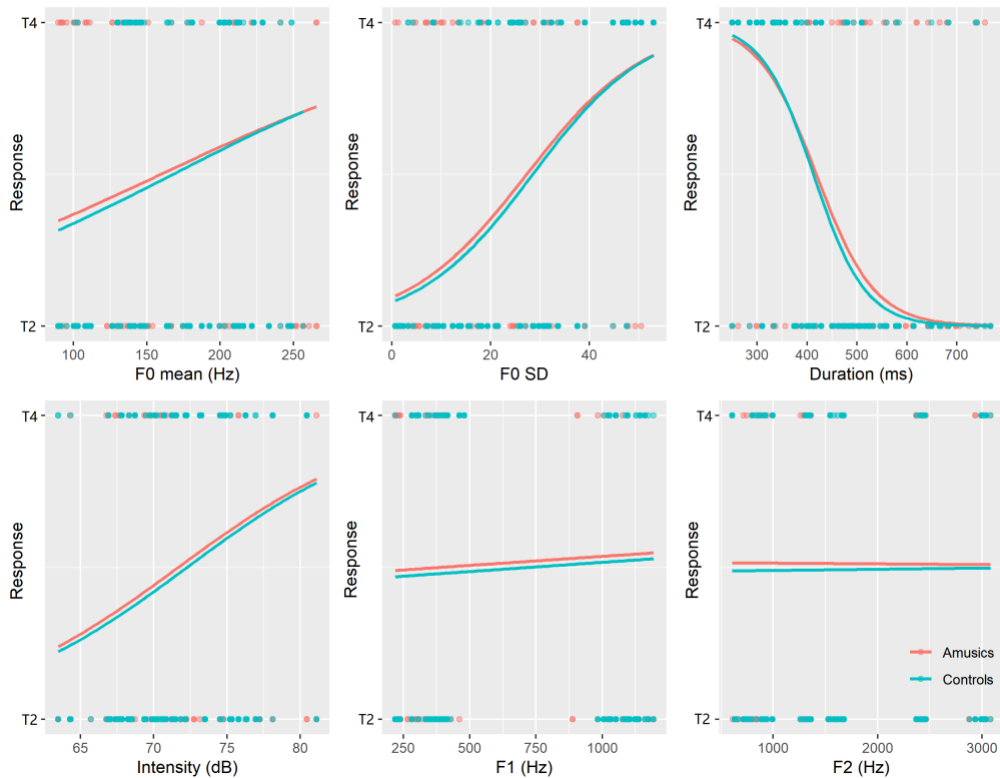
(c) The relationship between the acoustic cues and phonated T4 and T1 identification.



(d) The relationship between the acoustic cues and phonated T3 and T2 identification.



(e) The relationship between the acoustic cues and phonated T4 and T2 identification.



(f) The relationship between the acoustic cues and phonated T4 and T3 identification.

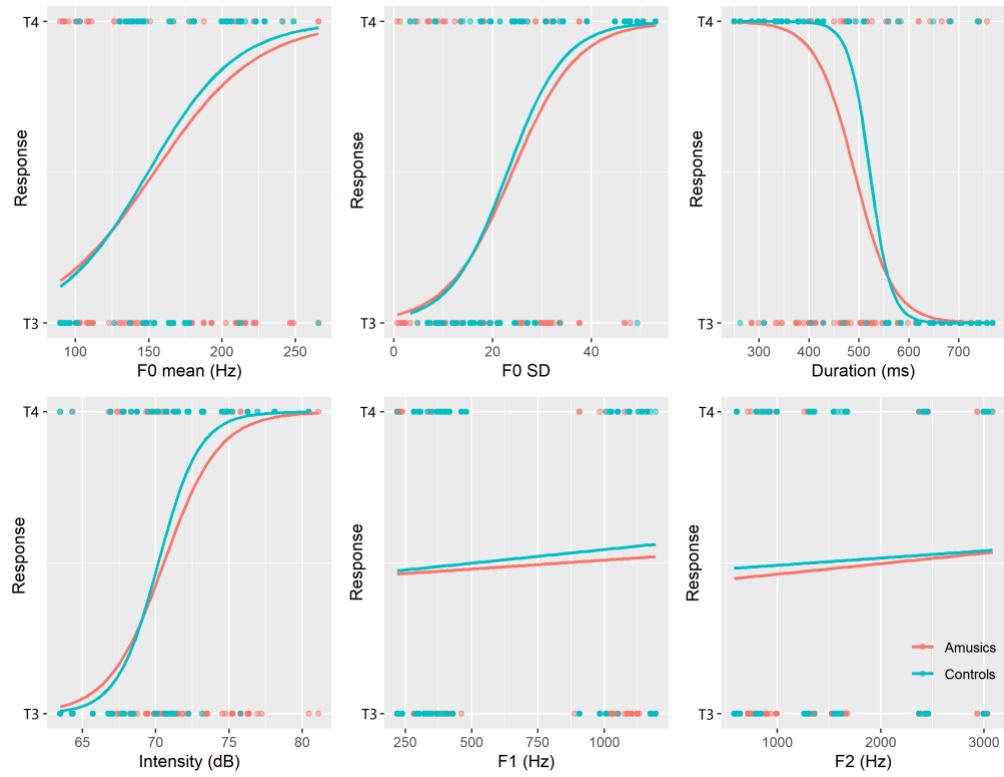


Figure 3. The relationship between the acoustic cues and phonated intonation identification.

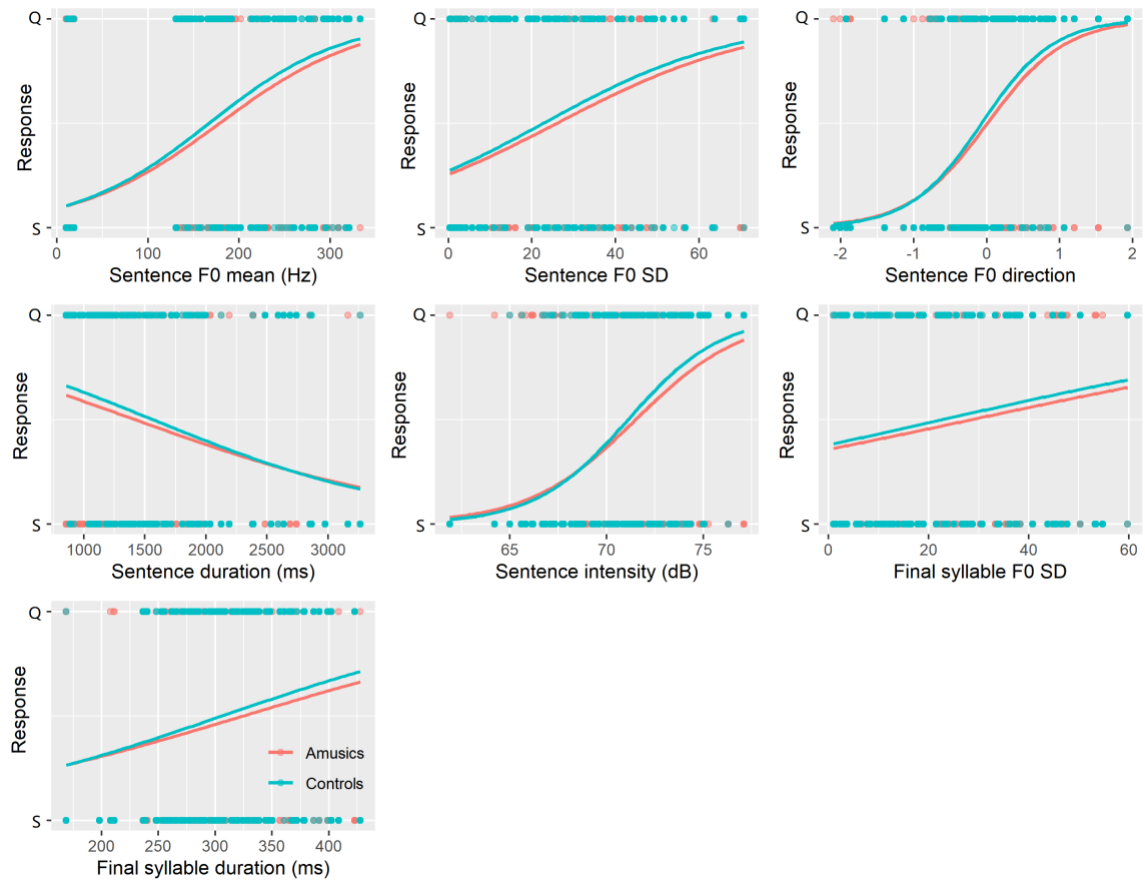
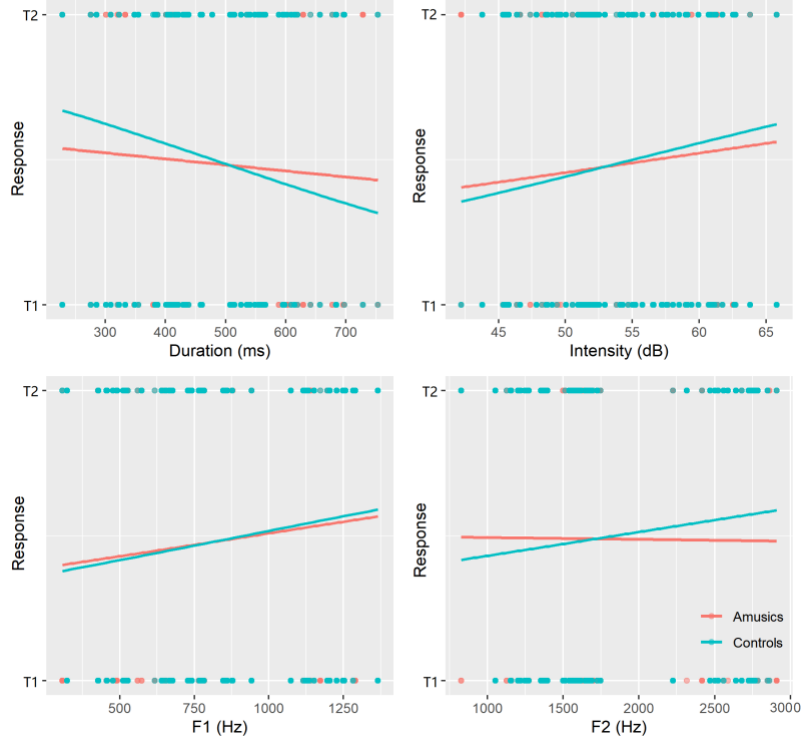
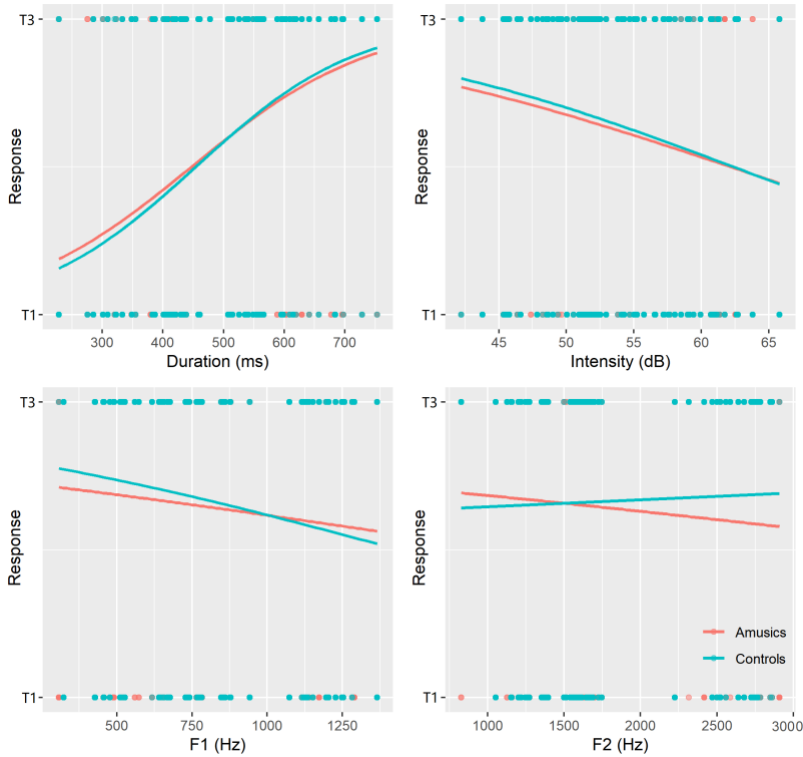


Figure 4. The relationship between the acoustic cues and whispered tone identification.

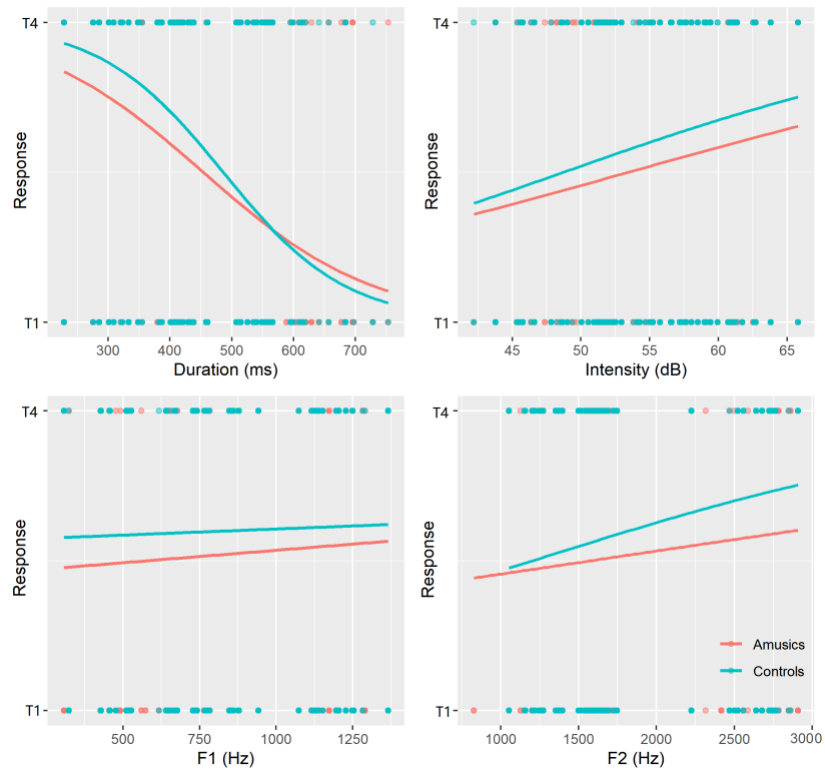
(a) The relationship between the acoustic cues and whispered T2 and T1 identification.



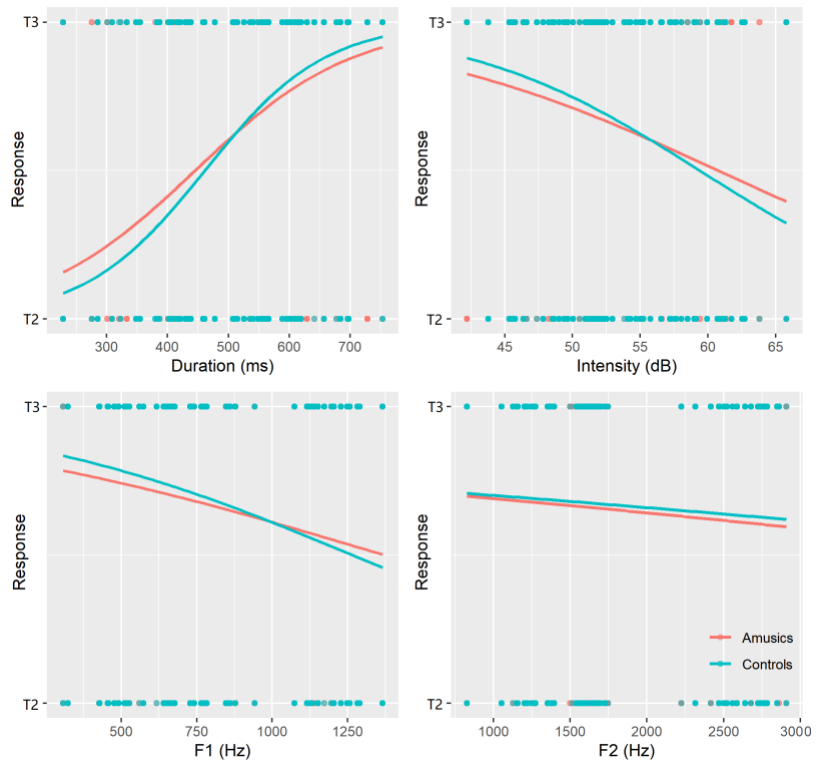
(b) The relationship between the acoustic cues and whispered T3 and T1 identification.



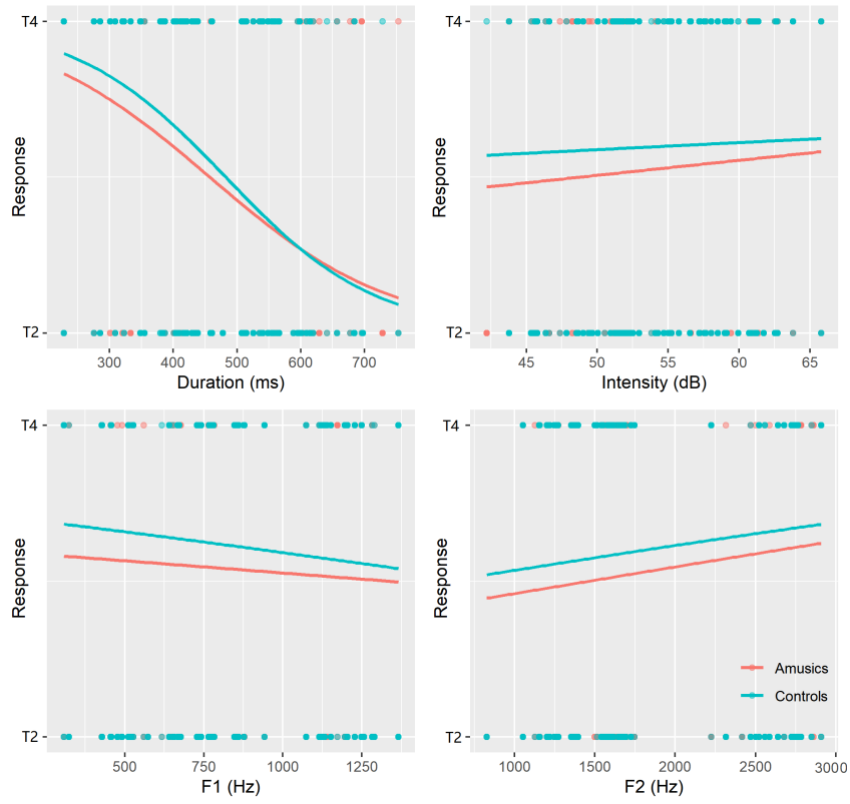
(c) The relationship between the acoustic cues and whispered T4 and T1 identification.



(d) The relationship between the acoustic cues and whispered T3 and T2 identification.



(e) The relationship between the acoustic cues and whispered T4 and T2 identification.



(f) The relationship between the acoustic cues and whispered T4 and T3 identification.

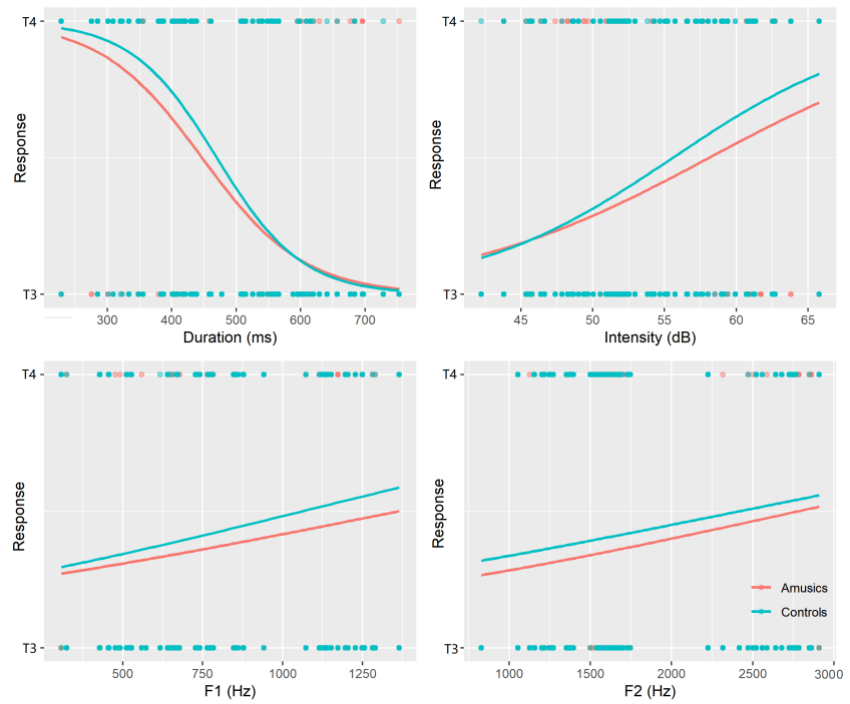


Figure 5. The relationship between the acoustic cues and whispered intonation identification.

