

The following publication Shao, J., Lau, R. Y. M., Tang, P. O. C., & Zhang, C. (2019). The effects of acoustic variation on the perception of lexical tone in Cantonese-speaking congenital amusics. *Journal of Speech, Language, and Hearing Research*, 62(1), 190-205 is available at https://dx.doi.org/10.1044/2018_JSLHR-H-17-0483. The journal web site is located at <https://pubs.asha.org/journal/jslhr>.

1 **The effects of acoustic variation on the perception of lexical tone in**
2 **Cantonese-speaking congenital amusics**

3
4 Jing Shao^{a,b}, Rebecca Yick Man Lau^a, Phyllis Oi Ching Tang^a, and Caicai Zhang^{a,b,*}

5
6 ^aDepartment of Chinese and Bilingual Studies, The Hong Kong Polytechnic University, Hong
7 Kong SAR, China

8 ^bShenzhen Institutes of Advanced Technology, Chinese Academy of Sciences, Shenzhen, China

9
10 *Corresponding author:

11 Caicai Zhang: Room EF741, Department of Chinese and Bilingual Studies, The Hong Kong
12 Polytechnic University, Hung Hom, Hong Kong SAR, China. Tel: (852) 34008465. Email
13 address: caicai.zhang@polyu.edu.hk.

14
15 Funding Statement

16 This work was supported by grants from the Research Grants Council of Hong Kong (ECS:
17 25603916), the National Natural Science Foundation of China (NSFC: 11504400), and the
18 PolyU Start-up Fund for New Recruits.

19
20 Conflict of Interest Statement

21 The authors have declared that no competing interests existed at the time of publication.

22
23 Running head: Effects of acoustic variation on amusia

1 **Abstract**

2 *Purpose:* Congenital amusia is an inborn neurogenetic disorder of fine-grained pitch processing.

3 This study attempted to pinpoint the impairment mechanism of speech processing in tonal
4 language speakers with amusia. We designed a series of perception tasks aiming at selectively
5 probing low-level pitch processing and relatively high-level phonological processing of lexical
6 tones, with an aim to illuminate the deficiency mechanism underlying tone perception in amusia.

7 *Method:* Sixteen Cantonese-speaking amusics and 16 matched controls were tested on the effects
8 of acoustic (talker/syllable) variations on the identification and discrimination of Cantonese
9 tones in two conditions. In the low variation condition, tones were always associated with the
10 same talker or syllable; in the high variation condition, tones were associated with either
11 different talkers (with the syllable controlled) or different syllables (with the talker controlled).

12 *Results:* Largely similar results were obtained in talker and syllable variation conditions.
13 Amusics exhibited overall poorer performance than controls in tone identification. While
14 amusics also demonstrated poorer performance in tone discrimination, the group difference was
15 more obvious in low variation conditions, where more acoustic constancy was provided. Besides,
16 controls exhibited a greater increase in discrimination sensitivity from high to low variation
17 conditions, implying a stronger benefit of acoustic constancy.

18 *Conclusions:* The findings suggested that amusics' lexical tone perception abilities, in terms of
19 both low-level pitch processing and high-level phonological processing, as measured in low and
20 high variation conditions, are impaired. Importantly, amusics were more impaired in taking
21 advantage of low acoustic variation contexts, and thus less efficiently sharpened their perception
22 of tones when perceptual anchors in talker/syllable were provided, suggesting a possible
23 'anchoring deficit' in congenital amusia.

1

2 **Keywords:** congenital amusia, lexical tone perception, acoustic variation, Cantonese, anchoring

3 deficit.

4

5

1 **Introduction**

2 Congenital amusia (amusia hereafter), also known as tone or tune deafness, is a lifelong
3 disorder of musical pitch processing. It occurs without brain damage and affects about 1.5-4% of
4 the population (Peretz & Hyde, 2003; Peretz & Vuvan, 2017). Individuals with amusia have
5 difficulties detecting mistuned and out-of-key tones (Ayotte, Peretz, & Hyde, 2002). The primary
6 deficit in congenital amusia lies in the processing of the pitch dimension (Foxton Dean, Jennifer
7 L., Gee, Rosemary, Peretz, Isabelle, and Griffiths, Timothy D., 2004; Hyde & Peretz, 2004;
8 Peretz et al., 2002) and impaired short-term memory for pitch (Tillmann, L  v  que, Forni,
9 Albouy, & Caclin, 2016; Tillmann, Schulze, & Foxton, 2009).

10 Empirical evidence has revealed that the pitch deficit in amusia is domain-general and
11 influences speech processing negatively (Patel, Foxton, & Griffiths, 2005). This domain-general
12 pitch-processing deficit hypothesis has been supported by several studies. For example, amusics
13 showed inferior performance on intonation processing and identification of emotion status,
14 where pitch is used extensively as a cue (Jiang et al., 2012; Jiang, Hamm, Lim, Kirk, & Yang,
15 2010; Liu, Patel, Fourcin, & Stewart, 2010; Lu, Ho, Liu, Wu, & Thompson, 2015; Thompson,
16 Marin, & Stewart, 2012).

17

18 *Deficits of lexical tone perception in amusia*

19 In addition to intonation and emotion states, pitch is also used to systematically distinguish
20 word meanings in tonal languages. Lexical tone is of both low-level pitch properties and high-
21 level linguistic properties (Gu, Zhang, Hu, & Zhao, 2013; Zatorre & Gandour, 2008).
22 Investigation on the perception of lexical tone in amusia may reveal different levels of
23 processing and thus are of particular importance in understanding the speech processing

1 mechanism in amusia. Several studies have found that non-tonal language speakers with amusia
2 showed degraded discrimination of non-native tone pairs in Mandarin and Thai, and also inferior
3 performance on the corresponding musical analogs, suggesting a domain-general pitch-
4 processing deficit (Nguyen, Tillmann, Gosselin, & Peretz, 2009; Tillmann et al., 2011). In a
5 similar vein, in tonal languages, previous studies have found that the impairment in lexical tone
6 perception in amusia might be attributed to a domain-general pitch-processing deficit (Liu et al.,
7 2012, 2016; Shao, Zhang, Peng, Yang, & Wang, 2016). To intentionally reduce the demand for
8 phonological processing of lexical tone, Liu et al. (2012) designed tone identification and
9 discrimination tasks in which tones were always associated with the same segments, and words
10 were auditorily presented with corresponding Chinese characters. In contrast to controls,
11 Mandarin-speaking amusics showed impaired performance on word discrimination in natural
12 speech and their gliding analogs, which were complex tones with the same pitch patterns as the
13 speech stimuli. However, they performed as well as controls on word identification in natural
14 speech. It thus indicates that the disorder in amusia is likely to be a domain-general pitch-
15 processing deficit. Although it is possible that native Mandarin speakers may engage
16 phonological processing of lexical tones in word discrimination, they can primarily rely on low-
17 level pitch processing, especially when the tones were constantly carried by the same segments.
18 Altogether, similar results from both non-tonal and tonal language speakers with amusia suggest
19 that the impairment in lexical tone perception may be primarily due to a domain-general pitch-
20 processing deficit in amusia, in that their pitch-processing deficit in the musical domain transfers
21 to the low-level auditory pitch processing in the speech domain.

22 Different from the view that the deficit of amusia in lexical tone perception primarily lies at
23 the low-level auditory pitch processing, some studies have suggested that high-level

1 phonological processing of lexical tones is impaired. For instance, several studies provided
2 evidence that categorical perception of native tones in tonal language speakers with amusia is
3 impaired (Huang, Liu, Dong, & Nan, 2015; Jiang, Hamm, Lim, Kirk, & Yang, 2012; Zhang,
4 Shao, & Huang, 2017). Jiang et al. (2012) found that in contrast to the controls, Mandarin-
5 speaking amusics showed no benefit for between-category tone discriminations, suggesting
6 absence of categorical perception. Similar results were found in Huang et al. (2015), though
7 there were individual differences in that only a subgroup of amusics with most severe tone
8 identification and discrimination impairments showed reduced categorical perception. Zhang et
9 al. (2017) examined the categorical perception by Cantonese-speaking amusics in four stimulus
10 contexts: lexical tone, pure tone, vowel, and voice onset time (VOT). Amusics performed
11 consistently worse than controls in the discrimination of frequency-based stimuli (lexical tone,
12 pure tone and vowel), but preserved the ability of discriminating duration differences (VOT).
13 Interestingly, there were differences in the impairment of frequency/spectral discrimination in
14 speech and nonspeech contexts. Amusics exhibited less benefit in between-category
15 discriminations than controls in speech contexts (lexical tone and vowel), suggesting reduced
16 categorical perception; on the other hand, they performed inferiorly compared to controls
17 regardless of between- and within-category discriminations in the nonspeech context (pure tone),
18 suggesting impaired general auditory pitch processing. These differences imply that amusia may
19 have a deficit of higher-level phonological processing in speech sounds, and also a deficit of
20 lower-level auditory pitch processing in nonspeech sounds.

21 Further corroborating the phonological processing deficit hypothesis, several studies have
22 reported reduced phonological awareness in amusics (Jones, Lucker, Zalewski, Brewer, &
23 Drayna, 2009; Sun, Lu, Ho, & Thompson, 2017a). Jones et al. (2009) examined phonological

1 processing ability in amusics and controls in several tests of phonological processing, including
2 auditory word discrimination, syllable segmentation, and the Comprehensive Test of
3 Phonological Processing (CTOPP). The results showed that amusics demonstrated lower
4 phonological and phonemic awareness abilities on all measures compared to controls. Sun et al.
5 (2017) found that a subgroup of amusics with severe pitch impairment exhibited significantly
6 worse performance than others in phonological awareness, which was further related to their
7 pitch discrimination abilities.

8 Taken together, the above studies suggest that there are two main yet not mutually exclusive
9 hypotheses on the deficit of amusia in speech processing. One hypothesis is that the deficit of
10 amusia primarily lies in domain-general pitch processing (Liu et al., 2012, 2016; Nguyen et al.,
11 2009; Tillmann et al., 2011), which is manifested as an impairment in low-level auditory pitch
12 processing of lexical tones in non-tonal as well as tonal language speakers with amusia,
13 especially when the segment variation was controlled. Another hypothesis is that the deficit
14 already prevails to higher-level phonological processing, affecting categorical perception of
15 lexical tones in native tonal language speakers (Huang et al., 2015; Jiang, Hamm, Lim, Kirk, &
16 Yang, 2012; Zhang, Shao, et al., 2017) and phonological awareness (Jones, Lucker, et al., 2009;
17 Sun et al., 2017a).

18 However, due to the scarcity of the studies that directly compare domain-general pitch
19 processing and higher-level phonological processing of lexical tones, the impairment mechanism
20 underlying the deficient tone perception performance in amusia is not well understood. The
21 current study attempted to probe the impairment mechanism of tone perception in tonal language
22 speakers with amusia. To this end, we designed a series of perception tasks aiming at selectively

1 tapping into relatively low- and high-level lexical tone perception, and by doing so aimed to shed
2 light on the nature of the impairment mechanism in amusia.

3

4 *Influence of acoustic variations on speech perception*

5 Syllable and talker variations are two primary sources of variations in speech perception. It
6 has been found that different degrees of syllable and talker variations employed in the stimulus
7 set generate different degrees of task difficulty and can tap into relatively low- and high-level
8 speech processing (Lee, Tao, & Bond, 2008, 2010; Mullennix & Pisoni, 1990; Shu, Peng, &
9 McBride-Chang, 2008). Nan, Sun, & Peretz (2010) examined the performance of Mandarin-
10 speaking amusics and controls in Mandarin tone discrimination, in which half of the tone pairs
11 were associated with the same syllables, while the other half were associated with different
12 syllables. Results revealed that amusics performed similarly to controls in discriminating tone
13 pairs associated with the same syllables, whereas their performance was impaired in the different
14 syllable condition. These results suggested that Mandarin-speaking amusics were particularly
15 impaired in the different syllable condition, where the demand of phonological processing was
16 presumably higher. Indeed, when the base segments were identical, tone discrimination can be
17 conducted primarily through low-level pitch comparison, without necessarily engaging
18 phonological processing of lexical tones. Thus lexical tone perception under this condition may
19 primarily reflect low-level pitch processing. When the base segments were varied, it increased
20 the demand for the extraction of tone categories independently from segmental processing, so
21 that even typical Mandarin speakers found tone identification and discrimination tasks
22 challenging (Lee et al., 2008; Shu et al., 2008). Lexical tone perception under this condition may
23 mainly reveal the high-level phonological processing of lexical tones.

1 On the other hand, it is well known that talker variation influences word recognition and
2 speech perception (Antoniou & Wong, 2015; Bradlow, Nygaard, & Pisoni, 1999; Peng, Zhang,
3 Zheng, Minett, & Wang, 2012; Pisoni, 1993; Wong & Diehl, 2003; Zhang et al., 2016; Zhang,
4 2018; Zhang et al., 2016; Zhang, Peng, & Wang, 2012, 2013; Zhang & Chen, 2016). For
5 example, Cantonese speakers identified Cantonese tones more accurately when the tone stimuli
6 were blocked by the talker than when they were mixed across talkers (Wong & Diehl, 2003).
7 When multi-talker speech stimuli were presented, it requires the normalization of talker variation,
8 including the processing of talker-specific voice characteristics and the computing of mapping
9 from talker-specific acoustic signals to phonological representations, thus increasing the demand
10 for talker voice processing and phonological processing (Mullennix & Pisoni, 1990). Listeners
11 were found to be slower and less accurate in speech recognition when exposed to multiple talkers
12 as opposed to a single talker (Green, Tomiak, & Kuhl, 1997; Mullennix & Pisoni, 1990;
13 Nusbaum & Morin, 1992). High talker variability also has an impact on the memory of words,
14 such that listeners recalled fewer words from multi-talker lists than from single-talker lists
15 (Martin, Mullennix, Pisoni, & Summers, 1989).

16 In sum, when tone stimuli were carried by identical segments or from a single talker, the
17 acoustic variation is limited and lexical tone perception can primarily rely on low-level pitch
18 processing. When the talker/syllable is varied, the task is more demanding, and lexical tone
19 perception may require the extraction of abstract tone representations or the operation of talker
20 normalization, and tap more into the high-level phonological processing of lexical tones. These
21 hypotheses are in line with the theory underneath the high variability perceptual training (HVPT)
22 paradigm in second language (L2) speech learning (Bradlow, 1999; Bradlow Pisoni, 1997; Pruitt,
23 Jenkins, & Strange, 2006), which exposes L2 learners to speech stimuli embedded in multiple

1 phonetic contexts and produced by multiple talkers. HVPT is found to be more effective than
2 low variability training in L2 speech learning, as the high-variation stimuli encourage and
3 stimulate the learners to form more abstract phonological representations (Bradlow, 2008).

4 In light of the aforementioned findings, low and high talker/syllable variation conditions
5 offer an excellent scenario to examine the *low-level pitch processing* versus *higher-level*
6 *phonological processing deficit* in tone perception in amusics. In the current study, we compared
7 Cantonese-speaking amusics and musically intact controls on the effects of low versus high
8 talker/syllable variation conditions, largely following the design of previous studies (Magnuson
9 & Nusbaum, 2007; Nan et al., 2010; Sjerps, Zhang, & Peng, 2017; Wong & Diehl, 2003). In the
10 high variation condition, tone stimuli produced by four different talkers (talker variation) or
11 carried by four different base syllables (syllable variation) were inter-mixed in a single block, for
12 the talker and syllable variation condition respectively. In the low variation condition, the same
13 set of multi-talker/syllable stimuli were presented, but blocked by the talker for the talker
14 variation condition and by the syllable for the syllable variation condition. We hypothesize that
15 lexical tone perception in low and high talker/syllable variation contexts is supported by different
16 processing mechanisms, with tone perception in the low talker/syllable variation context
17 primarily mediated by low-level pitch processing, whereas that in the high talker/syllable
18 variation context primarily driven by higher-level phonological processing. We further
19 hypothesize that the amusics' performance would be impaired in both conditions, and that the
20 high talker/syllable variation conditions, which require more phonological processing, would be
21 particularly challenging for amusics. Since listeners have to extract abstract tonal representations
22 from the highly variable speech stimuli through higher-level phonological operations for tone

1 identification and discrimination, amusics are expected to show greater impairment in the high
2 talker/syllable variation conditions.

3

4 **Method**

5 *Participants*

6 16 congenital amusics and 16 musically intact controls participated in this experiment.
7 Control participants were matched with amusic participants one by one in age, gender, and years
8 of education. All participants were native speakers of Hong Kong Cantonese, right-handed, with
9 no hearing impairment, and no reported history of formal musical training. Amusics and controls
10 were identified using the Montreal Battery of Evaluation of Amusia (MBEA) (Peretz, Champod,
11 & Hyde, 2003). The MBEA consists of six subtests: three of them are pitch-based tests (scale,
12 contour, and interval), two of them are duration-based tests (rhythm and meter), and the last one
13 is a memory test. All amusic participants scored below 71% (Nan et al., 2010) in the global score,
14 which is the mean of all six subtests, whereas all control participants scored higher than 80%.
15 Demographic characteristics of the participants are summarized in Table 1. The experimental
16 procedures were approved by the Human Subjects Ethics Sub-committee of The Hong Kong
17 Polytechnic University. Informed written consent was obtained from participants in compliance
18 with the experiment protocols.

19

20 *Stimuli*

21 Talker variation condition

22 The stimuli were six words contrasting six Cantonese tones on the syllable “/ji/”: high level
23 tone (T1) – /ji55/ (醫 ‘a doctor’), high rising tone (T2) – /ji25/ (椅 ‘chair’), mid level tone (T3) –

1 /ji33/ (意 ‘meaning’), extra low level/low falling tone (T4) – /ji21/ (兒 ‘son’), low rising tone (T5)
2 – /ji23/ (耳 ‘ear’), and low level tone (T6) – /ji22/ (二 ‘two’) (Bauer & Benedict, 1997;
3 Matthews & Yip, 2013). Two female and two male native Cantonese speakers were recorded
4 reading aloud the words in a carrier sentence, 呢個字係 /li55 ko33 tsi22 hei22/ (‘This word is’)
5 for six times. For each word, one clearly produced token was selected and segmented out of the
6 carrier sentence for each talker. All selected words were normalized in duration to 650 ms, which
7 was the mean duration of the selected tokens, and in mean intensity to 70 dB using Praat
8 (Boersma & Weenink, 2014). Figure 1 displays the F0 contours of the six tones produced by the
9 four talkers.

10

11 Syllable variation condition

12 The stimuli in the syllable variation condition were 24 meaningful words contrasting six
13 Cantonese tones on syllables /ji/, /fɛn/, /fu/ and /wɛi/ (see Appendix). These four syllables were
14 selected for the reasons that their syllable structures are relatively simple and that they occupy a
15 large vowel space (e.g., /i/, /u/ and /ɐ/) that could maximize the pitch variation, given the well-
16 established relationship between vowel height and pitch height (Fischer-Jorgensen, 1990; Lehiste,
17 1970). One female native Cantonese speaker was recorded reading aloud these words in a carrier
18 sentence, 呢個字係 /li55 ko33 tsi22 hei22/ (‘This word is’) for six times. For each word, one
19 clearly produced token was selected and segmented out of the carrier sentence. All selected
20 words were normalized in duration to 650 ms and in mean intensity to 70 dB using Praat. Figure
21 2 shows the F0 trajectories of the six tones carried by the four syllables.

22

23 *Procedure*

1 Talker variation condition

2 The same set of stimuli was presented in two conditions, the low variation and high
3 variation condition. The critical difference between the two conditions was that the stimuli from
4 multiple talkers were presented in a blocked-talker manner in the low variation condition and in
5 a mixed-talker manner in the high variation condition (Magnuson & Nusbaum, 2007; Sjerps et
6 al., 2017; Wong & Diehl, 2003). In other words, there was no talker variation in the low
7 variation condition. Each condition included an identification task and a discrimination task. The
8 stimuli were presented using E-prime 2.0.

9 In the low variation condition, the stimuli produced by the four talkers (M01, M02, F01 and
10 F02) were blocked by the talker and presented in four sub-blocks. In the *identification* task, each
11 set of six words (/ji/) from one talker was presented in a sub-block, generating four sub-blocks.
12 Six words in a talker set were repeated twice and presented randomly within the sub-block,
13 generating a total of 48 trials in the identification task. In each trial, a fixation first occurred on
14 the computer screen for 500 ms, followed by the presentation of a 650-ms speech stimulus via
15 headphones. Subjects were instructed to identify the tone of the word by pressing buttons 1-6 on
16 a keyboard. There was a break after each sub-block. Subjects were given a list of words
17 contrasting the six tones beforehand to get familiar with the tone distinctions and to facilitate
18 tone identification in the subsequent experiment. In the *discrimination* task, six words (/ji/) in
19 each talker set were grouped into 15 different tone pairs and six same tone pairs. Each talker set
20 was presented in a sub-block, generating four sub-blocks in total. Within a sub-block, 15
21 different tone pairs were presented twice in forward and reversed order and six same tone pairs
22 were repeated five times, generating equal number of different and same tone pairs, which were
23 intermixed and randomly presented. The total number of trials in the discrimination task was 240

1 (4 talker blocks \times 60 tone pairs). In each trial, a fixation first occurred on the computer screen for
2 500 ms, followed by the presentation of two 650-ms stimuli separated by an inter-stimulus-
3 interval of 500 ms via the headphones. Subjects were instructed to judge whether the two words
4 carried the same tone or different tones by pressing "left arrow" (same) or "right arrow"
5 (different) on a keyboard within 3 seconds. There was a break after every 20 trials within a sub-
6 block and also a break after each sub-block.

7 In the high variation condition, the stimuli produced by the four talkers (M01, M02, F01 and
8 F02) were intermixed in a block. In the *identification* task, the four talker sets (6 tones \times 4 talkers
9 = 24 trials) were repeated twice and presented randomly within a block, generating a total of 48
10 trials. The trial procedure was the same as that in the low variation condition. In the
11 *discrimination* task, two words in each tone pair (same or different) were always from two
12 different talkers. In order to keep the experiment duration short, the four talkers were grouped
13 into two sets, with set A including six talker pairs (M01-M02, F01-M01, M01-F02, F01-F02,
14 M02-F01, and F02-M02), and set B including the same talker pairs in reversed order (M02-M01,
15 M01-F01, F02-M01, F02-F01, F01-M02, and M02-F02). Each talker pair carried 15 different
16 and six same tone pairs. The 15 different tone pairs were repeated twice and six same tone pairs
17 were repeated five times, generating 60 trials for each talker pair. Within each group, half of the
18 participants were randomly assigned to the set A and the other half to the set B. The subjects
19 were informed that in the high variation condition the talkers were always different, and they
20 were asked to ignore the voice difference and pay close attention to the tones. The trial procedure
21 was the same as that described above. There was a break after every 20 trials within a block. The
22 total number of trials was 360 in the discrimination task (6 different talker pairs \times 60 tone pairs).
23 Note that the total number of trials in low and high variation discrimination tasks was not the

1 same. However, this difference is not expected to have much influence on the results, as the
2 variation difference was a within-subjects factor, and all participants did the same tasks.

3 In both identification and discrimination tasks, the stimuli were presented binaurally through
4 earphones to the participants in a soundproof room. The stimuli were presented at a comfortable
5 listening level to the subjects and the volume level was kept constant within an experiment and
6 across all subjects. The presentation order of the identification and discrimination tasks was
7 counterbalanced across the participants. The block order was kept identical between each amusic
8 participant and the accordingly matched control participant. Before each task, a practice block
9 was given to the participants to familiarize them with the procedure. In both identification and
10 discrimination tasks, accuracy and reaction time (RT) were recorded.

11

12 Syllable variation condition

13 Similar to the talker variation condition, the same set of stimuli with syllable variations was
14 presented in low and high variation conditions. Note that there was no syllable variation in the
15 low variation condition. Each condition included an identification task and a discrimination task.
16 The stimuli were presented using E-prime 2.0.

17 In the low variation condition, tone stimuli carried by the four syllables (/ji/, /fɛn/, /fu/ and
18 /wɛi/) were blocked by the syllable and presented in four sub-blocks. In the *identification* task,
19 each syllable set was presented in a sub-block, in which six words were repeated twice and
20 presented randomly. The trial procedure was the same as that described above. In the
21 *discrimination* task, six words in each syllable set were grouped into 15 different tone pairs and
22 six same tone pairs. Each syllable set was presented in a sub-block, in which 15 different tone
23 pairs were presented twice in forward and reversed order and six same tone pairs were repeated

1 five times, generating equal number of different and same tone pairs, which were intermixed and
2 randomly presented. The trial procedure was the same as that described above.

3 In the high variation condition, tone stimuli carried by the four syllables (/ji/, /fɛn/, /fu/ and
4 /wɛi/) were intermixed in a block. In the *identification* task, the four syllable sets (6 tones × 4
5 syllables = 24 trials) were repeated twice and presented randomly within a block. In the
6 discrimination task, each tone pair (same or different) was always associated with different
7 syllables. In order to keep the experiment duration short, the four syllables were grouped into
8 two sets, with set A including six syllable pairs (/fɛn/-/ji/, /ji/-/wɛi/, /fu/-/wɛi/, /fɛn/-/wai/, /fɛn/-
9 /fu/ and /fu/-/ji/) and set B including the same syllable pairs in reversed order (/ji/-/fɛn/, /wɛi/-/ji/,
10 /wɛi/-/fu/, /wɛi/-/fan/, /fu/-/fɛn/ and /ji/-/fu/). Each syllable pair carried 15 different and six same
11 tone pairs. The 15 different tone pairs were repeated twice and six same tone pairs were repeated
12 five times, generating a total of 60 trials for each syllable pair (6 different syllable pairs × 60
13 tone pairs). Within each group, half of the participants were randomly assigned to the set A and
14 the other half to the set B. The subjects were also informed that in the high variation condition
15 the base syllables were always different, and they were asked to ignore the syllable difference
16 and pay close attention to the lexical tones in tone discrimination.

17 In both the identification and discrimination tasks, the stimuli were presented binaurally
18 through earphones to the participants in a soundproof room. The stimuli were presented at a
19 comfortable listening level to the subjects and the volume level was kept constant within an
20 experiment and across all subjects. The presentation order of the identification and
21 discrimination tasks was counterbalanced across the participants, and kept identical between
22 matched amusic and control participants. Before each task, a practice block was given to the
23 participants to familiarize them with the procedure. In all tasks, accuracy and RT were recorded.

1

2 *Data analysis*

3 For the identification tasks, accuracy and RT were analyzed. Response to each trial was
4 coded as 1 or 0 (correct or incorrect) for each participant. To compare the accuracy of amusics
5 and controls, generalized mixed-effects models were fitted on the responses to each trial (1 or 0)
6 with *group* (amusics and controls) and *variation* (low talker/syllable variation and high
7 talker/syllable variation) as two fixed effects, and with by-subject random intercept and slope as
8 random effects; two-way interaction was also included as a fixed effect in the models. In order to
9 test the significance of fixed effects, a simple model (m_0) with only the intercept as a factor was
10 first fitted, and the factors *group* and *variation* were added to the model consecutively. The
11 model with a fixed effect (e.g., *group*) was compared with a baseline model without it. Models
12 were compared by likelihood ratio tests and p-values were obtained from those tests.

13 RT was measured from the offset of the stimuli to the time that a response was made. RT
14 was analysed in addition to accuracy, for the reason that it provides an extra index to measure the
15 cognitive effort used in the speech perception task, by revealing subtle temporal differences
16 between different listener groups or different listening conditions. When the accuracy among the
17 listeners groups is comparable, shorter RT usually suggests more rapid, robust and automatic
18 speech processing and implies less cognitive effort involved in the task (Strange, 2011). Thus,
19 RT was calculated in the current study to provide another dimension for evaluating the
20 perceptual ability of amusics and controls. All trials were included in the RT analysis,
21 irrespective of whether the response was correct or not. Incorrect trials were not discarded for the
22 reasons that the accuracy was quite low in the identification tasks (41%-65%) and not evenly
23 distributed across different conditions. Excluding incorrect trials therefore leads to the removal

1 of a significant portion of the data and imbalanced number of trials across different conditions.
2 Statistical analyses had been conducted on the RT of correct and incorrect trials separately, and
3 incorrect trials yielded either similar effects as correct trials or no significant effect (perhaps due
4 to small number of trials) (please see the supplemental files S1 and S2). The results of RT
5 analyses combining correct and incorrect trials are reported here. Linear mixed-effects models
6 were fitted on the log-transformed RT data with *group* and *variation* as two fixed effects, and
7 with by-subject random intercept and slope as random effects; two-way interaction was also
8 included as a fixed effect in the models. Models were compared by likelihood ratio tests and p-
9 values were obtained from those tests. The above two sets of analyses were performed with R (R
10 Core Team, 2014), using the *lme4* package (Bates, Maechler, & Bolker, 2012), the *lmtest*
11 package (Zeileis & Hothorn, 2002) and the *lsmeans* package (Lenth, 2016).

12 For the discrimination tasks, the sensitivity index d' (Macmillan & Creelman, 2005) and RT
13 were analyzed. The discrimination responses were analyzed in terms of d' rather than accuracy
14 for the reason that d' takes into consideration the response bias. Following previous studies on
15 speech discrimination (Zhang, 2018; Zhang, Shao, et al., 2017), the d' was computed as the z-
16 score value of the hit rate ("different" responses to different tone pairs) minus that of the false
17 alarm rate ("different" responses to same tone pairs) for each tone pair per subject. Note that this
18 formula for calculating d' ($z(H)-z(F)$) follows the assumption of an independent-observation
19 strategy in fixed designs (Macmillan & Creelman, 2005). Nonetheless, the tone pairs were
20 intermixed and randomly presented in a block in the current study, which is not unlike a roving
21 design. Macmillan & Creelman (2005) suggested that in roving designs, listeners might adopt a
22 differencing strategy, applying a threshold of difference to decide if two stimuli are different
23 enough to be judged as different. But since the stimuli were all speech materials with

1 recognizable categories in the current study, there is reason to assume that listeners adopted the
2 independent-observation strategy, which justifies the calculation of d' as $z(H)-z(F)$. Trials with
3 null responses were disregarded (talker variation condition: 0.5%; syllable variation condition:
4 2.6%). *Group* \times *variation* repeated measures ANOVAs were conducted on the d' scores, using
5 the Statistical Package for the Social Sciences (SPSS) (SPSS, 2011). Corrections for violations
6 of sphericity were made, where appropriate, using the Greenhouse-Geisser method. Tukey's
7 HSD post-hoc test was applied to make pairwise comparisons when necessary. All effects were
8 reported as significant at $p < .05$.

9 In the discrimination tasks, RT was measured from the offset of the second stimulus in a pair
10 to the time that a response was made. All trials with a response made within the time limit were
11 included in the RT analysis, irrespective of whether the response was correct or not. Linear
12 mixed-effects models were fitted on the log-transformed RT data with *group* and *variation* as
13 two fixed effects, and with by-subject random intercept and slope as random effects, following
14 the procedure described above. The results of RT analyses separated by correct and incorrect
15 trials are reported in the supplemental files S1 and S2.

16

17 **Results**

18 *Talker variation condition*

19 Figure 3(a) illustrates the identification accuracy in each condition for the amusic and
20 control group. Generalized mixed-effects model revealed a significant main effect of *group* ($\chi^2(1)$
21 $= 7.88, p = .004$). Amusics showed significantly lower accuracy compared to controls (amusics:
22 $M = 0.41, SD = 0.22$; controls: $M = 0.61, SD = 0.18$). No other effects were significant. Although
23 the two-way interaction was not significant ($\chi^2(1) = 0.60, p = .438$), there was a trend of larger

1 group differences in the low variation condition (amusics: $M = 0.41$, $SD = 0.24$; controls: $M =$
2 0.63 , $SD = 0.18$) compared to the high variation condition (amusics: $M = 0.41$, $SD = 0.20$;
3 controls: $M = 0.60$, $SD = 0.18$). Furthermore, controls exhibited a trend of benefiting more from
4 low talker variations (low variation: $M = 0.63$, $SD = 0.18$; high variation: $M = 0.60$, $SD = 0.18$)
5 compared to amusics (low variation: amusics: $M = 0.41$, $SD = 0.24$; high variation: amusics: $M =$
6 0.41 , $SD = 0.20$) in tone identification, though the difference was small.

7 Figure 3(b) displays the RT in the identification task. No significant effects were found.

8 Figure 3(c) shows the discrimination sensitivity d' . *Group* \times *variation* repeated-measures
9 ANOVA found significant main effects of *group* ($F(1, 30) = 24.413$, $p < .001$, $\eta_p^2 = 0.449$) and
10 *variation* ($F(1, 30) = 119.974$, $p < .001$, $\eta_p^2 = 0.800$) and a significant two-way interaction ($F(1,$
11 $30) = 10.453$, $p = .003$, $\eta_p^2 = 0.258$). Independent-samples t-tests showed that amusics exhibited
12 significantly lower d' than controls in both low variation ($t(30) = -4.762$, $p < .001$, $d = 1.684$) and
13 high variation conditions ($t(30) = -3.430$, $p = .002$, $d = 1.212$), but the group difference was
14 larger in the low variation condition. Controls demonstrated significantly higher d' in the low
15 variation condition than in the high variation condition ($t(30) = -6.505$, $p < .001$, $d = 2.299$), with
16 a large effect size. In contrast, while amusics also exhibited larger d' scores in the low variation
17 condition than in the high variation condition ($t(30) = -2.047$, $p = .05$, $d = 0.723$), this effect was
18 only marginally significant. These results indicate that amusics showed inferior performance
19 than controls under both low and high talker variation conditions, but the group difference was
20 larger in the low variation condition. Besides, controls benefited more from low talker variations
21 in tone discrimination.

22 Figure 3(d) shows the RT in the discrimination task. Linear mixed-effects model found a
23 significant main effect of *variation* ($\chi^2(1) = 394$, $p < .001$), where RT in the high talker variation

1 condition was significantly longer than in the low talker variation condition. No other effects
2 were significant.

3

4 *Syllable variation condition*

5 Figure 4(a) illustrates the identification accuracy. Generalized mixed-effects models found a
6 significant main effect of *group* ($\chi^2(1) = 7.17, p < .001$), where amusics showed significantly
7 lower accuracy compared to controls (amusics: $M = 0.44, SD = 0.21$; controls: $M = 0.62, SD =$
8 0.16). No other effects were significant. Although the *group* by *variation* interaction was not
9 significant ($\chi^2(1) = 1.529, p = .216$), the group difference appeared to be larger in the low
10 syllable variation condition (amusics: $M = 0.44, SD = 0.23$; controls: $M = 0.65, SD = 0.18$) than
11 that in the high variation condition (amusics: $M = 0.44, SD = 0.20$; controls: $M = 0.59, SD =$
12 0.13). Furthermore, controls appeared to benefit more from low syllable variations (low variation:
13 $M = 0.65, SD = 0.18$; high variation: $M = 0.59, SD = 0.13$) compared to amusics (low variation:
14 $M = 0.44, SD = 0.23$; high variation: $M = 0.44, SD = 0.20$) in tone identification.

15 Figure 4(b) shows the RT in the identification task. Linear mixed-effects model found a
16 significant main effect of *group* ($\chi^2(1) = 4.012, p = .04$), *variation* ($\chi^2(1) = 126.3, p < .001$), and a
17 significant two-way interaction ($\chi^2(1) = 6.036, p = .014$). Pairwise comparisons revealed that in
18 the control group, RT in the high variation condition was significantly longer than that in the low
19 variation condition ($z = 6.198, p < .001$); such effect was also significant in the amusic group, but
20 the difference was smaller ($z = 2.669, p = .01$). Under the high variation condition, controls
21 exhibited significantly longer RT than amusics ($z = -2.255, p = .03$), but such difference cannot
22 be found in the low variation condition ($z = -1.581, p = .124$). These results suggest that the

1 control group was under greater influences of syllable variations in terms of RT compared to the
2 amusic group.

3 Figure 4(c) displays the discrimination sensitivity d' . *Group* \times *condition* repeated-measures
4 ANOVA found significant main effects of *group* ($F(1, 30) = 33.637, p < .001, \eta_p^2 = 0.529$) and
5 *variation* ($F(1, 30) = 133.709, p < .001, \eta_p^2 = 0.817$) and a significant two-way interaction ($F(1,$
6 $30) = 4.926, p = .034, \eta_p^2 = 0.141$). Post hoc analyses showed that although both amusic and
7 control groups exhibited significantly higher d' scores in the low variation condition than in the
8 high variation condition (amusics: $t(30) = -5.462, p < .001, d = 1.931$; controls: $t(30) = -9.017, p$
9 $< .001, d = 3.188$), the effect was larger in the control group. Within each variation condition, the
10 amusic group demonstrated significantly lower d' than controls (high variation: $t(30) = -4.585, p$
11 $< .001, d = 1.621$; low variation: $t(30) = -4.595, p < .001, d = 1.624$). These results demonstrate
12 that controls benefited more from low syllable variation compared to amusics.

13 Figure 4(d) shows RT in the discrimination task. Linear mixed-effects model found a
14 significant interaction between *group* and *condition* ($\chi^2(1) = 196.16, p < .001$). Pairwise
15 comparisons revealed that in the high variation condition, RT in the control group was
16 significantly longer than the amusic group ($z = -2.187, p = .03$); in the low variation condition,
17 RT in the two groups was similar ($z = -0.142, p = .88$). Within each group, RT in the high
18 variation condition was significantly longer than that in the low variation condition, but the
19 effect was more evident in the control group ($z = 32.886, p < .001$) than in the amusic group ($z =$
20 $12.983, p < .001$). These RT results suggest that controls tended to show greater benefit from low
21 syllable variation, where they employed significantly shorter RT in the low variation condition
22 than in the high variation condition.

1 To summarize, the amusics performed worse than the controls in both tone identification
2 and discrimination tasks. In the talker variation condition, there was a significant group effect on
3 the identification accuracy, where controls outperformed than amusics. Although the two-way
4 interaction was not significant, there were trends that controls exhibited a greater increase of
5 accuracy from high to low variation conditions. In the discrimination task, there was a significant
6 *group* by *variation* interaction effect on the *d'* score, where only the controls demonstrated
7 significantly higher *d'* in the low variation condition than in the high variation condition, and the
8 group difference was more evident in the low variation condition. The discrimination RT mainly
9 revealed an effect of variation, where RT in the low variation condition was shorter than the high
10 variation condition. In the syllable variation condition, the perceptual patterns were similar,
11 except that the RT data also revealed significant *group* by *variation* interactions. In the
12 identification task, controls exhibited significantly shorter RT in the low variation condition than
13 in the high variation condition, but this effect was smaller in amusics. Similarly, in the
14 discrimination task, controls again tended to show greater benefit from the low variation
15 condition, by employing shorter RT in the low variation condition than in the high variation
16 condition.

17

18 **Discussion**

19 While amusia has been consistently reported to influence tone perception negatively (Liu et
20 al., 2016; Nan et al., 2010; Shao et al., 2016), the mechanism underlying the deficient tone
21 perception is not well understood. In this study, we examined this issue through a comparison of
22 two conditions, *low* talker/syllable variation versus *high* talker/syllable variation. In the low
23 variation condition, only one talker/syllable was associated with the tones, and acoustic

1 variations among the stimuli were low. It was hypothesized that tone perception in this condition
2 primarily tapped into relatively low-level pitch processing. In the high variation condition, tones
3 were always associated with different talkers or syllables, meaning that listeners had to extract
4 abstract tonal representations out of speech signals with large acoustic variations. When
5 associated with different talkers, tone perception additionally involved talker voice processing. It
6 was hypothesized that tone perception in this condition primarily tapped into relatively high-
7 level phonological processing of lexical tones. It was predicted that amusics would show greater
8 impairment in the high talker/syllable variation conditions given the hypothesized phonological
9 processing deficit.

10 We found that the general patterns in talker and syllable variation conditions were largely
11 similar. Amusics exhibited overall lower accuracy than controls in tone identification tasks. In
12 tone discrimination tasks, there were significant *group* by *variation* interactions, in that the
13 controls exhibited a greater increase of d' scores from high to low variation conditions, and the
14 group difference was larger in the low variation conditions. The results deviated somewhat from
15 the predicted patterns of greater group differences in the high variation conditions. In the text
16 below, we first discussed the observed impairment of amusics in low and high variation
17 conditions, respectively, in relation to the hypothesized low-level pitch-processing and high-level
18 phonological processing deficits of amusia. We then tried to explain the larger group differences
19 observed in the low variation conditions, together with an exploration of the difference between
20 identification and discrimination tasks.

21 We found that amusics performed significantly worse than the controls in *low* talker/syllable
22 variation conditions, confirming that amusics are generally impaired in tone perception,
23 compatible with the *low-level* pitch-processing deficit. These results are consistent with previous

1 findings of impaired tone perception in tonal and non-tonal language speakers with amusia when
2 the segments were fixed (Liu et al., 2016; Nguyen et al., 2009; Shao et al., 2016; Tillmann et al.,
3 2011). In Liu et al. (2016), four Cantonese tone pairs were selected and presented to amusics and
4 controls. The four tone pairs were always carried by the same syllables, which corresponds to the
5 low syllable variation condition in the current study. Liu and colleagues found that Cantonese-
6 speaking amusics demonstrated impaired tone perception compared to controls. The current
7 results also echo with Shao et al. (2016), where Cantonese-speaking amusics showed deficient
8 tone perception when the syllable variation was low.

9 However, the results that amusics were impaired in tone discrimination when the tones were
10 carried by the same syllables are inconsistent with those reported in Nan et al. (2010). Nan et al.
11 (2010) found that Mandarin-speaking amusics did not score significantly below the controls in
12 tone discrimination when the carrying syllables were identical, but they did exhibit poorer
13 performance when the syllables were different. This discrepancy between Cantonese- and
14 Mandarin-speaking amusics may be due to the different tonal systems of the two languages.
15 Each of the four tones in Mandarin has a unique F0 direction, which may facilitate successful
16 tone discrimination (Gandour, 1983; Gandour & Harshman, 1978; Peng et al., 2012). In contrast,
17 the tonal space of Cantonese is much denser (Bauer & Benedict, 1997; Matthews & Yip, 2013;
18 Peng, 2006; Peng et al., 2012). Discrimination of the six tones in Cantonese relies on both F0
19 level and F0 direction, which increases difficulties in tone discrimination (Gandour, 1983;
20 Gandour & Harshman, 1978; Peng et al., 2012). It is thus not surprising that Mandarin-speaking
21 amusics showed comparable performance to controls when the same syllables were associated to
22 the tones. In contrast, Cantonese-speaking amusics demonstrated poorer performance in the low
23 variation condition.

1 Our results also revealed that Cantonese-speaking amusics showed significantly worse tone
2 perception performance than controls in *high* talker/syllable variation conditions, suggesting a
3 deficit in *high-level* phonological processing of lexical tones. Tones carried by different
4 talkers/syllables were more acoustically different, increasing the difficulty of reliably extracting
5 tone categories from acoustical signals and the cognitive load (Johnson, 2005). Perception of
6 tones carried by different talkers/syllables involves phonological processing, i.e., extracting
7 abstract phonological representations of tones from different talkers and syllables, in addition to
8 talker voice processing (talker condition) and segmental processing (syllable condition). High
9 talker/syllable variations could decrease the “quality” of extracted tone categories, affecting tone
10 perception accuracy. The worse performance of amusics in the high variation condition therefore
11 may suggest a deficit in high-level phonological processing of lexical tones. These findings are
12 in line with previous studies which reported that categorical perception of native tones is
13 impaired in Chinese amusics (Huang et al., 2015; Jiang, Hamm, Lim, Kirk, & Yang, 2012;
14 Zhang et al., 2017). They also coincide partly with the findings of Wang & Peng (2014). Wang
15 & Peng (2014) indicated that Mandarin-speaking controls discriminated non-native Cantonese
16 level tones more accurately when the tones were carried by native Mandarin syllables than by
17 Cantonese syllables, whereas amusics failed to show such benefit, suggesting a possible deficit
18 of amusics in benefiting from native syllables when processing non-native tones. Altogether,
19 these findings indicate that the deficit of amusics is not confined to auditory pitch discrimination,
20 but prevails to high-level phonological processing of lexical tones.

21 The impairment of amusia in phonological processing is also somewhat compatible with the
22 observation of poor phonological awareness in amusia. Several studies have indicated that
23 amusics show reduced phonemic awareness (Jones, Lucker, et al., 2009; Sun et al., 2017a),

1 which may account for their impoverished performance in high variation conditions in the
2 current study, where much phonological decoding was demanded. When breaking down the
3 structure of the syllables to extract phonological representations of tones independently from
4 representations of segments, the amusics may encounter greater difficulties compared to controls,
5 leading to worse tone perception performance.

6 Note that the high talker variation conditions may not only involve phonological processing,
7 but also talker voice processing (Johnson, 2005). Worse tone perception performance of amusics
8 in these conditions may imply an additional talker-processing deficit in amusia. It has been found
9 that encoding an acoustic signal into a phonological representation is closely related to the
10 process of encoding a talker's voice. Talker and phonetic information can not be selectively
11 ignored when participants are instructed to attend to a specific dimension of the speech stimuli,
12 suggesting that talker and phonetic information is processed in an integral fashion (Kaganovich,
13 Francis, & Melara, 2006; Mullennix & Pisoni, 1990; Perrachione, Pierrehumbert, & Wong, 2009;
14 Zhang et al., 2016). The current findings that amusics were impaired in tone perception in high
15 talker variation conditions may be partially affected by a deficit in talker processing, in addition
16 to the high-level phonological processing deficit.

17 It should be noted that we hypothesized that the group difference would be more
18 pronounced in the high variation conditions, which were phonologically more demanding tasks.
19 However, we observed the opposite patterns in that the group difference was generally larger in
20 the low variation conditions, where the talker/syllable was consistently repeated and the acoustic
21 variation was limited. This pattern was most noticeable in the discrimination tasks. For the
22 identification accuracy, there was a similar trend of larger group differences in the low variation
23 conditions but it was not statistically significant. This is probably because in the identification

1 task, there was no talker/syllable variability within a trial, unlike in the discrimination task where
2 listeners were confronted with two stimuli with different talkers or syllables when discriminating
3 the tone differences. Although there was talker/syllable variability across trials in the
4 identification task, this may somehow have limited the influence of talker/syllable variability.
5 This may explain why there was only an insignificant trend in the identification task.
6 Nonetheless, the identification RT did reveal a significant interaction of *group* and *variation* in
7 the syllable condition, in that the control group tended to exhibit significantly shorter RT in the
8 low variation condition than in the high variation condition, but the effect was reduced in
9 amusics. RT thus provides an additional index to evaluate the perceptual performance of amusics,
10 complementing the identification accuracy. Overall, these findings suggest that controls
11 benefited more from low acoustic variations, exhibiting more accurate and rapid tone perception,
12 especially when there was within-trial talker/syllable variability (i.e., in discrimination tasks).

13 The finding that amusics were more impaired in the low variation condition can be
14 explained by the ‘anchoring theory’ which is originally proposed to account for the phonological
15 deficit in dyslexia (Ahissar, 2007; Ahissar, Lubin, Putter-Katz, & Banai, 2006). This hypothesis
16 states that the deficit of dyslexics lies in the dynamics that link perception with perceptual
17 memory through the implicit formation of perceptual anchors. These anchors may guide the
18 perception of subsequent stimuli, help retain and retrieve the incoming stimuli. Empirical
19 evidence showed that typical listeners tune to, or ‘anchor to’, incoming stimuli quickly. When
20 the stimuli are subsequently repeated, they can react more accurately with the ‘perceptual
21 anchors’. On the contrary, individuals with dyslexia failed to benefit from such repetitions. For
22 example, Ahissar et al. (2006) investigated nonspeech tone discrimination by individuals with
23 dyslexia. In one condition, a standard tone was constantly presented within each trial; in another

1 condition, no standard tone was presented. The results showed that the dyslexics performed
2 poorly in the condition where the standard tone was presented, while their performance on the
3 difficult ‘no standard’ condition did not differ from the controls. Moreover, the dyslexics’
4 performance did not differ between the two conditions; in contrast, the controls’ performance in
5 the standard condition was considerably better than in the ‘no standard’ condition. It was
6 proposed that in the standard condition, listeners could form a perceptual anchor based on the
7 repeated standard tones, and thus could make judgment quickly and accurately. Individuals with
8 dyslexia failed to form such an ‘anchor’ and showed degraded performance in the standard
9 condition. Altogether, it suggests that dyslexic individuals have difficulties in dynamically
10 constructing “perceptual anchors”, deriving from deficient adaptation mechanisms.

11 Results of the current study can also be accounted for by the ‘anchoring theory’. In the low
12 variation condition, tone pairs were always associated with the same syllable and talker.
13 Anchoring to the consistent syllable or talker identity in the speech stimuli probably facilitates
14 tone perception. Controls benefited from the consistency in syllable and talker identity,
15 exhibiting higher accuracy and shorter RT in tone perception in the low variation condition than
16 in the high variation condition. In contrast to controls, the performance of amusics improved less
17 when the same talker or syllable was presented across trials, indicating that their ability to build
18 and tune to an internal talker/syllable anchor might be impoverished. Therefore, without a
19 properly functioning anchoring mechanism, the perceptual system of amusics may be less
20 resilient to external variation, similar to the case of dyslexics. This explains why under
21 conditions with low talker/syllable variation, the perception of controls greatly sharpened
22 compared with the high talker/syllable variation conditions, whereas the perception of amusics
23 showed less improvement. The hypothesis of a potential ‘anchoring deficit’ in amusia is also

1 consistent with the result of a previous functional MRI study that the brain of amusics showed
2 abnormally strong activations towards repeated pitch stimuli, while the activation in the controls’
3 brain attenuated (Zhang, Peng, Shao, & Wang, 2017).

4 The fact that amusics might exhibit a similar ‘anchoring deficit’ to individuals with dyslexia
5 could be partly tracked down to the attention deficit and impaired short-term memory. For
6 example, Hari & Renvall (2001) indicated that the core deficit of dyslexia involves poor
7 attentional skills. Ahissar (2007) suggested that the anchoring deficit hypothesis in dyslexia can
8 be viewed as a type of attention deficit hypothesis. Furthermore, the short-term memory is found
9 to be closely related to the ‘anchoring’ performance in dyslexics (Ahissar et al., 2006). Amusics,
10 on the other hand, are also found to display impaired short-term memory (Albouy et al., 2013;
11 Gosselin, Jolicœur, & Peretz, 2009; Tillmann et al., 2016; Williamson & Stewart, 2010).
12 Approximately 40% of the amusic individuals are found to have an attention deficit (Jones,
13 Zalewski, Brewer, Lucker, & Drayna, 2009). With poor working memory and attention deficits,
14 amusics may have greater difficulty benefiting from the clean and variation-free conditions.
15 Altogether, it is likely that the potential anchoring deficit in amusia contributes to the
16 impoverished lexical tone perception performance in the low acoustic variation conditions,
17 beyond their pitch-processing deficit.

18 It is worth noting that while the results were largely similar between the talker and syllable
19 variation conditions, some differences can be observed, especially in the RT results. In general,
20 significant group differences or *group by condition* interactions were found in the RT of the
21 syllable variation condition, but not in the talker variation condition, implying a stronger impact
22 of syllable variations on revealing the disparity between amusics and controls in response time
23 during tone perception. It is possible that syllables carried more linguistic information, and when

1 the syllables varied it was more distracting than when the same syllable varied in different
2 talkers' voices. As a result, syllable variations may place a greater demand on the anchoring
3 system with a stronger influence on tone perception.

4 Lastly, the tone perception impairment in amusics appears to be superficially similar to that
5 of cochlear implant (CI) users, who also exhibit tone recognition deficiencies (Chang, Chang,
6 Lin, & Luo, 2016). It has been found that Mandarin-speaking CI users performed worse in the
7 mixed-talker condition than in the blocked-talker condition in tone and vowel recognition
8 (Chang et al., 2016). While our results suggested that amusics only showed marginally better
9 performance in the blocked-talker condition than the mixed-talker condition in Cantonese tone
10 perception, suggesting that they failed to benefit as much from acoustic constancy in the
11 blocked-talker condition. This discrepancy suggests that the deficiency mechanism of amusia is
12 likely to be different from that the CI users, who appeared to have a normal anchoring
13 mechanism, and impaired low-level pitch processing that is similar to the amusics. It also implies
14 that the deficit of amusia is more cortically originated (Liu, Maggu, Lau, & Wong, 2015). Future
15 studies can further explore the potentially different deficiency mechanisms of amusia and CI
16 users.

17 To conclude, we found that Cantonese-speaking amusics were impaired in both low-level
18 pitch processing and higher-level phonological processing of lexical tones. Furthermore, amusics
19 were more impaired than controls in low variation contexts, where acoustic constancy in
20 syllables/talkers can be used to construct perceptual anchors. These findings shed some light on
21 the mechanism of the deficient tone perception in amusia. We propose that it may reflect
22 impaired dynamics of the stimulus anchoring mechanism, suggesting a possible 'anchoring
23 deficit' in amusia. A few unsolved issues arise from the current study, which can be further

1 examined in future studies. First, to explicitly test the anchoring deficit hypothesis in amusics,
2 future studies can examine the performance of amusics using study designs identical to those in
3 previous studies of the “anchoring theory” (Ahissar, 2007; Ahissar et al., 2006), together with
4 tests of attention and short-term memory. Second, while the current study focused on impaired
5 speech processing (lexical tone perception) in tonal language speakers with amusia, it is possible
6 that the potential anchoring deficit revealed in lexical tone perception (to be further verified) is
7 domain-general. Analogous to the manipulation of syllable and talker variations in the current
8 study, future studies may examine the effect of timbre variation (e.g., piano, violin, etc.) on pitch
9 processing in amusical speakers of a non-tonal language. Third, it is worth examining whether
10 the amusics’ ability to develop a perceptual anchor would improve with more exposure to the
11 speech stimuli. Fourth, it is possible the potential anchoring deficit in amusics was manifested
12 most strongly when the task was demanding and the stimulus set was highly variable and
13 complex. In current study, there were four base syllables and four talkers’ voices; in the high
14 variation discrimination tasks, tones carried by different base syllables or different talkers’
15 voices were always paired together to maximize the variability. The task difficulty under the
16 current experiment settings was considerably high. It is possible that in less demanding
17 experiment settings, the group difference might be compressed and the anchoring deficit might
18 be less clearly demonstrated. Future studies may examine the effect of task demand and stimulus
19 variability on the perception of lexical tones in amusia. Fifth, future studies with a larger sample
20 size in each participant group may want to include tone/tone pair as a third factor in statistical
21 analyses, to further examine whether certain tone or tone pairs may suffer more from high
22 talker/syllable variations. Sixth, some aspects of stimulus presentation in the high variation
23 conditions can be improved. For example, stimulus presentation may be better controlled by not

1 presenting stimuli with the same talker/syllable in consecutive trials, in order to maximize the
2 variability across trials. Furthermore, the high variation conditions may be broken down into four
3 sub-blocks, to be more comparable to the low variation conditions. Last, the number of trials in
4 the high variation discrimination tasks was larger than that in the low variation discrimination
5 tasks, which may have increased fatigue. Although this presumably did not affect the comparison
6 between amusics and controls, future studies may try to equalize the number of trials by
7 randomly selecting a subset of four talker/syllable pairs in the high variation discrimination tasks.

8

9 **Acknowledgements**

10 We thank the anonymous reviewers for constructive comments. We thank Mr. Yubin Zhang
11 for help with data analysis. The data of the syllable variation condition were presented at the
12 Speech Prosody 2018 in Poznań, Poland, during 13-16 June 2018.

13

14 **References**

15 Ahissar, M. (2007). Dyslexia and the anchoring-deficit hypothesis. *Trends in Cognitive Sciences*,
16 11(11), 458–465. <http://doi.org/10.1016/j.tics.2007.08.015>

17 Ahissar, M., Lubin, Y., Putter-Katz, H., & Banai, K. (2006). Dyslexia and the failure to form a
18 perceptual anchor. *Nature Neuroscience*, 9(12), 1558–1564. <http://doi.org/10.1038/nn1800>

19 Albouy, P., Mattout, J., Bouet, R., Maby, E., Sanchez, G., Aguera, P.-E., ... Tillmann, B. (2013).

20 Impaired pitch perception and memory in congenital amusia: The deficit starts in the
21 auditory cortex. *Brain*, 136(5), 1639–1661. <http://doi.org/10.1093/brain/awt082>

22 Antoniou, M., & Wong, P. C. M. (2015). Poor phonetic perceivers are affected by cognitive load
23 when resolving talker variability; Poor phonetic perceivers are affected by cognitive load

1 when resolving talker variability. *Journal of Acoustical Society of America*, 138(2), 571–
2 574. <http://doi.org/10.1121/1.4923362>

3 Ayotte, J., Peretz, I., & Hyde, K. (2002). Congenital amusia: A group study of adults afflicted
4 with a music-specific disorder. *Brain*, 125(2), 238–251. <http://doi.org/10.1093/brain/awf028>

5 Bates, D., Maechler, M., & Bolker, B. (2012). lme4: linear mixed-effects models using S4
6 classes. R package version 0.999375-42. 2011.

7 Bauer, R., & Benedict, P. K. (1997). *Modern cantonese phonology*. Berlin: Mouton de Gruyter.

8 Boersma, P., & Weenink, D. (2014). Praat: Doing phonetics by computer.

9 Bradlow, A. (1999). Training Japanese listeners to identify English /r/and /l/: Long-term
10 retention of learning in perception and production. *Perception & Psychophysics*.

11 Bradlow, A. (2008). Training non-native language sound patterns:Lessons from training
12 Japanese adults on the English/p/–/l/ contrast.

13 Bradlow, A. R., Nygaard, L. C., & Pisoni, D. B. (1999). Effects of talker, rate, and amplitude
14 variation on recognition memory for spoken words. *Perception & Psychophysics*, 61(2),
15 206–219. <http://doi.org/10.3758/BF03206883>

16 Bradlow Pisoni, A. (1997). Training Japanese listeners to identify English /r/ and /l/: IV.ome
17 effects of perceptual learning on speech production. *JASA*.

18 Chang, Y., Chang, R. Y., Lin, C.-Y., & Luo, X. (2016). Mandarin tone and vowel recognition in
19 cochlear implant users: Effects of talker variability and bimodal hearing. *Ear and Hearing*,
20 37(3), 271.

21 Fischer-Jorgensen, E. (1990). Intrinsic F0 in tense and lax vowels with special reference to
22 German. *Phonetica*, 47(3–4), 99–140.

23 Foxton Dean, Jennifer L., Gee, Rosemary, Peretz, Isabelle, and Griffiths, Timothy D., J. M.

- 1 (2004). Characterization of deficits in pitch perception underlying “tone deafness.” *Brain*,
2 127, 801–810. Journal Article.
- 3 Gandour, J. (1983). Tone perception in Far Eastern languages. *J Phon*, 11, 49–175.
- 4 Gandour, J., & Harshman, R. (1978). Cross-language difference in tone perception: A
5 multidimensional scaling investigation. *Lang Speech*, 21, 1–33.
- 6 Gosselin, N., Jolicœur, P., & Peretz, I. (2009). Impaired memory for pitch in congenital amusia.
7 *Annals of the New York Academy of Sciences*, 1169(1), 270–272. JOUR.
8 <http://doi.org/10.1111/j.1749-6632.2009.04762.x>
- 9 Green, K. P., Tomiak, G. R., & Kuhl, P. K. (1997). The encoding of rate and talker information
10 during phonetic perception. *Attention, Perception, & Psychophysics*, 59(5), 675–692.
- 11 Gu, F., Zhang, C., Hu, A., & Zhao, G. (2013). Left hemisphere lateralization for lexical and
12 acoustic pitch processing in Cantonese speakers as revealed by mismatch negativity.
13 *NeuroImage*, 83(0), 637–645. Journal Article.
14 <http://doi.org/http://dx.doi.org/10.1016/j.neuroimage.2013.02.080>
- 15 Hari, R., & Renvall, H. (2001). Impaired processing of rapid stimulus sequences in dyslexia.
16 *Trends in Cognitive Sciences*, 5(12), 525–532. [http://doi.org/10.1016/S1364-](http://doi.org/10.1016/S1364-6613(00)01801-5)
17 [6613\(00\)01801-5](http://doi.org/10.1016/S1364-6613(00)01801-5)
- 18 Huang, W. T., Liu, C., Dong, Q., & Nan, Y. (2015). *Categorical perception of lexical tones in*
19 *Mandarin-speaking congenital amusics. Frontiers in Psychology*.
20 <http://doi.org/10.3389/fpsyg.2015.00829>
- 21 Hyde, K. L., & Peretz, I. (2004). Brains that are out of tune but in time. *Psychological Science*,
22 15(5), 356–360. <http://doi.org/10.1111/j.0956-7976.2004.00683.x>
- 23 Jiang, C., Hamm, J. P., Lim, V. K., Kirk, I. J., Chen, X., & Yang, Y. (2012). Amusia results in

1 abnormal brain activity following inappropriate intonation during speech comprehension.
2 *PLoS ONE*, 7(7), e41411. <http://doi.org/10.1371/journal.pone.0041411>

3 Jiang, C., Hamm, J. P., Lim, V. K., Kirk, I. J., & Yang, Y. (2010). Processing melodic contour
4 and speech intonation in congenital amusics with Mandarin Chinese. *Neuropsychologia*,
5 48(9), 2630–2639. <http://doi.org/http://dx.doi.org/10.1016/j.neuropsychologia.2010.05.009>

6 Jiang, C., Hamm, J. P., Lim, V. K., Kirk, I. J., & Yang, Y. (2012). Impaired categorical
7 perception of lexical tones in Mandarin-speaking congenital amusics. *Memory & Cognition*,
8 40(7), 1109–21. <http://doi.org/10.3758/s13421-012-0208-2>

9 Johnson, K. (2005). Speaker normalization in speech perception. In D. B. Pisoni and Remez,
10 Robert E. (Ed.), *The handbook of speech perception* (pp. 363–389). Blackwell Publishing.

11 Jones, J. L., Lucker, J., Zalewski, C., Brewer, C., & Drayna, D. (2009). Phonological processing
12 in adults with deficits in musical pitch recognition. *Journal of Communication Disorders*,
13 42(3), 226–234. <http://doi.org/10.1016/j.jcomdis.2009.01.001>

14 Jones, J. L., Zalewski, C., Brewer, C., Lucker, J., & Drayna, D. (2009). Widespread auditory
15 deficits in tone deafness. *Ear and Hearing*, 30(1), 63.

16 Kaganovich, N., Francis, A. L., & Melara, R. D. (2006). Electrophysiological evidence for early
17 interaction between talker and linguistic information during speech perception. *Brain*
18 *Research*, 1114(1), 161–172. <http://doi.org/http://dx.doi.org/10.1016/j.brainres.2006.07.049>

19 Lee, C.-Y., Tao, L., & Bond, Z. S. (2008). Identification of acoustically modified Mandarin
20 tones by native listeners. *Journal of Phonetics*, 36(4), 537–563.
21 <http://doi.org/10.1016/j.wocn.2008.01.002>

22 Lee, C.-Y., Tao, L., & Bond, Z. S. (2010). Identification of multi-speaker Mandarin tones in
23 noise by native and non-native listeners. *Speech Communication*, 52(11), 900–910.

- 1 Lehiste, I. (1970). *Suprasegmentals*. Cambridge, M.A.: MIT Press.
- 2 Lenth, R. V. (2016). Least-squares means: the R package lsmeans. *Journal of Statistical*
3 *Software*, 69(1), 1–33. <http://doi.org/10.18637/jss.v069.i01>
- 4 Liu, F., Chan, A. H. D., Ciocca, V., Roquet, C., Peretz, I., & Wong, P. C. M. (2016). Pitch
5 perception and production in congenital amusia: Evidence from Cantonese speakers. *The*
6 *Journal of the Acoustical Society of America*, 140(1), 563–575.
7 <http://doi.org/10.1121/1.4955182>
- 8 Liu, F., Jiang, C., Thompson, W. F., Xu, Y., Yang, Y., & Stewart, L. (2012). The mechanism of
9 speech processing in congenital amusia: Evidence from Mandarin speakers. *PLoS ONE*,
10 7(2), e30374. <http://doi.org/10.1371/journal.pone.0030374>
- 11 Liu, F., Maggu, A. R., Lau, J. C. Y., & Wong, P. C. M. (2015). Brainstem encoding of speech
12 and musical stimuli in congenital amusia: Evidence from Cantonese speakers. *Frontiers in*
13 *Human Neuroscience*, 8(1029). article. <http://doi.org/10.3389/fnhum.2014.01029>
- 14 Liu, F., Patel, A. D., Fourcin, A., & Stewart, L. (2010). Intonation processing in congenital
15 amusia: discrimination, identification and imitation. *Brain*, 133(6), 1682–1693.
16 <http://doi.org/10.1093/brain/awq089>
- 17 Lu, X., Ho, H. T., Liu, F., Wu, D., & Thompson, W. F. (2015). Intonation processing deficits of
18 emotional words among Mandarin Chinese speakers with congenital amusia: an ERP study.
19 *Frontiers in Psychology*, 6, 385. <http://doi.org/10.3389/fpsyg.2015.00385>
- 20 Macmillan, N. A., & Creelman, C. D. (2005). *Detection theory: A user's guide* (2nd ed.).
21 Mahwah, NJ, US: Lawrence Erlbaum Associates Publishers.
- 22 Magnuson, J. S., & Nusbaum, H. C. (2007). Acoustic differences, listener expectations, and the
23 perceptual accommodation of talker variability. *Journal of Experimental Psychology:*

1 *Human Perception and Performance*, 33(2), 391–409.

2 Martin, C. S., Mullennix, J. W., Pisoni, D. B., & Summers, W. V. (1989). Effects of talker
3 variability on recall of spoken word lists. *Journal of Experimental Psychology: Learning,*
4 *Memory, and Cognition*, 15(4), 676.

5 Matthews, S., & Yip, V. (2013). *Cantonese: A Comprehensive Grammar*. Routledge.

6 Mullennix, J. W., & Pisoni, D. B. (1990). Stimulus variability and processing dependencies in
7 speech perception. *Perception & Psychophysics*, 47(4), 379–390.

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

11 Nguyen, S., Tillmann, B., Gosselin, N., & Peretz, I. (2009). Tonal language processing in
12 congenital amusia. *Annals of the New York Academy of Sciences*, 1169(1), 490–493. JOUR.
13 <http://doi.org/10.1111/j.1749-6632.2009.04855.x>

14 Nusbaum, H. C., & Morin, T. M. (1992). Paying attention to differences among talkers. In Y. Y.
15 Tohkura Sagasaka, and E. Vatikiotis-Bateson (Ed.), *Speech Perception, Speech Production,*
16 *and Linguistic Structure* (pp. 113–134). Tokyo: OHM.

17 Patel, A. D., Foxton, J. M., & Griffiths, T. D. (2005). Musically tone-deaf individuals have
18 difficulty discriminating intonation contours extracted from speech. *Brain and Cognition*,
19 59(3), 310–313. JOUR. <http://doi.org/http://dx.doi.org/10.1016/j.bandc.2004.10.003>

20 Peng, G. (2006). Temporal and tonal aspects of Chinese syllables: A corpus-based comparative
21 study of Mandarin and Cantonese. *Journal of Chinese Linguistics*, 34(1), 135–154.

22 Peng, G., Zhang, C., Zheng, H.-Y., Minett, J. W., & Wang, W. S.-Y. (2012). The effect of inter-
23 talker variations on acoustic-perceptual mapping in Cantonese and Mandarin tone systems.

1 *J Speech Lang Hear Res*, 55(2), 579–595. [http://doi.org/10.1044/1092-4388\(2011/11-0025\)](http://doi.org/10.1044/1092-4388(2011/11-0025)

2 Peretz, I., Ayotte, J., Zatorre, R. J., Mehler, J., Ahad, P., Penhune, V. B., & Jutras, B. (2002).

3 Congenital amusia: A disorder of fine-grained pitch discrimination. *Neuron*, 33(2), 185–191.

4 [http://doi.org/http://dx.doi.org/10.1016/S0896-6273\(01\)00580-3](http://doi.org/http://dx.doi.org/10.1016/S0896-6273(01)00580-3)

5 Peretz, I., Champod, A. S., & Hyde, K. (2003). Varieties of musical disorders. *Annals of the New*

6 *York Academy of Sciences*, 999(1), 58–75.

7 Peretz, I., & Hyde, K. L. (2003). What is specific to music processing? Insights from congenital

8 amusia. *Trends in Cognitive Sciences*, 7(8), 362–367.

9 [http://doi.org/http://dx.doi.org/10.1016/S1364-6613\(03\)00150-5](http://doi.org/http://dx.doi.org/10.1016/S1364-6613(03)00150-5)

10 Peretz, I., & Vuvan, D. T. (2017). Prevalence of congenital amusia. *European Journal of Human*

11 *Genetics*.

12 Perrachione Pierrehumbert, Janet B., and Wong, Patrick C. M., T. K. (2009). Differential neural

13 contributions to native- and foreign-language talker identification. *Journal of Experimental*

14 *Psychology: Human Perception and Performance*, 35(6), 1950–1960.

15 Pisoni, D. B. (1993). Long-term memory in speech perception: Some new findings on talker

16 variability, speaking rate and perceptual learning. *Speech Communication*, 13(1–2), 109–

17 125. [http://doi.org/10.1016/0167-6393\(93\)90063-Q](http://doi.org/10.1016/0167-6393(93)90063-Q)

18 Pruitt, J. S., Jenkins, J. J., & Strange, W. (2006). Training the perception of Hindi dental and

19 retroflex stops by native speakers of American English and Japanese. *The Journal of the*

20 *Acoustical Society of America*, 119, 1684.

21 Shao, J., Zhang, C., Peng, G., Yang, Y., & Wang, W. S.-Y. (2016). Effect of noise on lexical

22 tone perception in Cantonese-speaking amusics. In *Proceedings of the Interspeech*. San

23 Francisco, U.S.A.

1 Shu, H., Peng, H., & McBride-Chang, C. (2008). Phonological awareness in young Chinese
2 children. *Developmental Science*, *11*(1), 171–181. <http://doi.org/10.1111/j.1467->
3 [7687.2007.00654.x](http://doi.org/10.1111/j.1467-7687.2007.00654.x)

4 Sjerps, M. J., Zhang, C., & Peng, G. (2017). Lexical tone is perceived relative to locally
5 surrounding context, vowel quality to preceding context. *Journal of Experimental*
6 *Psychology: Human Perception and Performance*, No Pagination Specified-No Pagination
7 Specified. <http://doi.org/10.1037/xhp0000504>

8 SPSS, I. B. M. (2011). IBM SPSS statistics for Windows, version 20.0. *New York: IBM Corp.*

9 Strange, W. (2011). Automatic selective perception (ASP) of first and second language speech:
10 A working model. *Journal of Phonetics*, *39*(4), 456–466.
11 <http://doi.org/10.1016/j.wocn.2010.09.001>

12 Sun, Y., Lu, X., Ho, H. T., & Thompson, W. F. (2017a). Pitch discrimination associated with
13 phonological awareness: Evidence from congenital amusia. *Scientific Reports*, *7*(March),
14 44285. <http://doi.org/10.1038/srep44285>

15 Sun, Y., Lu, X., Ho, H. T., & Thompson, W. F. (2017b). Pitch discrimination associated with
16 phonological awareness: Evidence from congenital amusia. *Scientific Reports*, *7*(March),
17 44285. <http://doi.org/10.1038/srep44285>

18 Team, R. C. (2014). R: A language and environment for statistical computing. Vienna, Austria:
19 R Foundation for Statistical Computing; 2014.

20 Thompson, W. F., Marin, M. M., & Stewart, L. (2012). Reduced sensitivity to emotional prosody
21 in congenital amusia rekindles the musical protolanguage hypothesis. *Proceedings of the*
22 *National Academy of Sciences of the United States of America*, *109*(46), 19027–19032.
23 <http://doi.org/10.1073/pnas.1210344109>

- 1 Tillmann, B., Burnham, D., Nguyen, S., Grimault, N., Gosselin, N., & Peretz, I. (2011).
2 Congenital amusia (or tone-deafness) interferes with pitch processing in tone languages.
3 *Frontiers in Psychology*, 2(JUN). <http://doi.org/10.3389/fpsyg.2011.00120>
- 4 Tillmann, B., Lévêque, Y., Fornoni, L., Albouy, P., & Caclin, A. (2016). Impaired short-term
5 memory for pitch in congenital amusia. *Brain Research*, 1640(Part B), 251–263.
6 <http://doi.org/http://dx.doi.org/10.1016/j.brainres.2015.10.035>
- 7 Tillmann, B., Schulze, K., & Foxtan, J. M. (2009). Congenital amusia: A short-term memory
8 deficit for non-verbal, but not verbal sounds. *Brain and Cognition*, 71(3), 259–264. JOUR.
9 <http://doi.org/http://dx.doi.org/10.1016/j.bandc.2009.08.003>
- 10 Wang, X., & Peng, G. (2014). Phonological processing in Mandarin speakers with congenital
11 amusia. *Journal of the Acoustical Society of America*, 136(6), 3360–3370.
- 12 Williamson, V. J., & Stewart, L. (2010). Memory for pitch in congenital amusia: Beyond a fine-
13 grained pitch discrimination problem. *Memory*, 18(6), 657–669. JOUR.
14 <http://doi.org/10.1080/09658211.2010.501339>
- 15 Wong, P. C. M., & Diehl, R. L. (2003). Perceptual normalization for inter- and intratalker
16 variation in Cantonese level tones. *Journal of Speech, Language, and Hearing Research*,
17 46(2), 413–421.
- 18 Zatorre, R. J., & Gandour, J. T. (2008). Neural specializations for speech and pitch: Moving
19 beyond the dichotomies. *Philosophical Transactions of the Royal Society B: Biological*
20 *Sciences*, 363(1493), 1087–1104.
- 21 Zeileis, A., & Hothorn, O. (2002). Diagnostic checking in regression relationships. *R News*, 2(3),
22 7–10.
- 23 Zhang, C. (2018). Online adjustment of phonetic expectation of lexical tones to accommodate

1 speaker variation: a combined behavioural and ERP study. *Language, Cognition and*
2 *Neuroscience*, 33(2), 175–195. <http://doi.org/10.1080/23273798.2017.1376752>

3 Zhang, C., & Chen, S. (2016). Toward an integrative model of talker normalization. *Journal of*
4 *Experimental Psychology: Human Perception and Performance*, 42(8), 1252–1268.

5 Zhang, C., Peng, G., Shao, J., & Wang, W. S. Y. (2017). Neural bases of congenital amusia in
6 tonal language speakers. *Neuropsychologia*, 97(July 2016), 18–28.
7 <http://doi.org/10.1016/j.neuropsychologia.2017.01.033>

8 Zhang, C., Peng, G., & Wang, W. S.-Y. (2012). Unequal effects of speech and nonspeech
9 contexts on the perceptual normalization of Cantonese level tones. *Journal of the Acoustical*
10 *Society of America*, 132(2), 1088–1099.

11 Zhang, C., Peng, G., & Wang, W. S.-Y. (2013). Achieving constancy in spoken word
12 identification: time course of talker normalization. *Brain and Language*, 126(2), 193–202.
13 <http://doi.org/10.1016/j.bandl.2013.05.010>

14 Zhang, C., Pugh, K. R., Mencl, W. E., Molfese, P. J., Frost, S. J., Magnuson, J. S., ... Wang, W.
15 S.-Y. (2016). Functionally integrated neural processing of linguistic and talker information:
16 An event-related fMRI and ERP study. *NeuroImage*, 124, 536–549.
17 <http://doi.org/http://dx.doi.org/10.1016/j.neuroimage.2015.08.064>

18 Zhang, C., Shao, J., & Huang, X. (2017). Deficits of congenital amusia beyond pitch: Evidence
19 from impaired categorical perception of vowels in Cantonese-speaking congenital amusics.
20 *PloS One*, 12(8), e0183151.

21
22
23

- 1 Supplemental file 1 (S1) shows the percentage of correct and incorrect trials in each condition.
- 2 Supplemental file 2 (S2) reports the RT results separated by correct and incorrect trials.
- 3

1 Figure 1. F0 curves of the six Cantonese tones used as stimuli in the talker variation condition.
2
3 Figure 2. F0 curves of the six Cantonese tones used as stimuli in the syllable variation condition.
4
5 Figure 3. Results in the talker variation condition for the amusic and control groups in low and
6 high variation conditions. (a) Identification accuracy. The dotted line indicates the chance level
7 accuracy ($1/6 = 0.167$; there were six tone choices in the identification task); (b) Identification
8 RT combining correct and incorrect trials; (c) Identification RT of correct trials only; (d)
9 Identification RT of incorrect trials only; (e) Discrimination sensitivity index d' ; (f)
10 Discrimination RT combining correct and incorrect trials; (g) Discrimination RT of correct trials
11 only; (h) Discrimination RT of incorrect trials only. HV = high variation; LV = low variation.
12
13 Figure 4. Results in the syllable variation condition for the amusic and control groups in low and
14 high variation conditions. (a) Identification accuracy. The dotted line indicates the chance level
15 accuracy ($1/6 = 0.167$; there were six tone choices in the identification task); (b) Identification
16 RT combining correct and incorrect trials; (c) Identification RT of correct trials only; (d)
17 Identification RT of incorrect trials only; (e) Discrimination sensitivity index d' ; (f)
18 Discrimination RT combining correct and incorrect trials; (g) Discrimination RT of correct trials
19 only; (h) Discrimination RT of incorrect trials only. HV = high variation; LV = low variation.
20

1 Table 1. Demographic characteristics of the amusic and control participants.

Subject information	Amusics	Controls
No. of participants	16 (8 M, 8 F)	16 (8 M, 8 F)
Age (range)	22.35 ± 2.8 years (19.1-27.5 years)	22.5 ± 3.1 years (18.7-28.5 years)
<i>MBEA (SD)</i>		
Scale	50.8 (17.7)	90.4 (5.6)
Contour	58.4 (19.6)	93.5 (4.9)
Interval	54.3 (18.1)	90.8 (4.3)
Rhythm	55.6 (15.0)	95.3 (3.6)
Meter	45.5 (10.4)	74.7 (14.3)
Memory	63.5 (23.3)	98.1 (2.9)
Global	54.7 (14.7)	90.5 (2.7)

2 *Note.* Amusics and controls were identified using the Montreal Battery of Evaluation of Amusia
 3 (MBEA) (Peretz et al., 2003). Amusics scored lower than 71% in the global score, which is the
 4 mean of all six subtests, whereas controls scored higher than 80%. M = male; F = female.

5

1 **S1.** Percentage of correct and incorrect trials in each condition.

Amusics	High variation		Low variation	
	Incorrect	Correct	Incorrect	Correct
Identification (talker)	59%	41%	59%	41%
Discrimination (talker)	27%	73%	20%	80%
Identification (syllable)	56%	44%	56%	44%
Discrimination (syllable)	46%	54%	24%	76%

Controls	High variation		Low variation	
	Incorrect	Correct	Incorrect	Correct
Identification (talker)	40%	60%	37%	63%
Discrimination (talker)	15%	85%	6%	94%
Identification (syllable)	40%	60%	35%	65%
Discrimination (syllable)	29%	71%	7%	93%

2

3

1 **S2.** RT results separated by correct and incorrect trials.

2 **Talker variation condition**

3 *Identification task*

4 For the identification task, no effects were significant for either correct or incorrect trials.

5

6 *Discrimination task*

7 Correct trials

8 There was only a significant main effect of *condition* ($\chi^2(1) = 8.702, p = 0.003$), where RT in the
9 high variation condition was longer than that in the low variation condition.

10

11 Incorrect trials

12 No effects were statistically significant.

13

14 **Syllable variation condition**

15 *Identification task*

16 Correct trials

17 There were a main effect of *condition* ($\chi^2(1) = 14.437, p < 0.001$), and significant two-way
18 interaction between *group* and *condition* ($\chi^2(1) = 8.769, p = 0.003$). Pairwise comparisons
19 revealed that the effect of group was not significant in either condition. The RT in the high
20 variation condition was longer than in the low variation condition in the control group ($z = 5.952,$
21 $p < 0.001$), but not in the amusic group ($z = 1.238, p = 0.224$)

22

23 Incorrect trials

1 There was a significant main effect of *group* ($\chi^2(1) = 5.371, p < 0.001$), where the RT of the
2 control group was longer than the amusic group. There was also a main effect of *condition* ($\chi^2(1)$
3 = 19.571, $p < 0.001$), where the RT in the high variation condition was longer than that in the
4 low variation condition.

5

6 *Discrimination task*

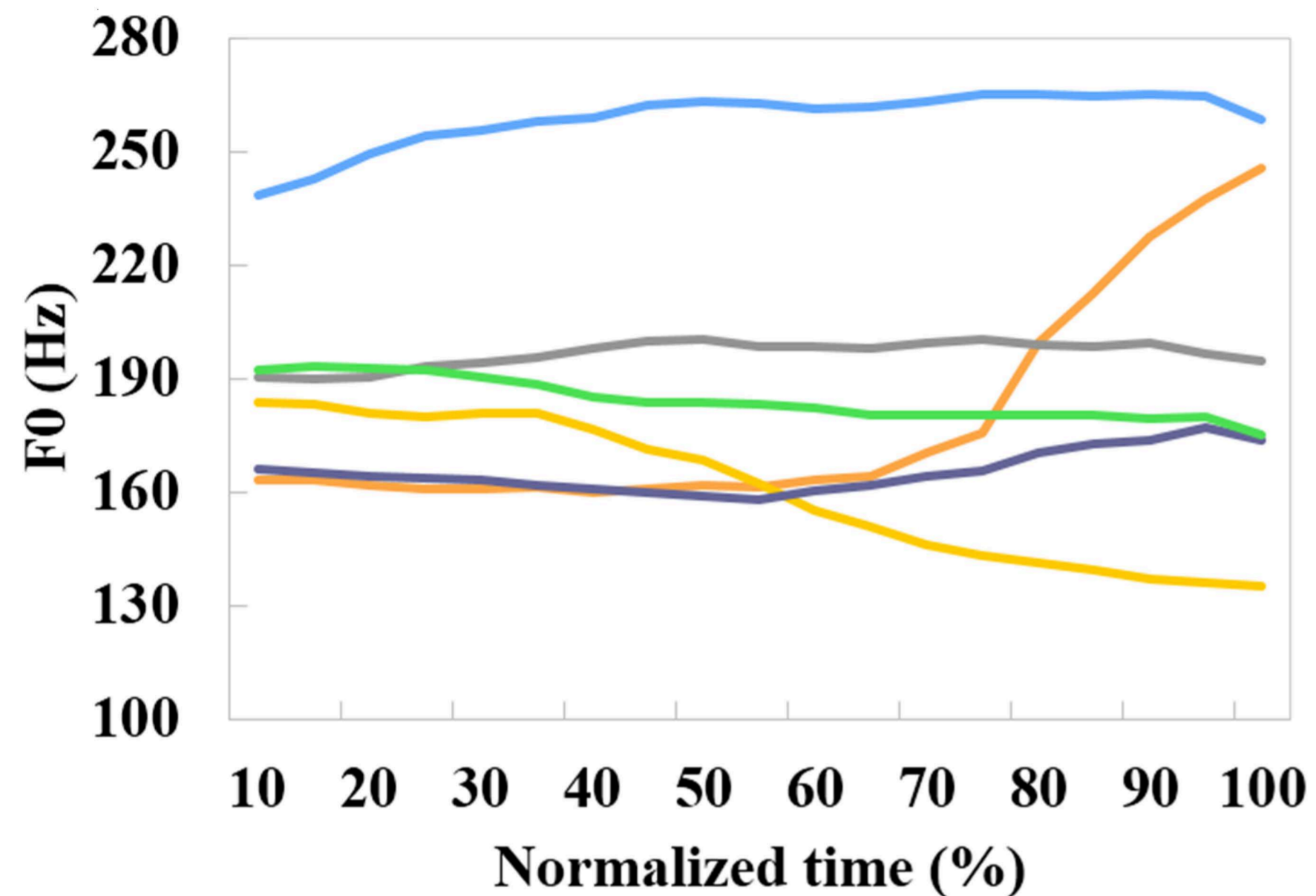
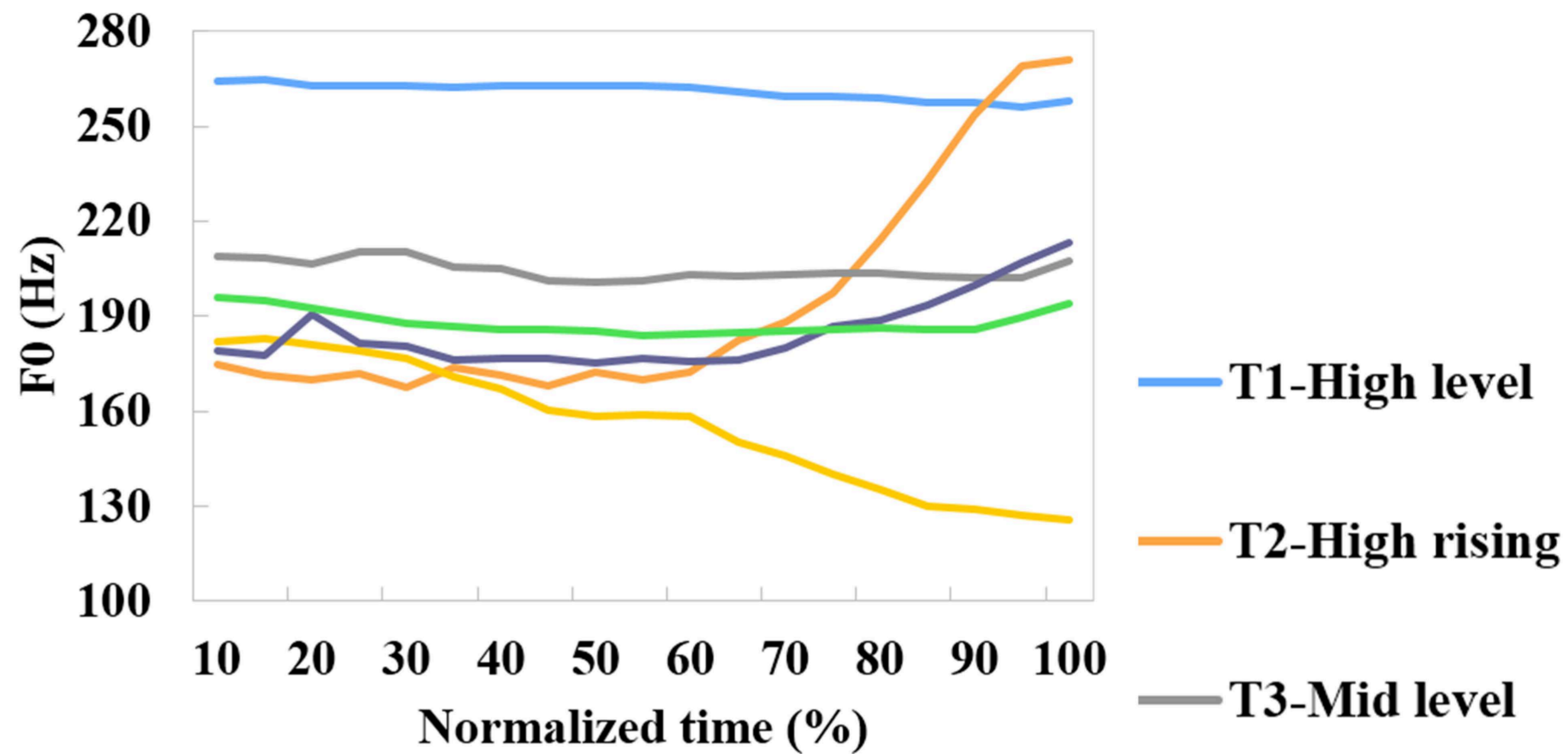
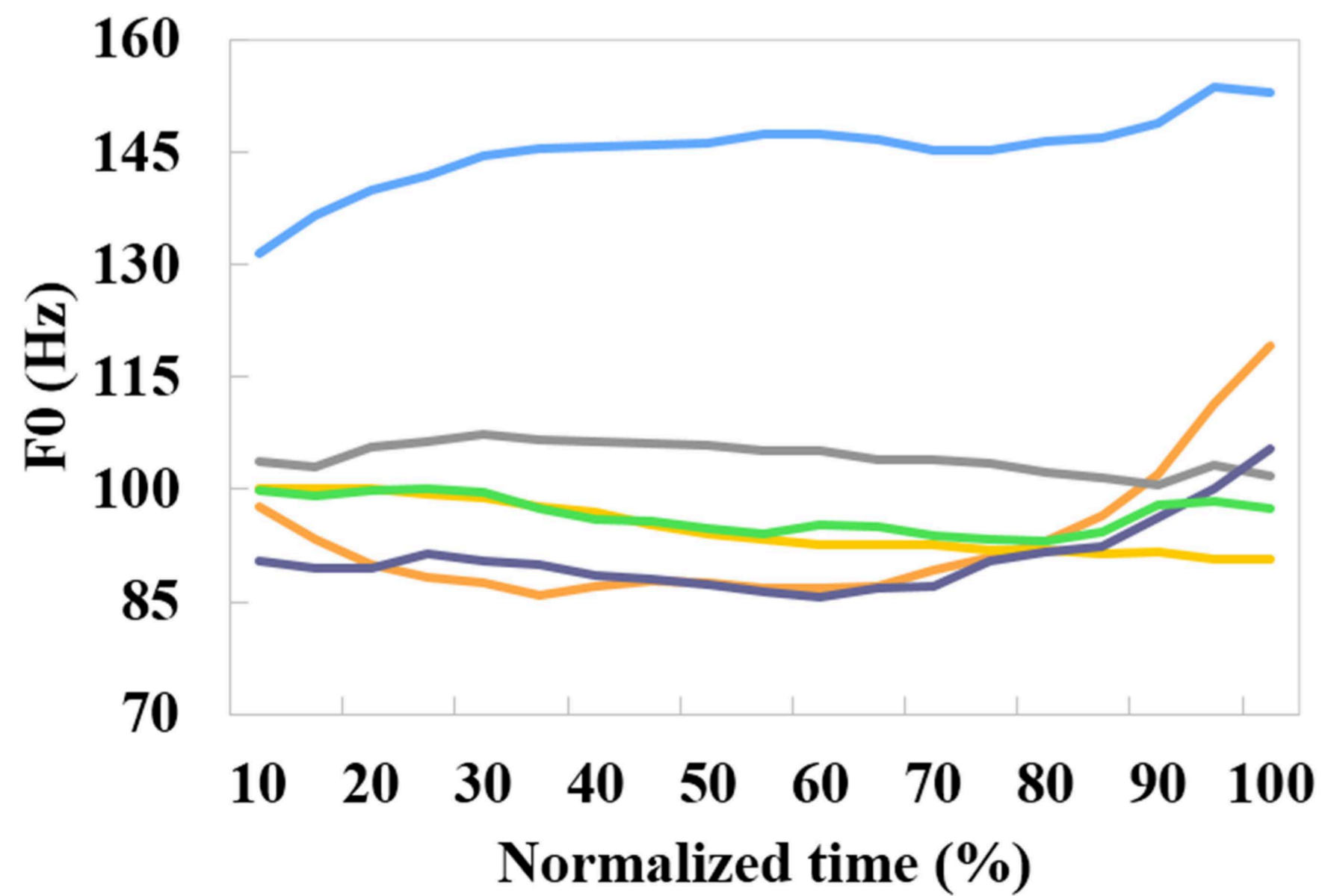
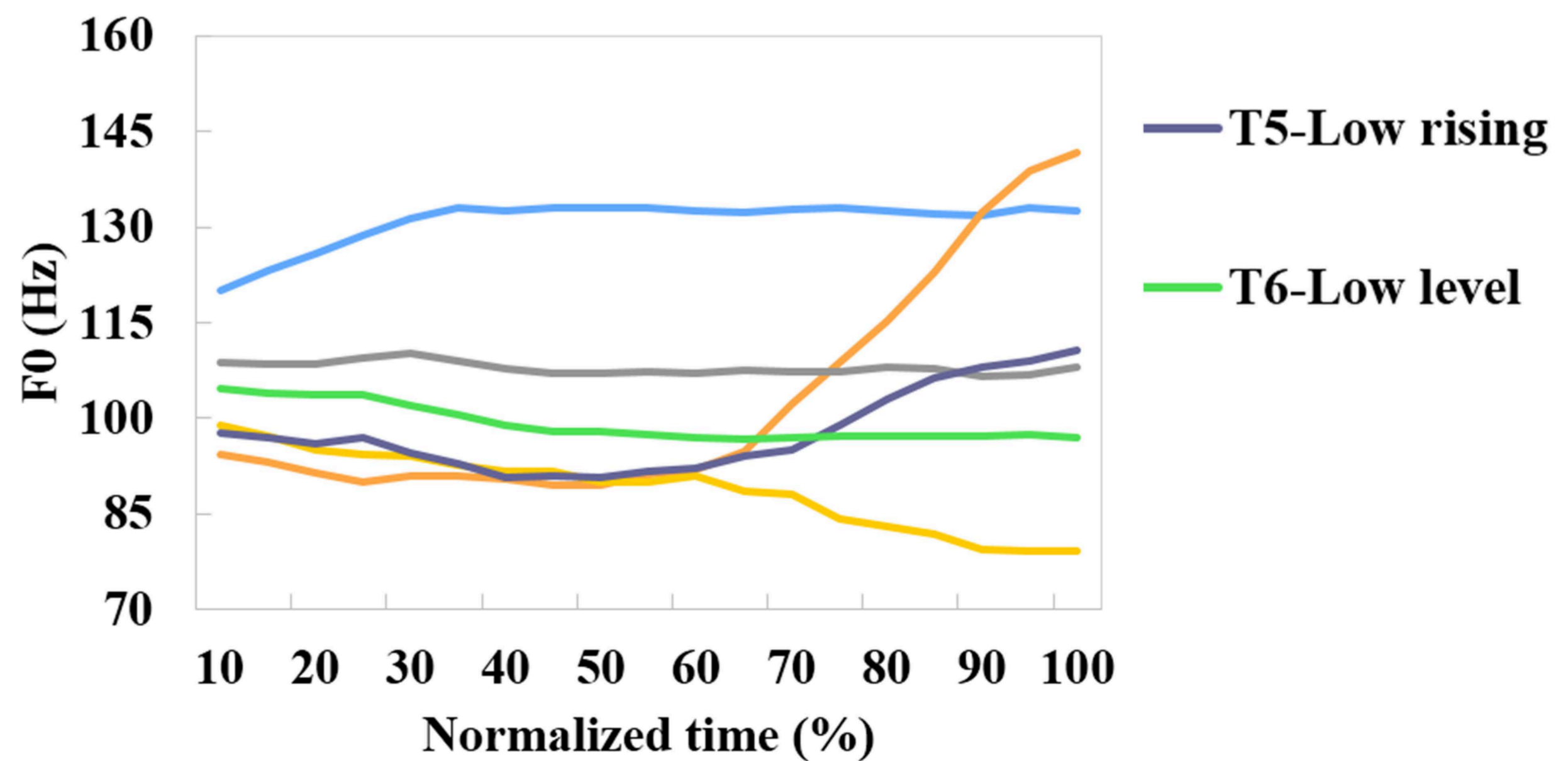
7 In the discrimination task, the RT of both correct and incorrect trials showed significant main
8 effects of *condition* (correct trials: ($\chi^2(1) = 13.614, p < 0.001$); incorrect trials ($\chi^2(1) = 7.395, p <$
9 0.001), where the RT in the high variation condition was longer than that in the low variation
10 condition. No other effects were significant.

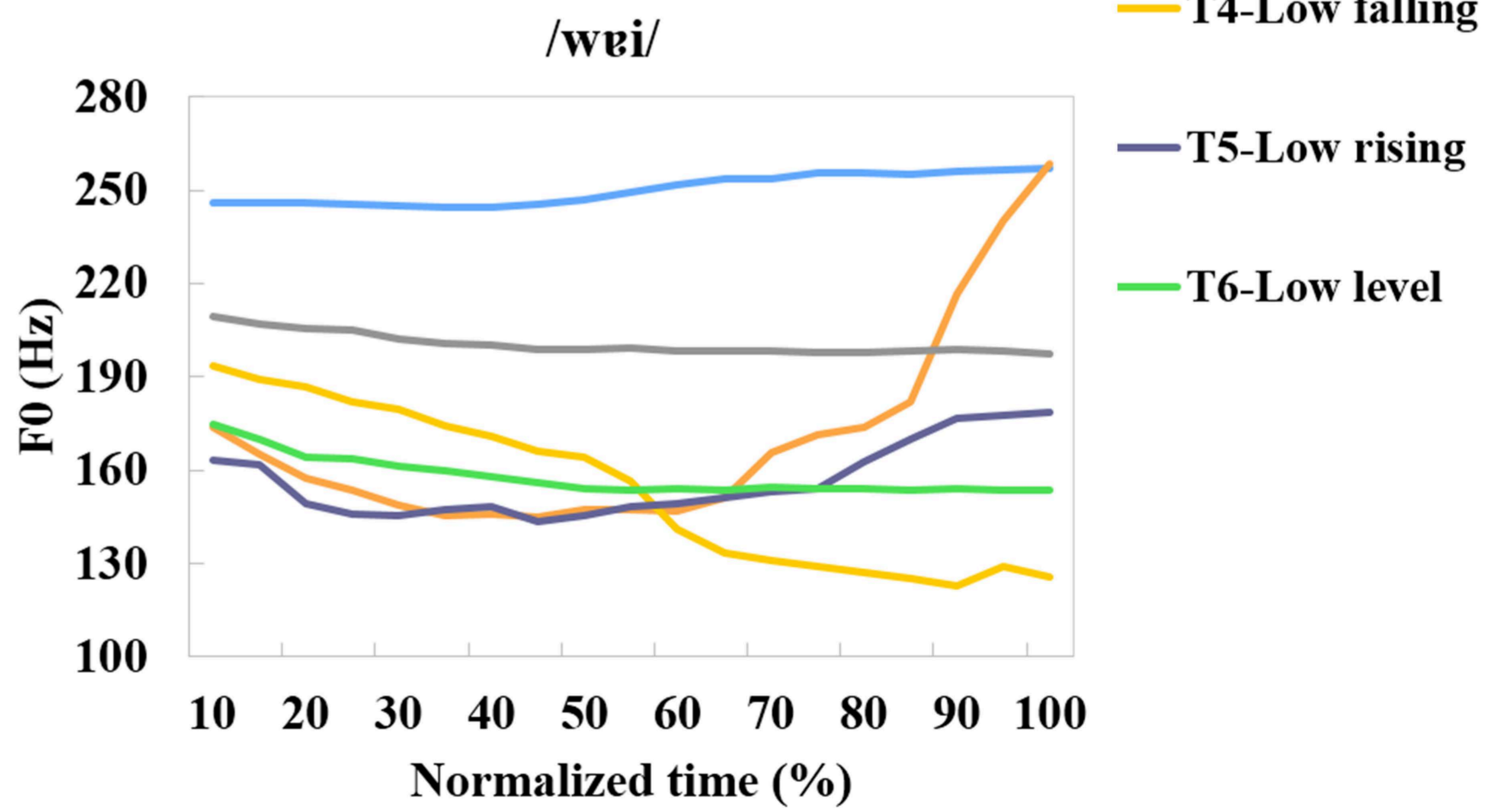
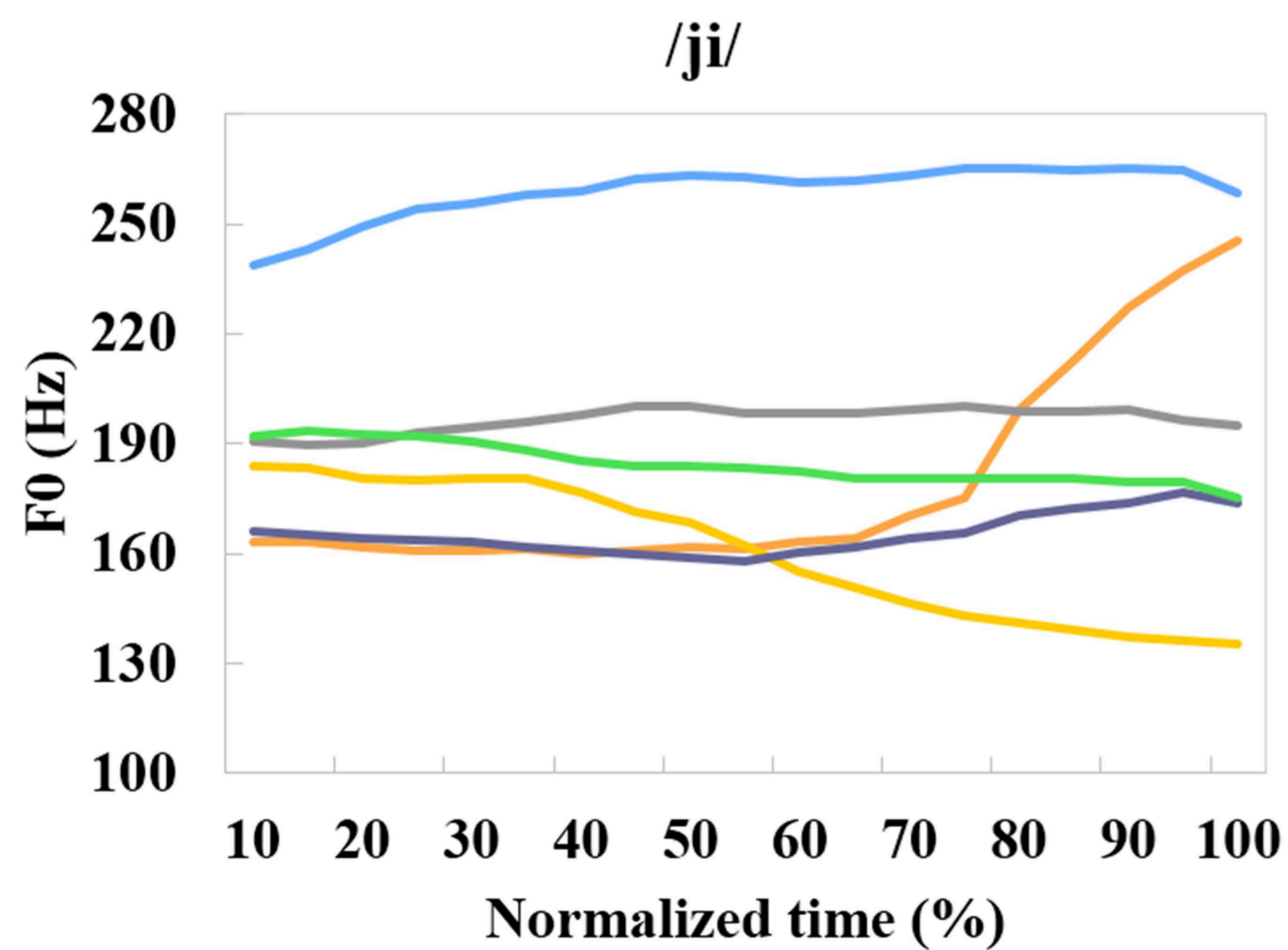
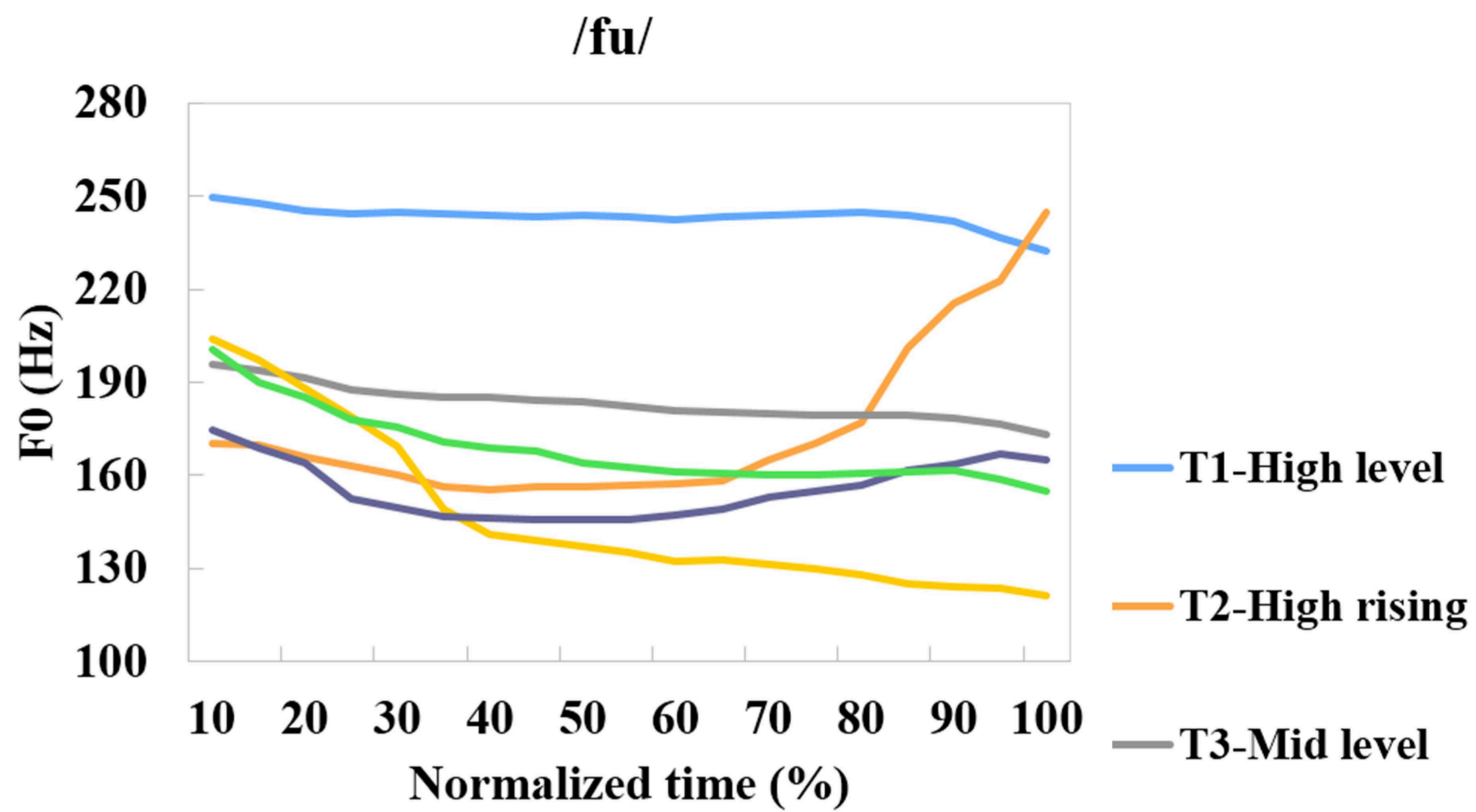
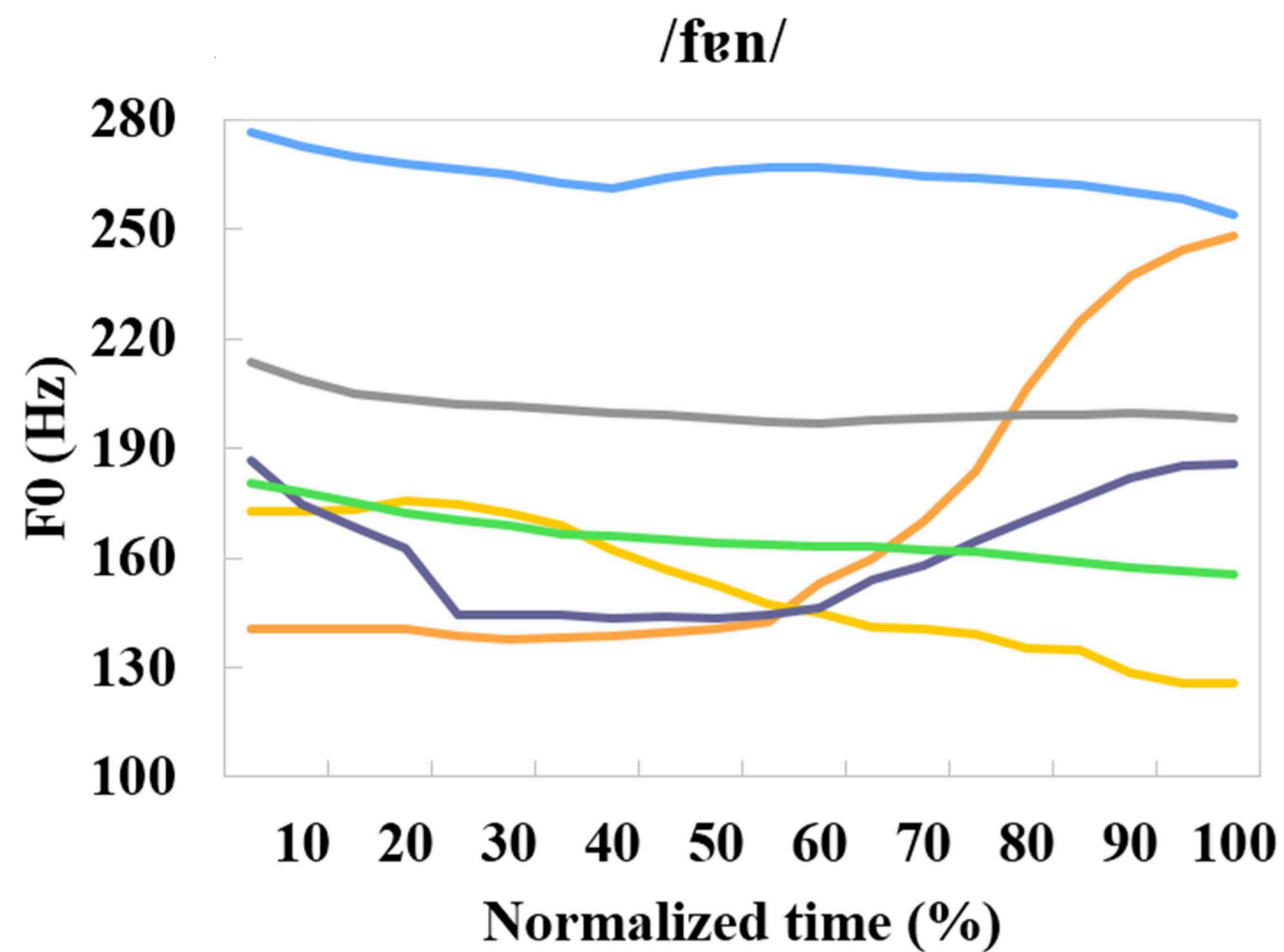
11

1 **Appendix.** The four sets of syllables used in the syllable variation condition.

Tone	/ji/	/fen/	/fu/	/wei/
T1 high level /55/	醫 “doctor”	婚 “marriage”	夫 “husband”	威 “power”
T2 high rising /25/	椅 “chair”	粉 “pink”	父 “father”	委 “council”
T3 mid level /33/	意 “meaning”	訓 “train”	褲 “trousers”	餵 “feed”
T4 low falling /21/	兒 “son”	焚 “burn”	符 “symbol”	圍 “surround”
T5 low rising /23/	耳 “ear”	奮 “strive”	婦 “woman”	偉 “grand”
T6 low level /22/	二 “two”	份 “part”	負 “negative”	彗 “comet”

2

Talker_F01**Talker_F02****Talker_M01****Talker_M02**



- T1-High level
- T2-High rising
- T3-Mid level
- T4-Low falling
- T5-Low rising
- T6-Low level

