

1 **Achieving perceptual constancy with context cues in second language speech perception**

2

3

Kaile Zhang^a, Defeng Li^b, Gang Peng^{c*}

4

5 * Corresponding author: gpeng@polyu.edu.hk

6 ^a Centre for Cognitive and Brain Sciences, University of Macau, Taipa, Macau Special
7 Administrative Region, China

8 ^b Centre for Studies of Translation, Interpreting and Cognition, Faculty of Arts and Humanities,
9 University of Macau, Taipa, Macau Special Administrative Region, China

10 ^c Research Centre for Language, Cognition, and Neuroscience, Department of Chinese and
11 Bilingual Studies, The Hong Kong Polytechnic University, Hung Hom, Kowloon, Hong Kong
12 Special Administrative Region, China

13

14 **Abstract**

15 Context cues are useful for listeners to normalize speech variability and achieve
16 perceptual constancy. It remains unknown whether this normalization strategy is language-
17 independent and can be generalized directly from the perception of first language (L1) to
18 second language (L2). To answer this question, Experiment 1 in the present study asked
19 Mandarin learners of Cantonese to perceive ambiguous Cantonese tones with context cues. The
20 results revealed a significant Cantonese-tone normalization process in Mandarin learners, but
21 the effect size was smaller than native speakers, suggesting that speech normalization required
22 language-specific knowledge and thus it was refined gradually during L2 acquisition. The
23 results also showed that even with effective context cues, Mandarin learners tended to give
24 more high level tone responses, a tone also in Mandarin, implying that L1 phonological system
25 interacts with immediate L2 context during L2 speech normalization. Experiment 2 revealed
26 that L2 immersion but not overall L2 proficiency or L2 phonological proficiency facilitated L2
27 normalization process, indicating that L2 speech normalization improved with perceptual
28 practice and needed more high-level L2 knowledge than L2 phonology.

29

30 **Keywords**

31 Speech variability, perceptual normalization, context cues, L2 acquisition, lexical tones,
32 language immersion

33

Achieving perceptual constancy with context cues in second language speech perception

1.0 Introduction.

Speech signals produced by different talkers vary a lot due to individual anatomic configurations of vocal tracts (Peterson & Barney, 1952). Even the same speaker's speech changes noticeably in different psychological conditions. The lack of one-to-one mapping between highly variable speech signals and the abstract linguistic representations is a fundamental difficulty for human speech perception (K. Johnson, 2005; Liberman et al., 1967; Sjerps et al., 2019). Listeners took more time to identify words presented in mixed-talker conditions than in blocked-talker conditions since the mixed-talker condition required frequent adaptation to new talkers (Nusbaum and Morin, 1992). Sometimes, talker variability even blurs the boundary of two similar phonemes. For example, fundamental frequency (f_0) is the main acoustic cue to identify lexical tones, but the f_0 of a female speaker's low tone could be even higher than the f_0 of a male speaker's high tone, resulting in more mistakes for tone perception in multiple-talker condition. Talker variability impacts not only native speakers but also second language (L2) learners. While there has been considerable research into how native speakers manage talker-related speech variability, few studies have examined the perceptual strategies L2 learners employ to cope with this variability—this gap will be addressed by the current study. The investigation on this question would enhance the understanding of the challenges inherent in L2 speech perception. The following part of the introduction would first review relevant studies in L1 and L2 speech normalization, and then present the research plan for the present study.

1.1 Normalizing talker variability with context cues

Although talker variability presents challenges for fast and accurate speech perception, in most cases, native listeners are able to communicate successfully with different talkers. One way to achieve perceptual constancy is to remove talker-related differences from each speech token and obtain a talker-independent linguistic pattern for phonemic identification (i.e., speech perception with normalization process) (K. Johnson, 2005). This perceptual strategy was supposed by recent ECoG studies which found that the talker-independent speech cues were represented at the auditory cortex (Ogania et al., 2023; Sjerps et al., 2019). Research has

1 shown that talker variability can be reduced through the use of contextual cues—a mechanism
2 referred to as extrinsic normalization (Ainsworth, 1975; Nearey, 1989). For instance, a level
3 tone from a multiple-speaker corpus can hardly be identified as either high tone or low tone,
4 but the identification rate improved dramatically once it was presented in a speech context
5 (Peng et al., 2012; Wong & Diehl, 2003). Contextual cues generally influence the perception
6 of the target speech in a contrastive manner. Specifically, an ambiguous tone is more likely to
7 be perceived as low when preceded by a high-pitched context and as high when preceded by a
8 low-pitched context, an effect known as the contrastive context effect (Moore & Jongman,
9 1997; K. Zhang et al., 2021).

10 Although extrinsic normalization has been widely demonstrated to be effective in
11 overcoming talker variability, it is inconclusive about the underlying cognitive mechanism.
12 Understanding the cognitive mechanism of extrinsic normalization is crucial for the present
13 study, as various theoretical models yield contrasting predictions regarding L2 learners’
14 normalization performance (refer to Section 1.2 for details). Prior to delving into L2 speech
15 normalization, the present study would first explore the different theoretical frameworks about
16 extrinsic normalization. The existing theories regarding extrinsic normalization fall into two
17 broad categories, distinguished by the level at which the process occurs: either at a general-
18 auditory level or at a speech-specific level.

19 Some researchers have posited that the extrinsic normalization relies on the acoustic
20 properties of speech signals and operates at the general-auditory level (e.g., Holt & Lotto, 2002).
21 In the extrinsic normalization process, the auditory system first assesses the acoustic properties
22 of the preceding context (e.g., long-term average spectrum, f_0 , and durations) and then
23 automatically encodes the target speech cue in contrast to the acoustic properties of the
24 preceding context (Watkins & Makin, 1994, 1996). The neural adaptation within the central
25 auditory system might be the biological basis for the acoustically contrastive encoding account
26 of extrinsic normalization. Neurons in the primary auditory cortex, attuned to different
27 frequency regions, become less responsive to frequencies similar to those in the preceding
28 context due to adaptation. Conversely, neurons that were less active during the context phase
29 become more sensitive to the frequencies of subsequent sounds, leading to the spectral contrast
30 perception (Stilp, 2019). Supporting evidence for this general-auditory level process comes
31 from studies like Holt et al. (2001) which demonstrated that even birds trained to recognize
32 speech sounds exhibit contrastive perceptual behavior.

1 Other scholars are of the view that extrinsic normalization necessitates language-
2 specific knowledge and functions at the level of speech-specific language processing. Nusbaum
3 and Morin, (1992); Nusbaum & Magnuson, (1997) proposed that listeners could use context
4 information to compute a talker-specific mapping between acoustic patterns and abstract
5 phonemic representations, and ambiguous target cues were identified by referring to the talker-
6 specific acoustic-phonemic mapping. For example, in the Cantonese greeting “早晨” (/zou 25
7 san 21/, good morning), the highest F0 aligns with the highest pitch (the end point of T25) and
8 lowest F0 matches the lowest pitch in Cantonese (the end point of T21), which outline a
9 speaker’s acoustic space for the phonemic system (the tonal system in this example) (K. Zhang
10 et al., 2023). The subsequent tones were perceived by referring to the talker-specific acoustic-
11 phonemic space. Based on the acoustic-phonemic mapping account of the extrinsic
12 normalization, language-specific knowledge is required since listeners at least need to know
13 phonological categories to construct the talker-specific acoustic-phonemic mapping. Beyond
14 phonological information, other higher-level linguistic knowledge, such as semantic or
15 syntactic information, is also helpful for the normalization process, as C. Zhang et al. (2015)
16 demonstrated that contexts composed of meaningful speech elicited better normalization
17 effects than meaningless word sequences.

18 *1.2. Speech normalization in L2 speech perception*

19 Talker variability also posed difficulties for L2 speech perception. Perceiving high-
20 variability L2 speech could be a dual challenge for L2 learners, as they have to cope with talker
21 variability and imperfect L2 knowledge (Lecumberri et al., 2010). Based on the cognitive
22 mechanisms of extrinsic normalization outlined previously, opposite predictions would be
23 drawn concerning L2 learners’ handling of L2 speech variability. Should extrinsic
24 normalization be a general-auditory level process, it would likely be language-independent,
25 allowing the normalization skills developed in L1 to be directly applied to L2. This would
26 result in comparable performances in L2 learners and native speakers in terms of managing L2
27 speech variability. Conversely, if extrinsic normalization operates at a speech-specific level
28 that relies on language-specific information, it would be language-dependent. In this case, L2
29 learners’ proficiency in extrinsic normalization is expected to incrementally reach a native-like
30 standard as their L2 proficiency increases.

1 Albeit the well-known challenge of speech variability in L2 speech perception, only a
2 limited amount of research investigated how L2 learners deal with speech variability, yielding
3 different results. The central debate surrounding these studies concerns whether the
4 normalization process is language-independent and can be applied directly to L2 based on prior
5 experience with L1 speech, echoing the dispute between two theoretical frameworks reviewed
6 above. Bradlow & Pisoni, (1999) found that talker variability had a comparable effect on
7 English word identification in native and non-native English speakers (7.2% vs. 7.9%) which
8 was quantified by calculating the difference in listeners' word identification rates between the
9 multiple-talker and single-talker conditions. Similar results were observed in studies examining
10 the perception of Mandarin tones (C. Y. Lee et al., 2013) and Mandarin fricatives (C.-Y. Lee
11 et al., 2012), where native and non-native speakers showed equivalent performance decrements
12 in mixed-talker condition (i.e., high talker variability) relative to the blocked-talker condition
13 (i.e., low talker variability). These findings led Bradlow and Pisoni (1999) and C. Y. Lee et al.
14 (2013) to conclude that the ability to deal with talker variability was a language-independent
15 skill that transferred easily from L1 to L2. However, Antoniou et al. (2015) reported a more
16 severe impact of speech variability on L2 learners than on native speakers in an English word-
17 monitoring task, in terms of signal detection accuracy and recognition time, suggesting a poorer
18 speech normalization ability of L2 learners. Additionally, Bent et al., (2010) found that Korean
19 listeners identifying English vowels from multiple speakers achieved a 22% lower accuracy
20 rate than native English speakers, with the error pattern analysis indicating a strong influence
21 of production variability within the normal range for native English talkers. This suggests that
22 non-native listeners may be more susceptible to a cross-talker variability due to difficulties in
23 forming native-like vowel categories. The limited cognitive resources available for L2 learners,
24 particularly those with low L2 proficiency, also matter since speech normalization is
25 constrained by cognitive faculties like attention and working memory (Kapadia & Perrachione,
26 2020; Nusbaum & Magnuson, 1997; Wong et al., 2004).

27 Bradlow and Pisoni (1999), C. Y. Lee et al. (2013), and Antoniou et al. (2015)
28 investigated L2 speech normalization using a blocked- vs. mixed- talker (or the single- vs.
29 multiple- talker) paradigm. Talker identity shifts more frequently in the mixed-talker condition,
30 requiring increased talker normalization. However, the stimuli in these studies were words
31 presented in isolation but not in contexts, which meant that no explicit external information
32 was available for speech normalization. As a result, the contrasts between blocked and mixed-
33 talker conditions in these studies demonstrate the extent to which talker variability impacts L2

1 speech perception, rather than how effectively L2 learners utilize contextual cues to mitigate
2 this variability. To better quantify the extrinsic normalization strategy employed by L2 learners,
3 language tasks such as the identification of multiple-talker speech within contexts could be
4 more appropriate, which would be adopted by the present study. So far, no study has directly
5 tested how well L2 learners could use context cues to overcome talker variability. The most
6 related study was done by C. Y. Lee et al. (2009) who tested if contexts could help native
7 Mandarin speakers and English learners of Mandarin to identify incomplete Mandarin tones,
8 such as missing the central part. The authors found that contexts were helpful for both groups,
9 but native speakers demonstrated a greater improvement in the presence of context cues
10 compared to L2 learners. Although C. Y. Lee et al. (2009) did not explicitly address the
11 question about speech normalization, it is reasonable to infer that L2 learners also cannot use
12 context cues as effectively as native speakers to normalize talker variability.

13 *1.3 The present study*

14 Extrinsic normalization has been extensively investigated in the context of native
15 speech perception. However, it remains unclear whether L2 learners can utilize contextual cues
16 to accommodate talker variability with the same effectiveness as native speakers. This inquiry
17 is particularly significant as it sheds light on whether extrinsic normalization, as a crucial
18 strategy for achieving perceptual constancy, is a general auditory-level process or a language-
19 dependent process that is gradually acquired during L2 learning, similar to L2 learners' cue
20 weight strategy in speech perception. Further investigation into this matter could also inform
21 L2 education by highlighting the importance of training this ability in curriculum design.

22 The present study aimed to probe this question by investigating Mandarin learners'
23 extrinsic normalization of Cantonese tones. This approach differed from C. Y. Lee et al. (2009)
24 and C. Y. Lee et al. (2013) who tested the English learners' perception of Mandarin tones in
25 several aspects. Firstly, the present study utilized Cantonese tones as testing material rather
26 than Mandarin tones. Mandarin has four lexical tones: high level (T55), high rising (T35),
27 falling rising (T214), and high falling (T51), each characterized by unique pitch heights and
28 contours, thus relying primarily on intrinsic cues for differentiation. When these intrinsic cues
29 are available, the influence of context is relatively minimal. In contrast, the discrimination of
30 Cantonese tones, particularly the three level tones, depends heavily on extrinsic contexts.
31 Cantonese has three level tones: high level (T55), mid level (T33), and low level (T22). A base
32 syllable with different level tones has different meanings. For example, the base syllable /ji/

1 means ‘doctor’ with T55, ‘meaning’ with T33, and ‘two’ with T22. Three level tones share a
2 similar pitch contour (i.e., level) but different pitch heights, and thus, they can only be
3 distinguished by pitch heights. However, individual talkers' vocal fold configurations can cause
4 them to produce identical tones at different pitch heights, compromising the reliability of this
5 intrinsic cue. Therefore, context cues are indispensable for interpreting Cantonese level tones,
6 making them more suitable than Mandarin tones for testing the extrinsic normalization process
7 (e.g., Wong & Diehl, 2003; C. Zhang et al., 2012; K. Zhang et al., 2017).

8 Second, Mandarin speakers instead of English speakers (tonal vs. nontonal language
9 speakers) were invited. C. Y. Lee et al. (2013) found that both English speakers and native
10 Mandarin speakers were equally affected by talker variability while perceiving Mandarin tones.
11 Consequently, they concluded that the ability to deal with talker variability could be transferred
12 easily from L1 to L2. However, since English speakers have no experience with lexical tone
13 normalization, it is not rigorous to conclude that the normalization strategies can be directly
14 transferred from L1 to L2 based on English speakers’ normalization of Mandarin tones.
15 Mandarin learners, as speakers of a tonal language, are proficient in the tone normalization
16 process, at least in Mandarin. Recruiting Mandarin learners of Cantonese instead of non-tonal
17 language speakers made it more feasible to test whether the normalization process for lexical
18 tones acquired in L1 can be directly generalized to L2 tone perception. Similar to Cantonese
19 speakers, Mandarin learners exhibit considerable confusion with Cantonese tones that have
20 similar contours but different pitch heights, such as the three level tones (T55, T33, and T22)
21 and two rising tones (T25 and T23) (K. Zhang et al., 2018). Owing to the influence of their
22 native language (i.e., Mandarin T55), Mandarin learners often misperceive Cantonese T33 and
23 T22 as T55 when the stimuli are presented in isolation. Therefore, extrinsic contexts may be
24 even more crucial for Mandarin learners to disambiguate the similar Cantonese tones.

25 Two experiments were carried out to test the L2 speech normalization process.
26 Experiment 1 aimed to test whether extrinsic normalization process is language-independent
27 and can be generalized easily to L2 based on L1 language use experience. Experiment 2
28 followed up on the results of Experiment 1 and aimed to identify the factors that contribute to
29 L2 normalization. L2 learners’ speech normalization performance was assessed in Experiment
30 1 by testing their ability to use context cues to accommodate multiple-talker speech, as opposed
31 to the blocked- vs. mixed- talker contrast utilized by Bradlow and Pisoni (1999), C. Y. Lee et
32 al. (2013), and Antoniou et al. (2015). Mandarin learners of Cantonese were asked to perceive

1 Cantonese level tones produced by multiple talkers in speech contexts with different pitch
2 heights. The preceding context typically influences speech perception in a contrastive manner.
3 For instance, a target tone is perceived as low if it follows a context spoken by a high-pitched
4 speaker, whereas the same target tone is perceived as high following a context from a low-
5 pitched speaker. Therefore, if Mandarin learners interpreted target tones by referring to context
6 cues, they should give more T55 responses (i.e., high tone) in low-f₀ contexts, more T33
7 responses in mid-f₀ contexts, and more T22 (i.e., low tone) responses in high-f₀ contexts. An
8 experimental condition with minimal context influence was included as a baseline. To keep
9 experimental conditions more comparable, the present study used a nonspeech-context
10 condition rather than the isolated condition to evaluate Cantonese tone perception without
11 effective context information, since Cantonese speakers' tone perception in nonspeech contexts
12 was similar as in isolation (Francis et al., 2006; C. Zhang et al., 2013; K. Zhang & Peng, 2021).
13 Native Cantonese speakers were recruited as the comparison to see if L2 learners could use the
14 extrinsic normalization strategy as effectively as native speakers did. Based on previous
15 findings (C. Y. Lee et al., 2009), the present study hypothesized that Mandarin learners of
16 Cantonese would exhibit a smaller normalization effect compared to native Cantonese speakers.
17 Most Mandarin participants from Experiment 1 also participated in Experiment 2 which
18 measured several potential factors related to L2 normalization. The combination of two
19 experiments aimed to shed light on the generalizability and factors influencing the extrinsic
20 normalization process in L2 speech perception.

21 **EXPERIMENT 1: MANDARIN LEARNERS' EXTRINSIC NORMALIZATION OF** 22 **CANTONESE TONES**

23 **2.0 Materials and Methods**

24 *2.1 Participants*

25 To ensure that the Mandarin participants in the study had adequate knowledge of
26 Cantonese and could accurately complete the tasks, a short screening test was conducted. The
27 screening test comprised two parts. In the first part, a native Cantonese speaker who specialized
28 in linguistics read six sentences, three of which were used as contexts in the present study. The
29 Mandarin participants were asked to translate these sentences into Mandarin. In the second part
30 of the screening test, the Mandarin participants were asked to produce six Cantonese words
31 containing three target syllables (/ji22/, /ji33/, /ji55/), and their pronunciation was evaluated by

1 the same native Cantonese speaker. Only Mandarin participants who correctly translated the
2 Cantonese sentences and correctly produced the Cantonese words were invited to participate
3 in the Cantonese-tone normalization task.

4 Ultimately, 30 young adult native Mandarin speakers (female = 16, male = 14) who
5 were learning Cantonese participated in the Cantonese-tone normalization task. All Mandarin
6 participants were born in Mainland China. While some spoke Chinese dialects in addition to
7 Mandarin, they all identified Mandarin as their native language. Of the 30 participants, five
8 were not born in Northern China (Guizhou, Jiangxi, Hubei, Chengdu, and Shenzhen,
9 respectively). However, none of their local Chinese dialects contains more than one level tone,
10 and the subject born in Shenzhen speaks Mandarin instead of Cantonese. Thus, they were
11 included as well. Additionally, 30 young adult native Cantonese speakers (female = 15, male
12 = 15) born in Hong Kong were invited as the comparison group. None of the participants had
13 received more than three years of professional music training. All participants reported no
14 hearing, speech, or language difficulties. All participants were right-handed which was
15 assessed using the Edinburgh Handedness Inventory (Oldfield, 1971). The study was approved
16 by the Human Subjects Ethics Sub-committee of The Hong Kong Polytechnic University.
17 Comprehensive insights into the experimental procedures were provided to all participants, and
18 informed consent was acquired prior to the commencement of the experiment.

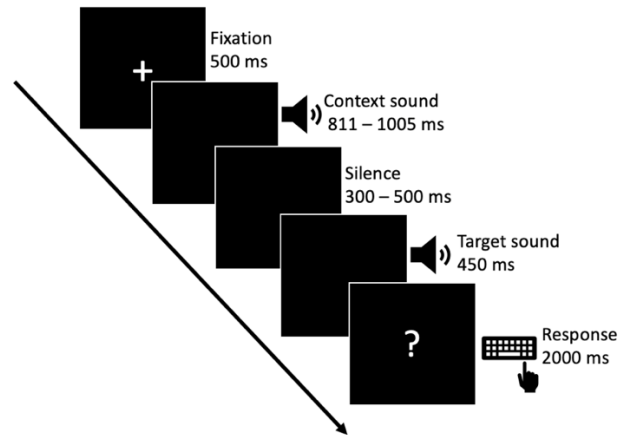
19 *2.2 Stimuli*

20 The stimuli used in the Cantonese tone normalization task were the same as those in
21 Tao et al. (2021) and K. Zhang et al. (2017). Each trial consists of either a speech context or a
22 nonspeech context and a speech target. The speech context is a Cantonese phrase 呢個字係
23 (/li55 ko33 tsi22 hei22/, this word is ...), which was produced by four native Cantonese
24 speakers of different pitch ranges (a female of high pitch, a female of low pitch, a male of high
25 pitch, and a male of low pitch). The f₀ trajectories of the original recordings were further raised
26 or lowered three semitones in Praat (Boersma & Weenink, 2023) to introduce the intra-talker
27 variability and also to trigger the contrastive context effect, resulting in three contexts of
28 different pitch heights: high-f₀ speech context, mid-f₀ speech context, and low-f₀ speech
29 context. The nonspeech contexts were composed of triangle waves. The pitch heights of
30 nonspeech contexts were manipulated to match the pitch heights of their speech counterparts,
31 resulting in three types of nonspeech contexts: high-f₀ nonspeech context, mid-f₀ nonspeech

1 context, and low-f₀ nonspeech context. The speech targets were the Cantonese syllable 意
2 (/ji33/, meaning) which were produced by the same four speakers. The task also included some
3 fillers produced by the female speaker of low pitch and the male speaker of high pitch, which
4 underwent the same pitch manipulation as the test stimuli. The context fillers were Cantonese
5 phrases: 我而家讀 (/ŋo23 ji21 ka55 tuk2/, now I will read...) and 請留心聽 (/ts^hiŋ25 ləu21
6 səm55 t^hiŋ55/, please listen to...carefully). The target fillers were Cantonese 意 (/ji33/,
7 meaning) or 二 (/ji22/, two). The speaker of the context and target stimuli was matched in each
8 trial. Detailed f₀ information of the auditory stimuli used in the present study can be found in
9 Table 1 in Tao et al. (2021). The durations of context stimuli were kept unchanged, while the
10 duration of the target stimuli was normalized to 450 ms. The intensity of speech stimuli was
11 normalized to 55dB, but the intensity of nonspeech stimuli was normalized to 75 dB to match
12 the perceived loudness of speech stimuli, which was evaluated by the native Cantonese
13 speakers.

14 *2.3 Experimental procedure*

15 The Cantonese tone normalization task which contained two blocks was conducted in
16 a sound-proofed booth: a speech-context block and a nonspeech-context block. Each block was
17 composed of 108 test trials (4 speakers × 3 pitch heights × 9 repetitions) and 36 fillers. The
18 trial procedure was illustrated in Figure 1. The audio stimuli were delivered bilaterally to
19 subjects via inserted earphones. In each trial, a 500-ms fixation was played first, which was
20 followed by a context stimulus. A period of silence jittered between 300 ms and 500 ms
21 appeared between context and target. After the target sound, a question mark appears on the
22 screen for 2000 ms. Subjects were asked to press the corresponding keys on the keyboard to
23 indicate whether the last word they heard was 醫 (/ji55/, doctor), 意 (/ji33/ meaning), or 二
24 (/ji22/, two). The next trial appears after the elimination of the question mark.



1

2

Figure 1. The trial procedure of the Cantonese-tone normalization task

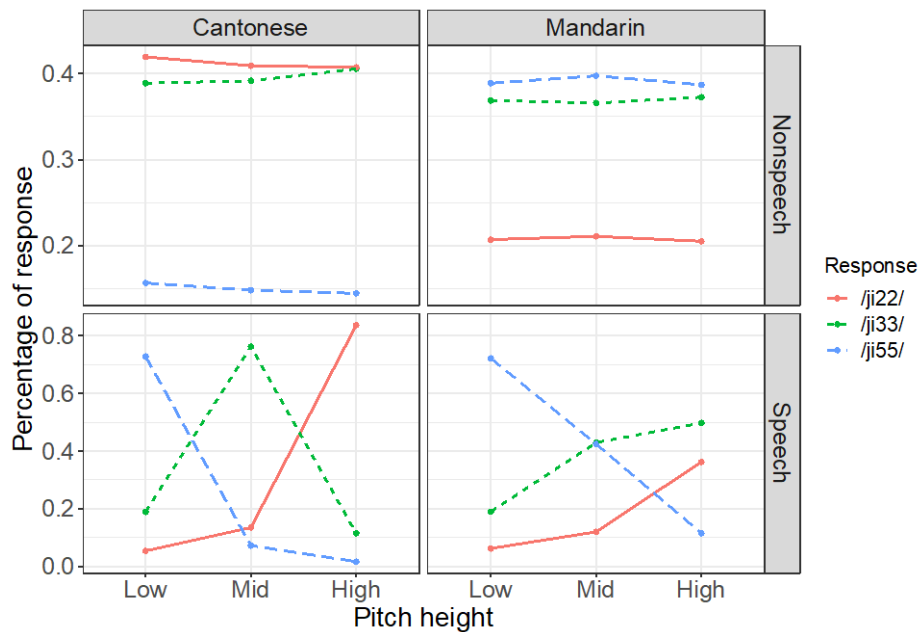
3

3.0 Results

4

Figure 2 illustrates the percentages of three response choices (i.e., 医 /ji55/ “doctor”, 意 /ji33/ “meaning”, and 二 /ji22/ “two”) in each experimental condition. To statistically evaluate how listeners used context cues to interpret target tones, a multinomial logistic regression model was fitted to all participants’ responses in the Cantonese-tone normalization task, using the *nnet* package (Venables & Ripley, 2002) in *R*. The dependent variable was *Response category* with three levels (/ji22/, /ji33/, and /ji55/). The predictors included *Pitch height* (low, mid, and high), *context* (nonspeech and speech), *group* (Cantonese and Mandarin), and their possible two-way and three-way interactions. The by-subject intercept was also included. All the predictors were centered around 0 (*group*: -1 for Cantonese and 1 for Mandarin; *pitch height*: -1 for low, 0 for mid, and 1 for high; *context*: -1 for nonspeech and 1 for speech). The reference level for *response category* was set to /ji33/. The fitted model could explain significantly more variance than the baseline model that only had by-subject intercept (AIC difference = 4533.441, $p < 0.001$).

16



1

2 **Figure 2. The percentage of three responses in each condition in the Cantonese-tone**
 3 **normalization task**

4

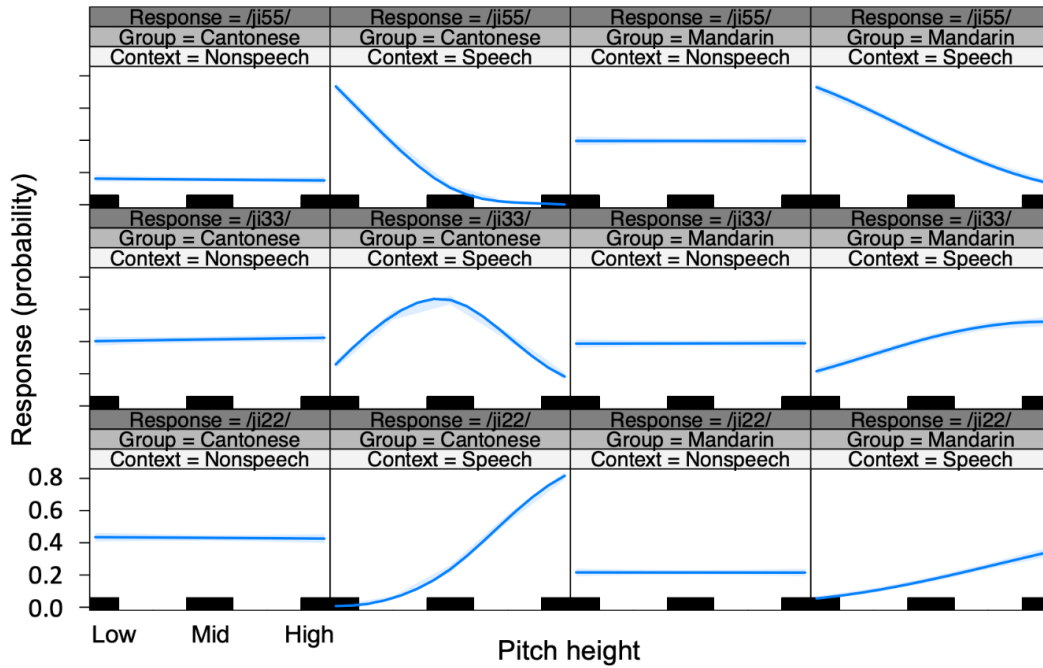
5 The statistics of the multinomial logistic regression model were summarized in Table
 6 1. The significant main effect of *context* indicated that participants gave significantly fewer
 7 /ji55/ responses (relative to /ji33/) ($\beta = -0.233, p < 0.001$) and fewer /ji22/ responses (relative
 8 to /ji33/) ($\beta = -0.347, p < 0.001$) in speech contexts than in nonspeech contexts. The significant
 9 main effect of *group* suggested that Mandarin learners, on average, provided more /ji55/
 10 responses (relative to /ji33/) ($\beta = 0.692, p < 0.001$) and fewer /ji22/ responses (relative to /ji33/)
 11 ($\beta = -0.127, p < 0.001$) than Cantonese speakers, indicating an influence of Mandarin T55. The
 12 significant main effect of *pitch height* indicated that participants on average gave significantly
 13 fewer /ji55/ responses (relative to /ji33/) ($\beta = -1.054, p < 0.001$) and more /ji22/ responses
 14 (relative to /ji33/) ($\beta = 0.739, p < 0.001$) when pitch height changed from low to high,
 15 illustrating a contrastive context effect. However, this effect varied between groups, as
 16 indicated by the significant *pitch height by group* interaction. With an increase in context pitch
 17 height from low to high, Mandarin learners gave more /ji55/ responses (relative to /ji33/) ($\beta =$
 18 $0.415, p < 0.001$) and fewer /ji22/ responses (relative to /ji33/) ($\beta = -0.507, p < 0.001$) than
 19 Cantonese speakers, suggesting a diminished contrastive context effect among Mandarin
 20 learners. Additionally, the contrastive context effect varied between speech and nonspeech
 21 contexts, as revealed by the significant *pitch height by context* interaction. When pitch height
 22 increased from low to high, participants gave fewer /ji55/ responses (relative to /ji33/) ($\beta = -$

1 1.021, $p < 0.001$) and more /ji22/ responses (relative to /ji33/) ($\beta = 0.759, p < 0.001$) in speech
 2 contexts compared to nonspeech contexts, indicating a more pronounced contrastive context
 3 effect in speech contexts. There was also a significant *context* by *group* interaction. When
 4 contexts changed from nonspeech to speech, Mandarin learners, compared to Cantonese
 5 speakers, gave more /ji55/ responses (relative to /ji33/) ($\beta = 0.197, p < 0.001$) and more /ji22/
 6 responses (relative to /ji33/) ($\beta = 0.188, p < 0.001$).

7 **Table 1** *The multinomial logistic regression model for the Cantonese-tone normalization*
 8 *task.*

	/ji55/ (relative to /ji33/)				
	Estimate	SE	95% CI	<i>z</i>	<i>p</i>
Intercept	-0.355	0.016	[-0.386 -0.323]	-22.255	<0.001
pitch height	-1.054	0.039	[-1.13 -0.977]	-27.014	<0.001
context	-0.233	0.032	[-0.296 -0.171]	-7.316	<0.001
group	0.692	0.032	[0.63 0.755]	21.717	<0.001
pitch height : context	-1.021	0.039	[-1.098 -0.945]	-26.184	<0.001
pitch height : group	0.415	0.039	[0.339 0.491]	10.639	<0.001
context : group	0.197	0.032	[0.134 0.259]	6.166	<0.001
pitch height : context : group	0.386	0.039	[0.31 0.463]	9.904	<0.001
	/ji22/ (relative to /ji33/)				
	Estimate	SE	95% CI	<i>z</i>	<i>p</i>
Intercept	-0.31	0.014	[-0.337 -0.283]	-22.525	<0.001
pitch height	0.739	0.036	[0.668 0.81]	20.378	<0.001
context	-0.347	0.028	[-0.401 -0.293]	-12.605	<0.001
group	-0.127	0.028	[-0.181 -0.073]	-4.606	<0.001
pitch height : context	0.759	0.036	[0.688 0.83]	20.933	<0.001
pitch height : group	-0.507	0.036	[-0.578 -0.436]	-13.995	<0.001
context : group	0.188	0.028	[0.134 0.242]	6.816	<0.001
pitch height : context : group	-0.523	0.036	[-0.594 -0.452]	-14.421	<0.001

9



1

2 **Figure 3. The visualization of pitch height by context by group interaction in the Cantonese-**
 3 **tone normalization task.**

4 The significant *pitch height by context by group* interaction was depicted in Figure 3,
 5 using the *effects* package (Fox & Hong, 2009) in R. It shows the predicted probability of each
 6 response in different experimental conditions. First of all, to ascertain how context pitch
 7 heights affect the tone perception in each group, a post hoc analysis with Bonferroni adjustment
 8 was conducted by comparing each response across three pitch heights for each group, using
 9 *emmeans* package (Lenth, 2019) in R. The analysis revealed significant context-dependent tone
 10 perception in both groups within speech contexts. Specifically, native Cantonese speakers gave
 11 more /ji55/ responses in the low-f₀ speech contexts (*prob* = 0.75; *SE* = 0.01) compared to high-
 12 f₀ (0.017; 0.004) and mid-f₀ speech contexts (0.0745; 0.008) (*ps* < 0.001). Native Cantonese
 13 speakers gave more /ji33/ responses in mid-f₀ speech contexts (0.786; 0.0126) than high-f₀
 14 (0.118; 0.01) and low-f₀ speech contexts (0.194; 0.0122) (*ps* < 0.001). Native Cantonese
 15 speakers gave more /ji22/ responses in high-f₀ speech contexts (0.8641; 0.01) than in mid-f₀
 16 (0.139; 0.01) and low-f₀ speech contexts (0.055; 0.007) (*ps* < 0.001). In speech contexts,
 17 Mandarin learners gave more /ji55/ responses in the low-f₀ condition (0.718; 0.013) than in
 18 high-f₀ (0.127; 0.01) and mid-f₀ condition (0.426; 0.015) (*ps* < 0.001). Mandarin learners gave
 19 more /ji22/ responses in high-f₀ speech contexts (0.359; 0.014) than in speech contexts with
 20 mid f₀ (0.128; 0.01) and low f₀ (0.076; 0.008) (*ps* < 0.001). Mandarin learners gave more /ji33/
 21 responses in mid-f₀ speech contexts (0.445; 0.015) contexts than in speech contexts with low

1 f0 (0.205; 0.012) ($p < 0.001$). However, Mandarin learners' /ji33/ responses did not differ
2 significantly between mid-f0 (0.445; 0.015) and high-f0 (0.5138; 0.015) speech contexts ($p >$
3 0.05), the only pair deviating from the contrastive context effect. In non-speech contexts, no
4 significant response differences were observed when context pitch heights varied from low to
5 high for either group ($ps > 0.05$).

6 Second, to elucidate group differences, post hoc analysis was conducted on the
7 significant *pitch height by context by group* interaction by comparing responses between the
8 two groups across each pitch height condition. Significant group differences were observed in
9 both speech- and nonspeech-context conditions. In speech contexts, Mandarin learners gave
10 fewer expected responses than Cantonese speakers. Specifically, in high-f0 speech contexts,
11 Mandarin speakers gave fewer /ji22/ responses (the expected response in high-f0 contexts) than
12 Cantonese speakers ($\beta = - 0.51, SE = 0.02, p < 0.001$). This was also the case in mid-f0 speech
13 contexts: Mandarin speakers gave fewer /ji33/ responses than Cantonese speakers ($\beta = - 0.34,$
14 $SE = 0.02, p < 0.001$). Only In low-f0 speech contexts, /ji55/ responses from Mandarin learners
15 and native Cantonese speakers were comparable ($\beta = - 0.032, SE = 0.02, p = 1$). In nonspeech
16 contexts, Mandarin learners more frequently perceived the target words (i.e, /ji33/ from
17 multiple speakers) as /ji55/ ($prob = 0.396; SE = 0.015$) and /ji33/ (0.388; 0.015) rather than
18 /ji22/ (0.216; 0.013) ($ps < 0.001$), but native Cantonese speakers were more likely perceive the
19 target words (i.e, /ji33/ from multiple speakers) as /ji33/ (0.403; 0.015) and /ji22/ (0.431; 0.015)
20 rather than /ji55/ (0.156; 0.011) ($ps < 0.001$).

21 In sum, the results of the Cantonese-tone normalization task suggested that both groups
22 were not affected by the pitch manipulation of nonspeech contexts, but in the speech-context
23 condition, both groups' tone perception changed significantly according to the pitch heights of
24 the contexts, indicating an extrinsic normalization process in Mandarin learners and Cantonese
25 speakers. However, Mandarin learners gave significantly fewer expected responses than native
26 Cantonese speakers especially in mid-f0 speech contexts (0.445 vs. 0.786), and in high-f0
27 speech contexts (0.359 vs. 0.8641), indicating a smaller normalization effect in Mandarin
28 learners.

29 **EXPERIMENT 2: THE FACTORS THAT AFFECT EXTRINSIC NORMALIZATION** 30 **OF L2 SPEECH**

1 The significant but smaller extrinsic normalization of Cantonese tones in Mandarin
2 learners revealed in Experiment 1 suggested that extrinsic normalization process is most likely
3 acquired gradually in L2 acquisition and cannot be directly transferred from L1 to L2 (see
4 section 6.1 for more discussion). This finding raises the question of which factors affect the
5 development of L2 extrinsic normalization. According to the acoustic-phonemic mapping
6 mechanism, learners' familiarity with the Cantonese tonal system is the most probable factor
7 influencing the development of L2 extrinsic normalization. Mandarin learners who were
8 sophisticated in Cantonese tones were more likely to establish a precise acoustic-phonemic
9 mapping to recalibrate ambiguous target speech cues. This is partially supported by studies
10 showing that the development of speech normalization in children is constrained by the
11 maturation of the phonological system. For examples, although 6-month-old children can
12 recognize phonemes from different talkers (Kuhl, 1976), this ability still do not achieve the
13 adult-like level at the age of six (Bent & Atagi, 2017). Mandarin-speaking children cannot use
14 context cues to identify ambiguous tones until the age of six (F. Chen et al., 2023), which is
15 consistent with their maturation in the categorical perception of Mandarin tones (F. Chen et al.,
16 2017; Feng & Peng, 2023).

17 Aside from acoustic and phonemic information, higher level linguistic knowledge also
18 contributes to the normalization process. Meaningful contexts triggered 11% higher
19 normalization effect than contexts composed of meaningless word sequence (C. Zhang et al.,
20 2015). Therefore, overall L2 proficiency may also affect the development of L2 speech
21 normalization process. L2 learners' insufficient knowledge of the sociolinguistic variation of
22 L2 was correlated with their poor performance in accommodating high-variability L2 speech
23 (Tamati & Pisoni, 2014). L2 immersion, in such a condition, is essential to boost the knowledge
24 of L2 sociolinguistic variation. Only 21-day Spanish exposure led to an improvement in
25 acquiring the L2 fine-phonetic details associated with stop voicing (Casillas, 2020). Long-term
26 Finnish exposure even enabled L2 learners to distinguish the Finnish vowel contrast at the pre-
27 attentive level (Winkler et al., 1999). The social contacts with native speakers significantly
28 affected how the regional variant of Spanish aspirated /s/ was identified (Schmidt, 2015). L2
29 learners of English who lived in America for several years showed much better talker
30 accommodation in identifying English words than L2 learners who were never exposed to a
31 native English environment (Antoniou et al., 2015). Therefore, L2 immersion probably also
32 boosts the accommodation of L2 speech variability at suprasegmental level.

1 Experiment 2 moved one step forward and investigated how L2 phonological
2 proficiency, overall L2 proficiency, and L2 immersion would affect the development of L2
3 extrinsic normalization process. The proficiency in the Cantonese tonal system was assessed
4 through a Cantonese tone identification task and a Cantonese tone production task. The overall
5 Cantonese proficiency and Cantonese immersion were evaluated using the language history
6 questionnaire (LHQ3) (Li et al., 2020). All these factors were hypothesized to significantly
7 affect Mandarin learners' normalization of Cantonese tones, with proficiency in the Cantonese
8 tonal system probably having a larger effect.

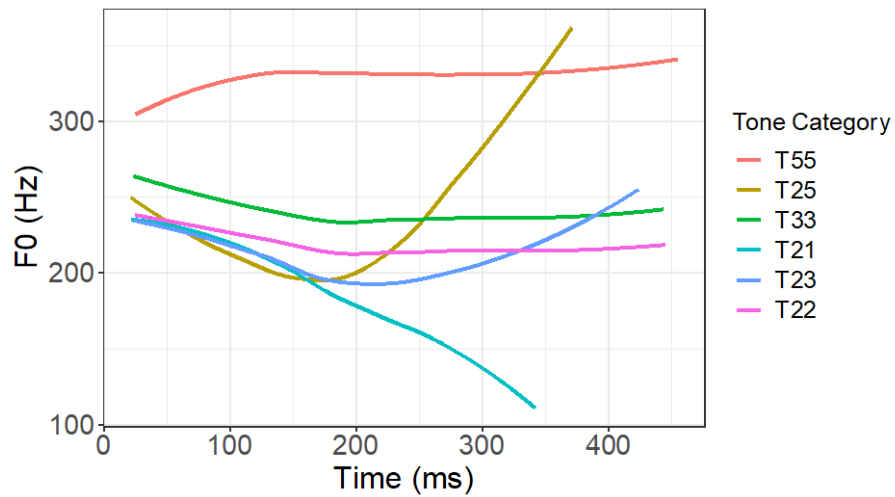
9 **4.0 Material and methods**

10 *4.1 Participants*

11 The Mandarin learners of Cantonese who took part in Experiment 1 were invited back
12 to participate in a Cantonese tone identification task and a Cantonese tone production task to
13 assess their proficiency in Cantonese tones. In addition, they were asked to complete the LHQ
14 3 to report their Cantonese acquisition and usage (Li et al., 2020). Finally, out of the 30
15 Mandarin participants in Experiment 1, 27 took part in Experiment 2.

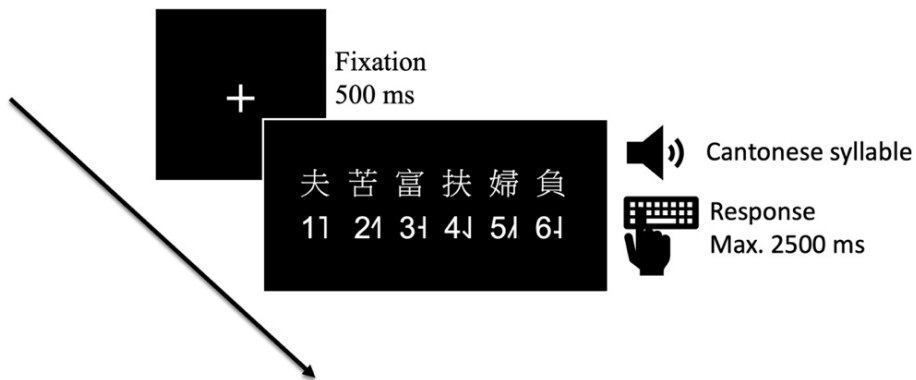
16 *4.2 The Cantonese tone identification task*

17 The stimuli used in the Cantonese tone identification task were identical to those used
18 in a previous study by K. Zhang et al., (2018). All stimuli (see Table 2) were produced by a
19 native Cantonese female speaker. The f0 trajectories of six Cantonese tones are depicted in
20 Figure 4, generated from stimuli using /jan/ as the base syllable. The trial procedure is
21 illustrated in Figure 5. During each trial, after a 500-ms fixation, a Cantonese syllable was
22 presented bilaterally to the subjects through inserted earphones. After the auditory stimuli, a
23 picture containing six Cantonese words, their corresponding tone letters, and tone numbers was
24 displayed on the screen. The six words on the screen had the same base syllable but different
25 tones, and the subjects were required to press the corresponding key from 1 to 6 on the
26 keyboard to indicate which word they heard. The next trial appeared after the subjects' response
27 or when the maximum response time of 2500 ms was reached. Thirty-six Cantonese syllables
28 (see Table 2) covering six Cantonese long tones were repeated five times and presented
29 randomly to subjects.



1
2

Figure 4. The f_0 trajectories of six Cantonese tones.



3
4

Figure 5. The trial procedure of the Cantonese tone identification task

5

Table 2 Thirty-six Cantonese tonal syllables

	/fan/	/fu/	/jan/	/ji/	/se/	/si/
T55	婚	夫	因	醫	些	詩
T25	粉	苦	隱	倚	寫	史
T33	訓	富	印	意	卸	嗜
T21	焚	扶	人	兒	蛇	時
T23	奮	婦	引	耳	社	市
T22	份	負	孕	二	射	事

6 4.3 The Cantonese tone production task

1 The Cantonese tone production task utilized in the present study was identical to the
2 Cantonese tone production task used in K. Zhang et al. (2018). The 36 Cantonese words listed
3 in Table 2 were repeated twice and presented in a random order. During each trial, a traditional
4 Chinese character, its Jyutping, and its tone letter were displayed on the screen. Mandarin
5 learners were instructed to produce the Cantonese words naturally and clearly, and their
6 production was recorded using Praat software installed on a Macbook Pro laptop in a
7 soundproof booth.

8 *4.4 LHQ 3*

9 The Mandarin learners were instructed to complete the LHQ3 in the laboratory (Li et
10 al., 2020). The LHQ3 is a widely used online tool developed by Li et al. (2020) for assessing
11 the linguistic background and language proficiency of L2 learners. It is based on frequently
12 asked questions from published studies, covering topics such as the age of acquisition, duration
13 of language usage, self-reported language proficiency, and scores from standardized language
14 proficiency tests. The full version of LHQ contains 27 questions, but it also offers users the
15 flexibility to select questions of interest and create customized versions. The present study
16 utilized this customized version, including questions related to the age of acquisition, years of
17 usage, self-rated language proficiency, self-rated strength of foreign accent, hours spent using
18 each language per day across various activities and with different people, and the frequency of
19 language mixing. The Mandarin learners were asked to report their acquisition and usage of
20 both their L1 (Mandarin) and Cantonese. Other L2 languages, such as English, were not
21 included.

22 **5.0 Results**

23 The objective of Experiment 2 was to investigate how various factors influenced the
24 Mandarin learners' normalization process. Consequently, the results section only reported
25 scores that assessed the Cantonese tone proficiency, overall L2 proficiency, and L2 immersion,
26 without conducting any further statistical analyses within each task.

27 *5.1 Cantonese tone identification task*

28 The accuracy rate of the Cantonese tone identification task was averaged across six
29 base syllables and plotted in Figure 5 (a). Each dot in Figure 5 (a) represents one Mandarin

1 learner's result. The confusion matrix for Mandarin learners' identification of Cantonese tones
 2 is presented in Table 3. The data indicate a strong L1 influence; Mandarin learners tend to
 3 identify T33 as T55, which is also present in Mandarin, rather than T22, which is acoustically
 4 more similar to T33 (refer to Figure 4). Even 17% of T22 stimuli were perceived as T55.

5 **Table 3** *The confusion matrix of Mandarin learners' Cantonese tone identification*

Response \ Stimuli	T55	T25	T33	T21	T23	T22
T55	83%	3%	30%	2%	2%	17%
T25	4%	62%	5%	3%	47%	6%
T33	9%	4%	49%	7%	5%	46%
T21	1%	7%	3%	61%	9%	6%
T23	1%	22%	6%	7%	32%	7%
T22	1%	1%	6%	18%	2%	16%

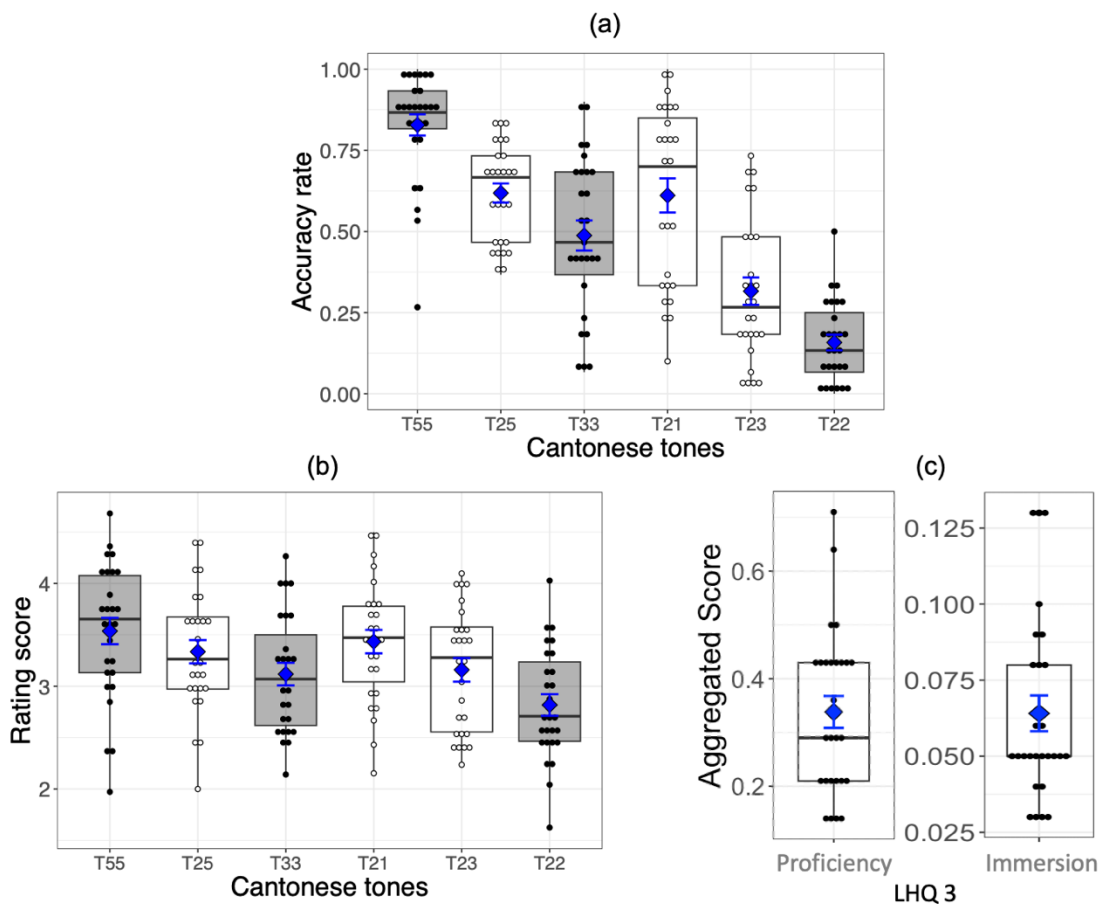
6 *5.2 Cantonese tone production task*

7 Mandarin learners' Cantonese tone production was evaluated by six native Cantonese
 8 speakers who specialized in linguistic-related subjects and had completed at least one course
 9 in phonetics. Tokens from the same Mandarin learner were presented in one block during the
 10 evaluation. Two tokens (/fan55/ and /fan25/) from native Cantonese speakers (one male and
 11 one female) were used as fillers (matched in gender in each block) to monitor the rater's
 12 performance. Finally, there were 27 blocks and each block contained 74 tokens (36 syllables ×
 13 2 repetitions + 2 fillers). Twenty-seven blocks were presented randomly, and 74 tokens within
 14 each block were also presented in random order. Six native Cantonese speakers rated each
 15 token on a 1-5 Likert scale, with 1 for the worst and most non-native-like pronunciation and 5
 16 for the best and most native-like pronunciation. They were asked to insist on their standards
 17 throughout the whole evaluation.

18 All raters completed the task with great care, as demonstrated by their high ratings of
 19 the two fillers produced by native Cantonese speakers (mean = 4.67 out of 5). The Interrater
 20 reliability was high among six raters [ICC (2, 6) = 0.96, $p < 0.01$], and thus all six raters' rating
 21 scores were included. The rating scores for each Cantonese tone were averaged across six base
 22 syllables and six raters. The result of each tone for each Mandarin learner were represented by
 23 a single dot in Figure 5 (b).

1 5.3 The aggregate scores from LHQ 3

2 LHQ 3 utilizes aggregated scores to represent the proficiency and immersion levels of
3 participants in the languages they acquired. The scores are derived from four components:
4 reading, writing, listening, and speaking (Li et al., 2020). Proficiency level is the weighted sum
5 of participants' self-evaluation of their proficiency levels on each component, while immersion
6 level is based on their ages, age of acquisition, and duration of use of the language in each
7 component. For the detailed calculation formula, please see [https://lhq-](https://lhq-blclab.org/static/docs/aggregate-scores.html)
8 [blclab.org/static/docs/aggregate-scores.html](https://lhq-blclab.org/static/docs/aggregate-scores.html). Although Mandarin and Hong Kong Cantonese
9 adopt distinct Chinese writing systems, people educated in one writing system can somewhat
10 understand another due to the similarities. Therefore, reading and writing were excluded while
11 calculating aggregated scores. Figure 5 (c) illustrates the overall Cantonese proficiency and
12 Cantonese immersion aggregated scores for each Mandarin learner.



13

14 **Figure 6. (a) Mandarin learners' accuracy rate in the Cantonese tone identification task, (b)**
15 **the rating score of Mandarin learners' Cantonese tone production, and (c) the aggregated**

1 scores from LHQ 3. The boxes in black in panel (a) and (b) represent the tonal categories
2 used in E1.

3 5.4. The effects of different factors on Mandarin learners' Cantonese tone normalization
4 process

5 A generalized linear mixed-effect model was fitted using the *lme4* package (Bates et
6 al., 2015) in R to evaluate how Mandarin learners' Cantonese tone normalization was affected
7 by their Cantonese tone proficiency, overall L2 proficiency and L2 immersion. According to
8 the contrastive context effect, listeners were anticipated to give /ji55/ responses for low-f0
9 contexts, /ji33/ responses for mid-f0 contexts, and /ji22/ responses for high-f0 contexts. Thus,
10 a normalization accuracy of 1 was assigned to each trial when participants gave the expected
11 response and 0 otherwise. The dependent variable in the model was normalization accuracy,
12 whereas the fixed factors in the model were Cantonese tone identification accuracy, Cantonese
13 tone production score, L2 proficiency score and L2 immersion score. All continuous
14 independent variables were centered around 0. Only the by-subject intercept was included as a
15 random effect due to convergence issues.

16 Considering the correlation among the predictors (Cantonese tone perception score and
17 Cantonese tone production score: *Pearson's r* = 0.4; *p* < 0.05; L2 immersion score and
18 Cantonese tone production score: *r* = 0.41; *p* < 0.05; L2 proficiency score and L2 immersion
19 score: *r* = 0.65; *p* < 0.01), the variance inflation factor (VIF) was calculated for each predictor
20 in the logistic regression model to detect the multilinearity problem. All VIF values were less
21 than 5 (1.6 for Cantonese tone identification accuracy, 1.86 for Cantonese tone production
22 score, 1.9 for L2 proficiency score, and 2.75 for L2 immersion score), indicating a low risk of
23 multicollinearity, and thus all factors were retained in the model (James et al., 2021).

24 Results from the logistic regression analysis revealed that L2 immersion score
25 significantly predicted L2 normalization accuracy ($\beta = 9.16$, 95% CI = [1.75, 16.63], *SE* = 3.61,
26 *z* = 2.53, *p* = 0.01). Specifically, Mandarin learners with more extensive L2 immersion
27 demonstrated more effective normalization of talker variability. However, Cantonese tone
28 identification ($\beta = 0.78$, 95% CI = [- 0.48, 2.06], *SE* = 0.62, *z* = 1.26, *p* = 0.21), Cantonese tone
29 production ($\beta = - 0.03$, 95% CI = [- 0.37, 0.3], *SE* = 0.17, *z* = - 0.21, *p* = 0.84), and overall L2
30 proficiency score ($\beta = -0.08$, 95% CI = [- 1.31, 1.15], *SE* = 0.6, *z* = - 0.13, *p* = 0.89) did not
31 significantly impact Mandarin learners' Cantonese-tone normalization.

1 **6.0 Discussion.**

2 *6.1. Extrinsic normalization process is a language-dependent process.*

3 Extrinsic contexts are useful for listeners to normalize speech variability and achieve
4 perceptual constancy. To investigate whether L2 learners could utilize context cues to
5 normalize L2 speech variability as they did in their native languages, Experiment 1 in this study
6 examined Mandarin learners' ability to identify Cantonese level tones produced by multiple
7 speakers in nonspeech and speech contexts. Neither Mandarin learners nor native Cantonese
8 speakers were affected by the pitch manipulation in nonspeech contexts, confirming the
9 ineffectiveness of nonspeech contexts in accommodating speech variability. In the absence of
10 effective context cues (i.e., in the nonspeech-context condition), both groups were severely
11 confused T33 produced by multiple speakers with another two Cantonese level tones (see the
12 upper panels in Figure 2), highlighting the importance of context cues in Cantonese tone
13 perception. When effective context cues were available (i.e., in the speech-context condition),
14 the perception of the target tone changed significantly in response to the pitch heights of the
15 contexts in both groups. Specifically, native Cantonese speakers gave more /ji55/ responses in
16 low-f0 speech contexts than in high-f0 speech contexts (0.75 vs. 0.017) and more /ji22/
17 responses in high-f0 speech contexts than in low-f0 speech contexts (0.86 vs. 0.06). More
18 importantly, Mandarin learners demonstrated similar perceptual patterns, giving more /ji55/
19 responses when contexts changed from high f0 to low f0 (0.72 vs. 0.13) and more /ji22/
20 responses when contexts changed from low f0 to high f0 (0.4 vs. 0.08). These results suggest
21 that Mandarin learners exhibited a significant extrinsic normalization process during
22 Cantonese tone perception, similar to native Cantonese speakers, when effective context cues
23 are available.

24 However, the normalization effect was significantly smaller in Mandarin learners than
25 in native Cantonese speakers. Except in low-f0 contexts, Mandarin learners gave significantly
26 fewer expected responses based on the contrastive contexts effect. Specifically, they gave
27 fewer /ji33/ responses in mid-f0 speech contexts than native Cantonese speakers did (0.45 vs.
28 0.79), and fewer /ji22/ responses in high-f0 speech contexts than native Cantonese speakers
29 did (0.36 vs. 0.86). The significantly smaller normalization effect suggested that Mandarin
30 learners could not use context cues as effectively as native Cantonese speakers could. This
31 finding supported that extrinsic normalization requires language-specific knowledge and thus
32 should be acquired gradually during L2 learning, rather than directly generalized from L1 to

1 L2. The finding of Experiment 1 was also consistent with previous research reporting that
2 contexts composed of native phonemes were more effective than those composed of nonnative
3 phonemes in consonant normalization (Kang et al., 2016) and vowel normalization (K. Zhang
4 & Peng, 2021). Together with the findings from the developmental study of child language
5 which showed that children only at the age of six started to show extrinsic normalization
6 process (F. Chen et al., 2023), the present study suggested that extrinsic normalization was a
7 language-dependent process and was developed gradually in language acquisition of both L1
8 and L2. The importance of language-specific knowledge in extrinsic normalization further
9 supports the notion that extrinsic normalization is implemented at the speech-specific level. To
10 our knowledge, this is the first study to address this dispute (i.e., the general-auditory level
11 process vs. the speech-specific level process) from the perspective of L2 acquisition, thereby
12 providing new evidence to shed light on the cognitive mechanisms underlying extrinsic
13 normalization.

14 *6.2. Language immersion facilitates the extrinsic normalization of L2 speech.*

15 Following up the results of Experiment 1, Experiment 2 aimed to examine how L2
16 phonological proficiency, overall L2 proficiency, and L2 immersion affected the extrinsic
17 normalization process in L2. The logistic regression model with all these factors showed that
18 only L2 immersion had a significant and positive impact on the extrinsic normalization
19 accuracy. This means that L2 learners with more extensive L2 immersion can use context cues
20 to cope with talker variability more effectively. Previous studies have found that language
21 immersion enhances L2 learners' ability to identify L2 words from multiple talkers (Antoniou
22 et al., 2015) and native speakers' ability to perceive nonnative-accented speech (Bradlow &
23 Bent, 2008; Xie et al., 2017; Zheng & Samuel, 2020). These findings have been explained by
24 the perceptual re-organization hypothesis, which proposes that language exposure leads to the
25 updating of mental representations of linguistic units based on speech inputs, by adjusting
26 category boundaries or internal category structure (Bradlow & Bent, 2008; Norris et al., 2003;
27 Xie et al., 2017). Apart from retuning the mental representations of intrinsic cues, the present
28 study suggests that the capacity to utilize extrinsic contextual cues in speech perception serves
29 as an additional factor modulating the interplay between language immersion and speech
30 normalization. The extrinsic normalization may benefit from L2 immersion due to perceptual
31 practice. Mandarin learners have to cope with talker variability in Cantonese tones when they
32 communicate with native speakers, so longer Cantonese immersion also implies more practice

1 in the extrinsic normalization process. This frequent practice enables L2 learners to improve
2 their ability to use context cues to accommodate speech variability and thus achieve better
3 normalization outcomes.

4 The regression analysis of Experiment 2 in this study showed that Cantonese tone
5 proficiency did not significantly predict the extrinsic normalization results. This contradicted
6 previous studies that found that familiarity with context phonemes facilitates the extrinsic
7 normalization process. For instance, Kang et al. (2016) demonstrated that the rounded vowel
8 /y/ in French did not influence English speakers' perception of the /s/-/ʃ/ continuum, but French
9 speakers exhibited the contrastive context effect in such a condition. K. Zhang & Peng (2021)
10 reported a significant reduction of the context effect on English speakers' vowel normalization
11 in nonnative speech contexts compared to native speech contexts. A possible explanation for
12 this discrepancy is that phonological knowledge is indeed necessary for speech normalization,
13 but it may have a limited impact. Speech normalization might involve multiple processing
14 stages (Xie et al., 2023). Besides perceptually adjusting the perceived signals against the
15 acoustic-phonemic mapping at the first stage, speech normalization also includes a second
16 stage — cognitive adjustment — where listeners make final decisions on the perceived signals
17 with high-level information (Bosker et al., 2017). Semantic and syntactic information might
18 play a role at this stage, as evidenced by C. Zhang et al. (2015), which showed that meaningful
19 contexts elicited an 11% higher normalization effect than contexts composed of meaningless
20 word sequences. The high-level information (including but not limited to semantic and
21 syntactic information) varied greatly among Mandarin learners in this study as indicated by
22 their overall L2 proficiencies. In addition, the cognitive adjustment also relied on listeners'
23 cognitive ability, which also varied individually (Carroll & Maxwell, 1979; Ou & Law, 2017).
24 Thus, the potential between-subject variability in the cognitive-adjustment stage might obscure
25 the effect of Cantonese tone proficiency on the final normalization results.

26 Overall L2 proficiency was also not a significant predictor of the L2 learners'
27 normalization process. However, an exploratory logistic regression analysis that only included
28 L2 proficiency score, Cantonese tone perception score, and Cantonese tone production score
29 as fixed factors and normalization accuracy as the dependent variable showed a significant
30 main effect of L2 proficiency score ($\beta = 0.96$, 95% CI = [- 0.02, 1.95], $SE = 0.49$, $z = 1.98$, p
31 = 0.048). This effect disappeared when L2 immersion score was added to the model. Despite
32 the significant correlation between L2 proficiency score and L2 immersion score ($r = 0.65$, p

1 < 0.01), the insignificance of L2 proficiency score in the larger model might not result from
2 multicollinearity as the VIF value for each factor was less than 5. It is more likely that L2
3 proficiency score and L2 immersion score measure a similar underlying construct, leading to a
4 situation where they share variance in explaining the dependent variable (i.e., normalization
5 accuracy). In this case, when L2 immersion score is added to the model, it may account for
6 some of the variance that was previously attributed to the L2 proficiency score, thus making
7 the latter non-significant. This study assessed the overall L2 proficiency using the LHQ 3,
8 which calculated the aggregated scores of language proficiency based on self-rated proficiency
9 levels on listening and speaking. The self-rated proficiency may not be as reliable as standard
10 language tests since different people has different standards. Future studies can use standard
11 language tests to better measure overall L2 proficiency, which may explain more underlying
12 variance of the L2 normalization process that was not captured by L2 immersion.

13 *6.3 L1 phonological representations interact with immediate context cues in the L2 speech* 14 *normalization process.*

15 Mandarin learners' performance in Experiment 1 was strongly influenced by their L1
16 tonal system. In the absence of effective context cues (i.e., in the nonspeech-context condition),
17 native Cantonese speakers mostly confused T33 with T22 (around 43%) and occasionally with
18 T55 (around 16%), which was mainly due to the intrinsic f₀ similarity between T22 and T33
19 (Peng, 2006). However, unlike native Cantonese speakers, Mandarin learners more frequently
20 misperceived T33 as T55 (around 40%), which is also a Mandarin tone, but relatively less as
21 T22 (around 21%), whose intrinsic f₀ is closer to T33, showing a strong L1 influence. Similar
22 results can also be observed in Table 3 which showed the confusion matrix of Mandarin
23 learners' perception of Cantonese tones in isolation. According to the Perceptual Assimilation
24 Model, L2 learners tend to assimilate L2 phonemes to similar phonemes in their L1 and use
25 the mental representations of L1 phonemes to perceive the similar L2 phonemes (Best & Tyler,
26 2007). Although the screening results suggested that all Mandarin participants could produce
27 three Cantonese level tones, the mental representations of three Cantonese level tones were not
28 robust. In such a condition, they probably relied on the mental representation of Mandarin T55
29 to perceive the similar Cantonese level tones, resulting in a high proportion of confusion
30 between T33 and T55 (around 40%).

31 When effective context cues were available (i.e., in the speech-context condition),
32 Mandarin learners' Cantonese tone perception was influenced by both immediate context cues

1 and L1 tonal representations. A clear contrastive context effect was observed in Mandarin
2 learners, as they gave more T22 in high-f0 (35.9%) than in low-f0 (12.8%) speech contexts,
3 and more T55 in low-f0 (71.8%) than in high-f0 (12.7%) speech contexts. Meanwhile,
4 noticeable L1 influences were also observed. First, Mandarin learners overall gave more T55
5 responses in speech contexts than native Cantonese speakers did (see Figure 2). The interaction
6 between L1 influence and immediate context cues was particularly evident in high-f0 speech
7 contexts. L1 influence typically resulted in a higher number of T55 responses, but the high-f0
8 speech contexts theoretically led to more T22 responses. The interplay of L1 influence and
9 immediate context cues produced the greatest proportion of T33 responses within the high-f0
10 speech contexts among the Mandarin group. To our knowledge, this was the first study that
11 examined the influence of L1 tone information on L2 tone processing within the framework of
12 speech normalization. The joint effect of L1 phonological influence and immediate context in
13 the present study is consistent with S. Chen et al. (2022) who modeled mental representations
14 in speech normalization of prosodic cues. S. Chen et al. (2022) first estimated the statistical
15 distributions of acoustic cues associated with Cantonese level tones from a speech corpus and
16 then tested how the distributional information that listeners established based on experience
17 and the immediate context cues affected native Cantonese speakers' tone normalization. Their
18 results suggest that speech normalization was a complex process that integrated both the
19 distributional information of target tones and the immediate context information.

20 **7.0 Conclusions**

21 The present study tested if L2 learners can generalize the extrinsic normalization
22 process that they used in L1 to L2 to accommodate L2 speech variability by investigating
23 Mandarin learners' perception of Cantonese tones. A significant but weaker normalization
24 process in Mandarin learners than native Cantonese speakers indicates that the extrinsic
25 normalization process cannot be generalized directly from L1 to L2 and requires gradual
26 learning. Such a finding suggests that L2 learners and teachers should pay attention to the
27 training of this skill in L2 acquisition. The significant positive impact of language immersion
28 on L2 learners' extrinsic normalization process highlights the importance of language
29 immersion for developing L2 speech perceptual strategies. Future studies may include speakers
30 of nontonal languages to explore whether having a tonal-language background is an advantage
31 for L2 tone normalization.

32

1 **Acknowledgements**

2 This work was supported by the Research Grants Council of Hong Kong [GRF:
3 15607518].

4

5

1 Reference List

- 2 Ainsworth, W. A. (1975). Intrinsic and extrinsic factors in vowel judgments. In G. Fant & M.
3 A. A. Tatham (Eds.), *Auditory Analysis and Perception of Speech* (pp. 103–113).
4 Academic Press. <https://ci.nii.ac.jp/naid/10012630848/en/>
- 5 Antoniou, M., Wong, P. C. M., & Wang, S. (2015). The effect of intensified language
6 exposure on accommodating talker variability. *Journal of Speech, Language, and*
7 *Hearing Research*, 58(3), 722–727. https://doi.org/10.1044/2015_JSLHR-S-14-0259
- 8 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models
9 using lme4. *Journal of Statistical Software*, 67(1), 1–48.
10 <https://doi.org/10.18637/jss.v067.i01>
- 11 Bent, T., & Atagi, E. (2017). Perception of nonnative-accented sentences by 5-to 8-year-olds
12 and adults: The role of phonological processing skills. *Language and Speech*, 60(1),
13 110–122. <https://doi.org/10.1177/0023830916645374>
- 14 Bent, T., Kewley-Port, D., & Ferguson, S. H. (2010). Across-talker effects on non-native
15 listeners' vowel perception in noise. *The Journal of the Acoustical Society of America*,
16 128(5), 3142–3151. <https://doi.org/10.1121/1.3493428>
- 17 Best, C. T., & Tyler, M. D. (2007). Nonnative and second-language speech perception:
18 Commonalities and complementarities. In O.-S. Bohn & M. J. Munro (Eds.), *Language*
19 *experience in second language speech learning: In honor of James Emil Flege* (pp. 13–
20 34). John Benjamins Publishing Company.
- 21 Boersma, P., & Weenink, D. (2023). *Praat: doing phonetics by computer [Computer*
22 *program]. Version 6.3.09, retrieved 2 March 2023 from <http://www.praat.org/>.*
- 23 Bosker, H. R., Reinisch, E., & Sjerps, M. J. (2017). Cognitive load makes speech sound fast,
24 but does not modulate acoustic context effects. *Journal of Memory and Language*, 94,
25 166–176. <https://doi.org/10.1016/j.jml.2016.12.002>
- 26 Bradlow, A. R., & Bent, T. (2008). Perceptual adaptation to non-native speech. *Cognition*,
27 106(2), 707–729. <https://doi.org/10.1016/j.cognition.2007.04.005>
- 28 Bradlow, A. R., & Pisoni, D. B. (1999). Recognition of spoken words by native and non-
29 native listeners: Talker-, listener-, and item-related factors. *The Journal of the*
30 *Acoustical Society of America*, 106(4), 2074–2085. <https://doi.org/10.1121/1.427952>
- 31 Campbell, K. L., & Tyler, L. K. (2018). Language-related domain-specific and domain-

- 1 general systems in the human brain. *Current Opinion in Behavioral Sciences*, 21, 132–
2 137. <https://doi.org/10.1016/j.cobeha.2018.04.008>
- 3 Carroll, J. B., & Maxwell, S. E. (1979). Individual Differences in Cognitive Abilities. *Annual*
4 *Review of Psychology*, 30(1), 603–640.
5 <https://doi.org/10.1146/annurev.ps.30.020179.003131>
- 6 Casillas, J. V. (2020). The Longitudinal Development of Fine-Phonetic Detail: Stop
7 Production in a Domestic Immersion Program. *Language Learning*, 70(3), 768–806.
8 <https://doi.org/10.1111/lang.12392>
- 9 Chen, F., Peng, G., Yan, N., & Wang, L. (2017). The development of categorical perception
10 of Mandarin tones in four-to seven-year-old children. *Journal of Child Language*, 44(6),
11 1413–1434. <https://doi.org/10.1017/S0305000916000581>
- 12 Chen, F., Zhang, K., Guo, Q., & Lv, J. (2023). Development of achieving onstancy in lexical
13 tone identification with contextual cues. *Journal of Speech, Language, and Hearing*
14 *Research*, 1–17. https://doi.org/10.1044/2022_JSLHR-22-00257
- 15 Chen, S., Zhang, C., Lau, P., Yang, Y., & Li, B. (2022). Modelling representations in speech
16 normalization of prosodic cues. *Scientific Reports*, 12(1), 1–21.
17 <https://doi.org/10.1038/s41598-022-18838-w>
- 18 Feng, Y., & Peng, G. (2023). Development of categorical speech perception in Mandarin-
19 speaking children and adolescents. *Child Development*, 94(1), 28–43.
20 <https://doi.org/10.1111/cdev.13837>
- 21 Fox, J., & Hong, J. (2009). Effect Displays in R for Multinomial and Proportional-Odds
22 Logit Models: Extensions to the effects Package. *Journal of Statistical Software*, 32(1).
23 <https://doi.org/10.18637/jss.v032.i01>
- 24 Francis, A. L., Ciocca, V., Wong, N. K. Y., Leung, W. H. Y., & Chu, P. C. Y. (2006).
25 Extrinsic context affects perceptual normalization of lexical tone. *The Journal of the*
26 *Acoustical Society of America*, 119(3), 1712–1726. <https://doi.org/10.1121/1.2149768>
- 27 Holt, L. L., & Lotto, A. J. (2002). Behavioral examinations of the level of auditory
28 processing of speech context effects. *Hearing Research*, 167(1), 156–169.
29 [https://doi.org/10.1016/S0378-5955\(02\)00383-0](https://doi.org/10.1016/S0378-5955(02)00383-0)
- 30 Holt, L. L., Lotto, A. J., & Kluender, K. R. (2001). Influence of fundamental frequency on
31 stop-consonant voicing perception: A case of learned covariation or auditory

- 1 enhancement? *The Journal of the Acoustical Society of America*, 109(2), 764–774.
2 <https://doi.org/10.1121/1.1339825>
- 3 James, G., Witten, D., Hastie, T., & Tibshirani, R. (2021). *An Introduction to Statistical*
4 *Learning*. Springer US. <https://doi.org/10.1007/978-1-0716-1418-1>
- 5 Johnson, K. (2005). Speaker Normalization in Speech Perception. In D. B. Pisoni & R. E.
6 Remez (Eds.), *The Handbook of Speech Perception* (pp. 363–389). Blackwell
7 Publishing. <https://doi.org/10.1002/9780470757024.ch15>
- 8 Kang, S., Johnson, K., & Finley, G. (2016). Effects of native language on compensation for
9 coarticulation. *Speech Communication*, 77, 84–100.
10 <https://doi.org/10.1016/j.specom.2015.12.005>
- 11 Kapadia, A. M., & Perrachione, T. K. (2020). Selecting among competing models of talker
12 adaptation: Attention, cognition, and memory in speech processing efficiency.
13 *Cognition*, 204(February), 104393. <https://doi.org/10.1016/j.cognition.2020.104393>
- 14 Kuhl, P. K. (1976). Speech perception in early infancy: Perceptual constancy for vowel
15 categories. *The Journal of the Acoustical Society of America*, 60(S1), S90–S91.
16 <https://doi.org/10.1121/1.2003600>
- 17 Lecumberri, M. L. G., Cooke, M., & Cutler, A. (2010). Non-native speech perception in
18 adverse conditions: A review. *Speech Communication*, 52(11–12), 864–886.
19 <https://doi.org/10.1016/j.specom.2010.08.014>
- 20 Lee, C.-Y., Zhang, Y., Li, X., Tao, L., & Bond, Z. S. (2012). Effects of speaker variability
21 and noise on Mandarin fricative identification by native and non-native listeners. *The*
22 *Journal of the Acoustical Society of America*, 132(2), 1130–1140.
23 <https://doi.org/10.1121/1.4730883>
- 24 Lee, C. Y., Tao, L., & Bond, Z. S. (2009). Speaker variability and context in the
25 identification of fragmented Mandarin tones by native and non-native listeners. *Journal*
26 *of Phonetics*, 37(1), 1–15. <https://doi.org/10.1016/j.wocn.2008.08.001>
- 27 Lee, C. Y., Tao, L., & Bond, Z. S. (2013). Effects of speaker variability and noise on
28 Mandarin tone identification by native and non-native listeners. *Speech, Language and*
29 *Hearing*, 16(1), 46–54. <https://doi.org/10.1179/2050571X12Z.0000000003>
- 30 Lenth, R. (2019). Emmeans: estimated marginal means. In *R package version 1.4.2*.
31 <https://cran.r-project.org/package=emmeans>

- 1 Li, P., Zhang, F., Yu, A., & Zhao, X. (2020). Language History Questionnaire (LHQ3): An
2 enhanced tool for assessing multilingual experience. *Bilingualism*, 23(5), 938–944.
3 <https://doi.org/10.1017/S1366728918001153>
- 4 Liberman, A. M., Cooper, F. S., Shankweiler, D. P., & Studdert-Kennedy, M. (1967).
5 Perception of the speech code. *Psychological Review*, 74(6), 431–461.
- 6 Moore, C. B., & Jongman, A. (1997). Speaker normalization in the perception of Mandarin
7 Chinese tones. *The Journal of the Acoustical Society of America*, 102(3), 1864–1877.
8 <https://doi.org/10.1121/1.420092>
- 9 Nearey, T. M. (1989). Static, dynamic, and relational properties in vowel perception. *The*
10 *Journal of the Acoustical Society of America*, 85(5), 2088–2113.
11 <https://doi.org/10.1121/1.397861>
- 12 Norris, D., McQueen, J. M., & Cutler, A. (2003). Perceptual learning in speech. *Cognitive*
13 *Psychology*, 47(2), 204–238. [https://doi.org/10.1016/S0010-0285\(03\)00006-9](https://doi.org/10.1016/S0010-0285(03)00006-9)
- 14 Nusbaum and Morin. (1992). Paying attention to difference among talkers. In Y. Tohkura, E.
15 Vatikiotis-Bateson, & Y. Sagisaka (Eds.), *Speech Perception, Speech Production, and*
16 *Linguistic Structure* (pp. 113–134). IOS Press, Amsterdam.
- 17 Nusbaum, H., & Magnuson, J. S. (1997). Talker Normalization : Phonetic Constancy as a
18 Cognitive Process. In K. A. Johnson & J. W. Mullennix (Eds.), *Talker variability and*
19 *speech processing* (pp. 109–132). Academic Press. <https://doi.org/10.1121/1.2028337>
- 20 Oganian, Y., Bhaya-Grossman, I., Johnson, K., & Chang, E. F. (2023). Vowel and formant
21 representation in the human auditory speech cortex. *Neuron*, 111(13), 2105-2118.e4.
22 <https://doi.org/10.1016/j.neuron.2023.04.004>
- 23 Oldfield, R. C. (1971). The assessment and analysis of handedness: The Edinburgh inventory.
24 *Neuropsychologia*, 9(1), 97–113. [https://doi.org/10.1016/0028-3932\(71\)90067-4](https://doi.org/10.1016/0028-3932(71)90067-4)
- 25 Ou, J., & Law, S. P. (2017). Cognitive basis of individual differences in speech perception,
26 production and representations: The role of domain general attentional switching.
27 *Attention, Perception, and Psychophysics*, 79(3), 945–963.
28 <https://doi.org/10.3758/s13414-017-1283-z>
- 29 Peng, G. (2006). Temporal and tonal aspects of Chinese syllables: A corpus-based
30 comparative study of mandarin and cantonese. *Journal of Chinese Linguistics*, 34(1),
31 134–154.

- 1 Peng, G., Zhang, C., Zheng, H., Minett, J. W., & Wang, W. S.-Y. (2012). The effect of
2 intertalker variations on acoustic – perceptual mapping in Cantonese. *Journal of Speech,
3 Language, and Hearing Research, 55*, 579–596. [https://doi.org/10.1044/1092-
4 4388\(2011/11-0025\)language](https://doi.org/10.1044/1092-4388(2011/11-0025)language)
- 5 Peterson, G. E., & Barney, H. L. (1952). Control methods used in a study of the vowels. *The
6 Journal of the Acoustical Society of America, 24*(2), 175–184.
7 <https://doi.org/10.1121/1.1906875>
- 8 Schmidt, L. B. (2015). Not all forms of dialect contact are the same: Effects of regional
9 media, travel, and social contacts on the perception of Spanish aspirated-/s/. *Borealis –
10 An International Journal of Hispanic Linguistics, 4*(1), 99.
11 <https://doi.org/10.7557/1.4.1.3284>
- 12 Sjerps, M. J., Fox, N. P., Johnson, K., & Chang, E. F. (2019). Speaker-normalized sound
13 representations in the human auditory cortex. *Nature Communications, 10*:2465.
14 <https://doi.org/10.1038/s41467-019-10365-z>
- 15 Stilp, C. (2019). Acoustic context effects in speech perception. *Wiley Interdisciplinary
16 Reviews: Cognitive Science, e1517*. <https://doi.org/10.1002/wcs.1517>
- 17 Tamati, T. N., & Pisoni, D. B. (2014). Non-native Listeners' Recognition of High-Variability
18 Speech Using PRESTO. *Journal of the American Academy of Audiology, 25*(09), 869–
19 892. <https://doi.org/10.3766/jaaa.25.9.9>
- 20 Tao, R., Zhang, K., & Peng, G. (2021). Music Does Not Facilitate Lexical Tone
21 Normalization: A Speech-Specific Perceptual Process. *Frontiers in Psychology,
22 12*(October), 1–14. <https://doi.org/10.3389/fpsyg.2021.717110>
- 23 Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with S* (Fourth Edi).
24 <https://www.stats.ox.ac.uk/pub/MASS4/>
- 25 Watkins, A. J., & Makin, S. J. (1994). Perceptual compensation for speaker differences and
26 for spectral-envelope distortion. *The Journal of the Acoustical Society of America,
27 96*(3), 1263–1282. <https://doi.org/10.1121/1.410275>
- 28 Watkins, A. J., & Makin, S. J. (1996). Effects of spectral contrast on perceptual
29 compensation for spectral-envelope distortion. *The Journal of the Acoustical Society of
30 America, 99*(6), 3749–3757. <https://doi.org/10.1121/1.414981>
- 31 Winkler, I., Kujala, T., Tiitinen, H., Sivonen, P., Alku, P., Lehtokoski, A., Czigler, I., Csépe,

- 1 V., Ilmoniemi, R. J., & Näätänen, R. (1999). Brain responses reveal the learning of
2 foreign language phonemes. *Psychophysiology*, *36*(5), 638–642.
3 <https://doi.org/10.1017/S0048577299981908>
- 4 Wong, P. C. M., & Diehl, R. L. (2003). Perceptual normalization for inter- and intratalker
5 variation in cantonese level tones. *Journal of Speech, Language, and Hearing Research*,
6 *46*(2), 413–421. [https://doi.org/10.1044/1092-4388\(2003/034\)](https://doi.org/10.1044/1092-4388(2003/034))
- 7 Wong, P. C. M., Nusbaum, H. C., & Small, S. L. (2004). Neural bases of talker
8 normalization. *Journal of Cognitive Neuroscience*, *16*(7), 1173–1184.
9 <https://doi.org/10.1162/0898929041920522>
- 10 Xie, X., Jaeger, T. F., & Kurumada, C. (2023). What we do (not) know about the
11 mechanisms underlying adaptive speech perception: A computational framework and
12 review. *Cortex*, *166*, 377–424. <https://doi.org/10.1016/j.cortex.2023.05.003>
- 13 Xie, X., Theodore, R. M., & Myers, E. B. (2017). More than a boundary shift: Perceptual
14 adaptation to foreign-accented speech reshapes the internal structure of phonetic
15 categories. *Journal of Experimental Psychology: Human Perception and Performance*,
16 *43*(1), 206–217. <https://doi.org/10.1037/xhp0000285>
- 17 Zhang, C., Peng, G., & Wang, W. S. Y. (2013). Achieving constancy in spoken word
18 identification: Time course of talker normalization. *Brain and Language*, *126*(2), 193–
19 202. <https://doi.org/10.1016/j.bandl.2013.05.010>
- 20 Zhang, C., Peng, G., Wang, X., & Wang, W. S. (2015). Cumulative effects of phonetic
21 context on speech perception. *Proceedings of the 18th International Congress of*
22 *Phonetic Sciences*.
- 23 Zhang, K., & Peng, G. (2021). The time course of normalizing speech variability in vowels.
24 *Brain and Language*, *222*(July), 105028. <https://doi.org/10.1016/j.bandl.2021.105028>
- 25 Zhang, K., Peng, G., Li, Y., Minett, J. W., & Wang, W. S. Y. (2018). The effect of speech
26 variability on tonal language speakers' second language lexical tone learning. *Frontiers*
27 *in Psychology*, *9*(OCT), 1–13. <https://doi.org/10.3389/fpsyg.2018.01982>
- 28 Zhang, K., Sjerps, M. J., & Peng, G. (2021). Integral perception, but separate processing: The
29 perceptual normalization of lexical tones and vowels. *Neuropsychologia*, *156*, 107839.
30 <https://doi.org/10.1016/j.neuropsychologia.2021.107839>
- 31 Zhang, K., Tao, R., & Peng, G. (2023). The Advantage of the Music-Enabled Brain in

1 Accommodating Lexical Tone Variabilities. *Brain and Language*, 247(October),
2 105348. <https://doi.org/10.1016/j.bandl.2017.10.014>
3 Zhang, K., Wang, X., & Peng, G. (2017). Normalization of lexical tones and nonlinguistic
4 pitch contours: Implications for speech-specific processing mechanism. *The Journal of*
5 *the Acoustical Society of America*, 141(1), 38–49. <https://doi.org/10.1121/1.4973414>
6 Zheng, Y., & Samuel, A. G. (2020). The Relationship Between Phonemic Category
7 Boundary Changes and Perceptual Adjustments to Natural Accents. *Journal of*
8 *Experimental Psychology: Learning, Memory, and Cognition*, 46(7), 1270–1292.
9 <https://doi.org/10.1037/xlm0000788>

10