



The more complex the better? Mandarin tone perception by Cantonese and Hakka speakers

Siyi Lian¹, Min Liu¹

¹College of Chinese Language and Culture, Jinan University, Guangzhou, China
liansiyi94@gmail.com, minliu@jnu.edu.cn

Abstract

Much research has been conducted to investigate how native tonal language experience, in comparison with non-tonal language experience, shapes the perception of non-native tones. Little is known about how tonal experience of different dialects within a tonal language affects tone perception of the standard variety. The present study aimed to investigate how the complexity of tonal system in two Chinese varieties (i.e., Cantonese and Hakka), and how the tonal correspondences between each of the two varieties and Mandarin (i.e., the standard variety), affect Mandarin tone perception by speakers of the two varieties. The tonal system of Cantonese is more complex than that of Hakka, in terms of both tonal inventories and tonal categories. However, the correspondences between Cantonese tones and the Mandarin level/falling tones are looser than those between Hakka and Mandarin tones. An identification experiment and a discrimination experiment of Mandarin level and falling tones were conducted among native speakers of Guangzhou Cantonese, Meizhou Hakka (both with high level of Mandarin) and Mandarin. Results showed that the more complex the native tonal system, the better perception of Mandarin tones. Surprisingly, the tonal correspondences between the language variety and Mandarin did not seem to affect Mandarin tone perception as expected.

Index Terms: tone perception, Mandarin, tonal language experience, tonal complexity, tonal correspondences

1. Introduction

With 60%-70% of the world's languages being tonal languages [1], an important area of linguistic research concerns how native language experience affects the perception of non-native tones. Most of the research compared the non-native tone perception between listeners with a native tonal language and those with a native non-tonal language background. Mixed results have been found. Some studies reported that native tonal language experience facilitated the perception of non-native tones, with tonal listeners outperformed non-tonal listeners in the identification and/or discrimination of non-native tones [2-4]. Other studies, in contrast, found that native tonal language experience, interfered with the perception of non-native tones [5-7]. A third group of studies did not find significant differences between listeners with a native tonal language and a native non-tonal language. Instead, they revealed language-specific tonal difficulty patterns [8-10].

The above conflicting results could result from many factors, such as the different degrees of cross-linguistic similarities and differences between the target language and the listeners' native language [10], different tonal contrasts investigated [11], different levels of non-native language proficiency [12], and different prosodic functions pitch played

in the listeners' native language [13] across studies. This calls for cross-linguistic tone perception research in more closely related languages where these factors can be controlled.

Chinese provides an ideal test case for its rich varieties. For most Chinese, they grew up speaking a regional dialect. And they learned the standard variety (i.e., Mandarin) through media or education. It is therefore of great interest to investigate how tonal experience of different dialects within Chinese affects the Mandarin tone perception of these dialectal speakers.

Relatively few studies have been conducted to investigate the cross-linguistic tone perception by listeners from different native tonal languages, among which much attention has been paid to the effect of the size of native tonal inventory. Again, a denser native tonal inventory has been found to either facilitate [14] or inhibit [5, 6] the non-native tone perception. However, an often-neglected fact in these studies is that the effect of the size of native tonal inventory, more often than not, is intertwined with the tonal similarities or differences between the target language and the different native tonal languages under investigation. The latter, according to the Perceptual Assimilation Model [15,16] and the Speech Learning Model [17], could affect the non-native tone perception and the learnability of a non-native sound. The present study aimed to separate the two factors by investigating how the complexity of tonal system in two Chinese varieties, and how the tonal correspondences between each of the two varieties and Mandarin affect Mandarin tone perception.

The two Chinese varieties examined in this study are Guangzhou Cantonese and Meizhou Hakka. There are 9 tones (with 3 entering tones) in Guangzhou Cantonese. On the five-scale notation system [18], six of them are level tones (55, 33, 22, 5, 3, 2); two are rising tones (25, 23). Falling tones are rare (21, 53 as allotone of 55 in some cases) [19]. As to Meizhou Hakka, there are 6 tones (with 2 entering tones). Three are level tones (44, 11, 4/5). The rest are all falling tones (31, 52, 21) [20]. Apparently, the tonal system of Cantonese is more complex than that of Hakka, in terms of both tonal inventories (9 vs. 6) and tonal categories (level/rising/falling vs. level/falling). Moreover, the target language Mandarin has a high-level tone (T1: 55) and a high-falling tone (T4: 51) [18]. Cantonese speakers are equipped with high-level tone experience, but not high-falling tone experience in their native tonal system. Hakka speakers, contrastively, are familiar with both the high-level and high-falling tones. That is, the correspondences between Cantonese tones and the Mandarin level/falling tones are looser than those between Hakka and Mandarin tones. The specific research question arises as to how the complexity of tonal system in Cantonese and Hakka, and how the tonal correspondences between Cantonese and Mandarin, and between Hakka and Mandarin, affect Mandarin tone perception by Cantonese and Hakka speakers. An

identification experiment (Experiment 1) and a discrimination experiment (Experiment 2) of Mandarin level and falling tones were conducted among native speakers of Guangzhou Cantonese, Meizhou Hakka and Mandarin (as control group) to address the issue. Two hypotheses were made: 1) The tonal correspondences between each dialect and Mandarin plays a more important role than the complexity of the tonal system itself in the perception of Mandarin level-falling tonal contrast; 2) The better the tonal correspondences between the dialect and Mandarin, the better Mandarin tone perception by listeners who speak the dialect. This should be reflected in better perception of Mandarin level-falling tonal contrast by Hakka speakers than by Cantonese speakers. We expect that the Mandarin native speakers would have the best performance, as a result of native language effect.

2. Experiment 1: Identification

2.1. Method

2.1.1. Participants

Thirty undergraduate students at local universities were paid to participate in the experiment. Ten were Cantonese native speakers, born and raised in Guangzhou, Guangdong, China. Ten were Hakka native speakers from Meizhou, Guangdong, China. The other ten were Mandarin native speakers from northern China with no dialect experience. Note that the Cantonese speakers and the Hakka speakers also reported a high level of Mandarin spoken comprehension on a scale of 1 to 10 (8.5 ± 1.1 vs. 8.0 ± 0.9). None of them had received any formal musical training or had reported any speech or hearing disorders. Informed consent was obtained from all the participants before the experiment.

2.1.2. Stimuli

Thirty monosyllabic minimal tone sets with full sets of all four Mandarin tones (T1: high-level; T2: mid-rising; T3: low-dipping; T4: high-falling) were selected. Each monosyllabic item was a frequent monosyllabic word with more than 4,500 occurrences in a corpus of 193 million words [21]. And the frequency of the four monosyllabic items within each minimal set was comparable. An exemplar set was $ma1[ma^{55}]$, $ma2[ma^{35}]$, $ma3[ma^{214}]$, $ma4[ma^{51}]$. In total, 120 monosyllabic items (30 Syllables * 4 Tones) were selected.

These Mandarin monosyllables were recorded by a female speaker, who was born and grew up in Beijing and had no knowledge of any other dialects. The recording took place in Beijing at a soundproof recording booth. All the stimuli were recorded at 16-bit resolution with a sampling rate of 44.1 KHz. The amplitude of all the speech items was normalized for perception in Praat [22].

2.1.3. Procedure

Participants were tested individually in the behavioral lab at Jinan University in Guangzhou. All the monosyllabic items (30 Syllables * 4 Tones) were randomly presented to the participants using the E-Prime 3.0 software through headphones at a comfortable listening level. The T1(level) and T4 (falling) items were experimental stimuli, whereas the T2 and T3 items served as fillers to avoid response strategies.

The experiment included a practice block and two experimental blocks. The practice block contained 5 trials,

which were not used in the experimental blocks. Each experimental block contained 60 trials. Between two blocks there was a self-paced break. Each trial started with a 250 ms fixation cross, followed by a 250 ms blank screen. A monosyllabic speech item was then auditory presented along with a visual task interface. Participants were requested to identify the tone of the Mandarin speech sound they heard from the four Mandarin tone marks (“ˉ ˊ ˋ ˋˊ” standards for T1, T2, T3 and T4, respectively) as accurately and as quickly as possible. They were given up to 3 s to respond. The inter-stimulus interval was 500 ms. Instructions were given both visually on screen and orally by the experimenter in Mandarin before the experiment.

2.1.4. Data Analysis

Identification accuracy and Reaction time (RT) were analyzed and compared among different groups of participants between different tones. We were particularly interested in the confusion matrix between Mandarin T1 and T4 by the Cantonese and Hakka speakers relative to the Mandarin speakers, if there is any. RT was calculated from the offset of the monosyllabic items for correct responses, since the duration of each tone differs. To normalize the distribution, all the raw RTs were transformed using the nature logarithm.

Statistical analyses were carried out with the package *lme4* [23] in R [24]. Analysis of Identification accuracy (Correct/Incorrect) was performed using binomial logistic regression models, and analysis of RT was performed using linear mixed-effects regression models. All models included Tone (T1/T2/T3/T4) and Native Language (Cantonese/Hakka/ Mandarin) as fixed factors, and Subjects and Items as random factors. The fixed factors were added in a stepwise fashion and their effects on model fits were evaluated via model comparisons based on log-likelihood ratios. For models of RT, trials with absolute standardized deviations exceeding 2.5 from the mean were considered as outliers and removed from further analysis.

2.2. Results

2.2.1. Identification accuracy

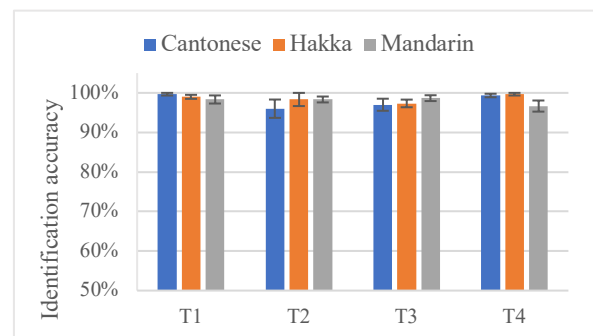


Figure 1: Identification accuracy of Mandarin tones by native speakers of Cantonese, Hakka and Mandarin. The error bars represents $\pm 1 SE$ of the means across participants.

Results (see Figure 1) showed a significant main effect of Tone ($\chi^2(3) = 158.85$, $p < 0.001$). However, neither a main effect of Native Language ($\chi^2(2) = 2.29$, $p = 0.32$) nor an interaction of Tone * Native Language ($\chi^2(6) = 4.54$, $p = 0.60$) were found. These suggest that the identification accuracy of Mandarin tones was not significantly different among Cantonese, Hakka

and Mandarin speakers across different tones. Post-hoc comparisons between tones showed that the identification accuracy of T1 was higher than that of T4 ($p < 0.001$). However, as identification accuracy was almost at ceiling level for each condition, the very few incorrect responses were likely motor-related errors due to the speed requirement of the task. More importantly, we did not find any significant confusion between Mandarin T1 and T4 for any group of participants.

2.2.2. Reaction time

Seventy-five trials (2.1%) were identified as outliers and removed from further analysis. Analysis of the remaining datapoints (see Figure 2) showed a significant main effect of Tone ($\chi^2(3) = 24.01, p < 0.001$), and a significant main effect of Native Language ($\chi^2(2) = 175.64, p < 0.001$). There was no interaction of Tone * Native Language ($\chi^2(6) = 7.01, p = 0.32$).

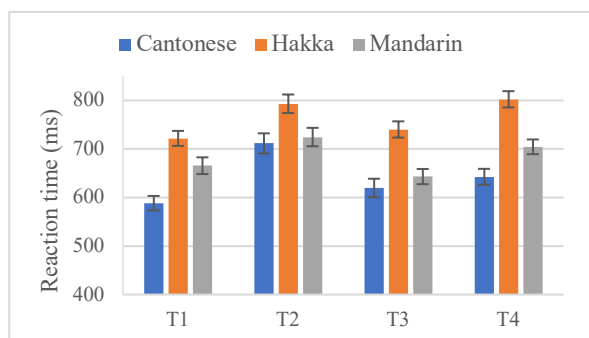


Figure 2: Average reaction time to identify Mandarin tones by native speakers of Cantonese, Hakka and Mandarin. The error bars represents $\pm 1 SE$ of the means across participants.

Post-hoc pairwise comparisons of reaction time between tones showed that Mandarin T1 was identified faster than T4 across participant groups ($\beta = -0.09, z = -3.57, p = 0.002$). And post-hoc pairwise comparisons of reaction time between participant groups showed that Cantonese speakers were significantly faster than Hakka ($\beta = -0.20, z = -13.28, p < 0.0001$) and Mandarin speakers ($\beta = -0.08, z = -4.88, p < 0.0001$) in Mandarin tone identification across tones. Hakka speakers, on the other hand, were notably slower than Mandarin speakers ($\beta = 0.13, z = 8.34, p < 0.0001$).

Overall, all three groups of participants were able to identify Mandarin tones with very high accuracy. However, in terms of reaction time, T1 seems to be easier to identify than T4 for all of them. Cantonese speakers were faster in Mandarin tone identification than Mandarin speakers, who were in turn faster than Hakka speakers. There seems to be an advantage for Cantonese speakers in Mandarin tone identification.

3. Experiment 2: Discrimination

3.1. Method

3.1.1. Participants

The same participants took part in Experiment 2. The order of the two experiments was counterbalanced across participants.

3.1.2. Stimuli

The same set of stimuli as in Experiment 1 was used. The experimental stimuli consisted of 30 monosyllabic level_falling

(T1_T4) minimal pairs in Mandarin. An AX discrimination task was adopted. Each T1_T4 pair was combined in four ways: AA, BB, AB and BA. The former two were the “same” combinations, whereas the latter two were the “different” combinations. In total, there were 120 (30 Root syllables * 4 Combinations) pairs of experimental stimuli for discrimination. The other half stimuli with the same setup for other minimal tone pairs (e.g., T2_T3) served as fillers. Recording and stimuli editing were the same as in Experiment 1.

3.1.3. Procedure

Procedure was similar as in Experiment 1, except that in Experiment 2 there were four experimental blocks. Each block contained 60 trials as well. A trial started with a 250 ms fixation cross, followed by a 250 ms blank screen. The first speech item was then played. After a 600 ms pause, the second speech item was played along with a visual task interface. Participants were requested to judge if the two tones of the speech items were the same or different as accurately and as quickly as possible. They had to respond in 3 s, or else the program moved on to the next trial automatically. The inter-stimulus interval was 500 ms.

3.1.4. Data Analysis

Discrimination accuracy and Reaction time (RT) were analyzed and compared among different groups of participants for the T1_T4 pair. RT was calculated from the onset of the second speech item. To normalize the distribution, all the raw RTs were transformed using the nature logarithm.

Similar statistical analyses as Experiment 1 were conducted for Discrimination accuracy and RT in Experiment 2, except that this time, all models included Combination Type (Same/Different) and Native Language (Cantonese/Hakka/ Mandarin) as fixed factors, and Subjects and Items as random factors.

3.2. Results

3.2.1. Discrimination accuracy

Results (see Figure 3) showed no main effect of Native Language ($\chi^2(2) = 1.64, p = 0.44$) or Combination Type ($\chi^2(1) = 2.54, p = 0.11$). There was, however, an interaction of Native Language * Combination Type ($\chi^2(2) = 6.04, p = 0.048$).

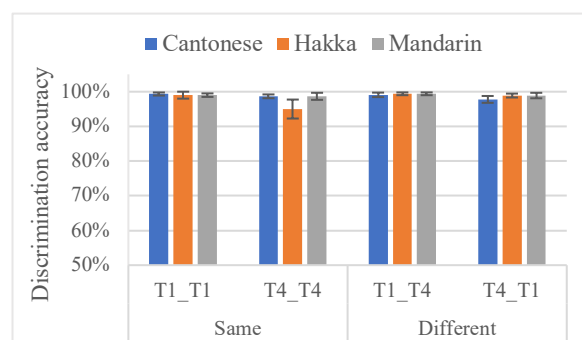


Figure 3: Discrimination accuracy of the Mandarin T1_T4 pair by native speakers of Cantonese, Hakka and Mandarin. The error bars represents $\pm 1 SE$ of the means across participants.

Separate models were constructed for subset data of the same tone combinations and different tone combinations. We found a significant main effect of Tone Combination for the same tone combinations ($\chi^2(1) = 6.67, p = 0.01$), but not for the

different tone combinations ($\chi^2(1) = 2.57, p = 0.11$). That is, the discrimination accuracy for the same T1 tones were higher than that for the same T4 tones. Nevertheless, as can be seen from Figure 3, the discrimination of Mandarin T1 and T4 almost reached ceiling level across all experimental conditions. All three groups of participants could discriminate between Mandarin T1 and T4 with very high accuracy.

3.2.2. Reaction time

Sixty-one trials (1.7%) were identified as outliers and removed from further analysis. Analysis of the remaining datapoints (see Figure 4) showed a significant main effect of Native Language ($\chi^2(2) = 47.91, p < 0.001$). No main effect of Combination Type ($\chi^2(1) = 0.63, p = 0.43$) or an interaction effect of Native Language * Combination Type ($\chi^2(2) = 1.60, p = 0.45$) was found. These suggest that the reaction time to discriminate between the same tones (T1_T1 & T4_T4) did not differ from that to discriminate between different tones (T1_T4 & T4_T1), regardless of the native language background of the participants.

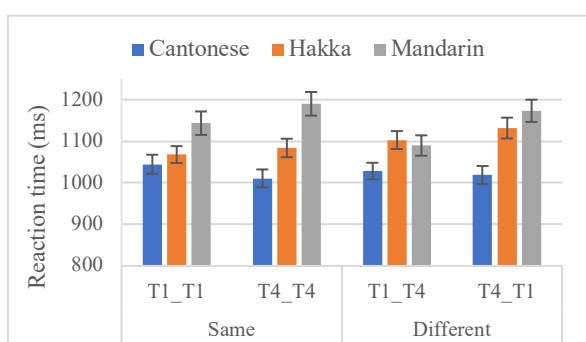


Figure 4: Average reaction time to discriminate between the Mandarin T1 and T4 by native speakers of Cantonese, Hakka and Mandarin. The error bars represents $\pm 1 SE$ of the means across participants.

Post-hoc pairwise comparisons of reaction time between participant groups revealed that Cantonese speakers took significantly shorter time than Hakka ($\beta = -0.07, z = -5.12, p < 0.0001$) and Mandarin speakers ($\beta = -0.09, z = -6.62, p < 0.0001$) to distinguish Mandarin T1 and T4. Hakka speakers, however, were not significantly different from Mandarin speakers as to the reaction time needed for Mandarin T1_T4 discrimination ($\beta = -0.02, z = -1.49, p = 0.30$).

Overall, all three groups of participants could distinguish Mandarin T1 and T4 with very high accuracy. Despite that, there were indeed differences in reaction time among the three groups. Cantonese speakers were significantly faster than Hakka and Mandarin speakers to discriminate between Mandarin T1 and T4. And Hakka speakers were at comparable speed with Mandarin speakers. There, again, appears to be an advantage for Cantonese speakers in the discrimination of Mandarin T1 and T4.

4. Discussion

The present study investigated how the complexity of tonal system in two Chinese varieties (i.e., Cantonese and Hakka), and how the tonal correspondences between each of the two varieties and Mandarin (i.e., the standard variety), affect Mandarin tone perception by speakers of the two varieties. An identification experiment and a discrimination experiment of

Mandarin level and falling tones were conducted among native speakers of Guangzhou Cantonese, Meizhou Hakka and Mandarin. Both the identification accuracy and discrimination accuracy almost reached ceiling level across different groups of participants, suggesting that Guangzhou Cantonese and Meizhou Hakka speakers could correctly identify and distinguish the Mandarin level and falling tones eventually. However, there were noticeable reaction time differences among the three groups of participants. Cantonese speakers were significantly faster than Hakka and Mandarin speakers to complete both tasks. As reaction time is often taken as a reliable indicator of the degree of difficulty of a perceptual decision [25], the reaction time results here suggest that there seems to be a processing advantage for Cantonese speakers to perceive the Mandarin level-falling tonal contrast. Overall, it appears that the complexity of the native tonal system plays a more important role than the tonal correspondences between the dialect and Mandarin in the perception of Mandarin tonal contrast for the dialectal speakers. A more complex native tonal system tends to facilitate the non-native tone perception.

These results, one might have noticed, are not in agreement with either of our hypotheses. We hypothesized that the tonal correspondences between each dialect and Mandarin should outweigh the complexity of the tonal system in the perception of Mandarin level-falling tonal contrast. If it were true, with better correspondences between Hakka tones and Mandarin level/falling tones than those between Cantonese and Mandarin tones, better perception of Mandarin level-falling tonal contrast would have been found for Hakka speakers than for Cantonese speakers. And yet we found opposing patterns. Cantonese speakers showed a processing advantage than Hakka speakers and even Mandarin speakers, as evidenced by the reaction time measure. This suggests that the more complex the native tonal system, the better capabilities of non-native tone perception [14]. However, the high accuracy of Cantonese speakers in distinguishing Mandarin level and falling tones contradicts with the discrimination difficulty reported in [9] and [10]. One plausible reason could be that the dialectal speakers we tested had native-like proficiency in Mandarin, while those Cantonese participants tested in [9] and [10] were either intermediate-level Mandarin learners or naïve to Mandarin. They were reported to assimilate the Mandarin level and falling tones to the same Cantonese tone, causing perceptual confusion of the two Mandarin tones [9]. It seems that our Cantonese participants have reshaped the perceptual space for Mandarin level and falling tones with growing Mandarin experience. A direct comparison research between the two groups is needed. All in all, the mechanism of cross-dialect tone perception can be different from that of the cross-language tone perception.

5. Conclusions

To conclude, the complexity of the native tonal system plays a more important role than the tonal correspondences between the dialect and Mandarin in the perception of Mandarin tonal contrast for the dialectal speakers. A more complex native tonal system tends to facilitate the non-native tone perception.

6. Acknowledgements

This research was supported by the Guangdong Planning Office of Philosophy and Social Science Grant (GD23CZY01) and the Fundamental Research Funds for the Central Universities (23JNQN15, 23JNLH10).

7. References

- [1] M. Yip, *Tone*. Cambridge: Cambridge University Press, 2002.
- [2] R. P. Wayland, and S. G. Guion, "Training English and Chinese Listeners to Perceive Thai Tones: A Preliminary Report," *Language Learning*, vol. 54, no. 4, pp. 681-712, 2004.
- [3] Y.-h. S. Chang, Y. Yao, and B. H. Huang, "Effects of linguistic experience on the perception of high-variability non-native tones," *The Journal of the Acoustical Society of America*, vol. 141, no. 2, pp. EL120-EL126, 2017.
- [4] L. Liu, R. Lai, L. Singh, M. Kalashnikova, P. C. M. Wong, B. Kasisopa, A. Chen, C. Onsuwan, and D. Burnham, "The tone atlas of perceptual discriminability and perceptual distance: Four tone languages and five language groups," *Brain and Language*, vol. 229, pp. 105106, 2022.
- [5] X. Wang, "Perception of Mandarin Tones: The Effect of L1 Background and Training," *The Modern Language Journal*, vol. 97, no. 1, pp. 144-160, 2013.
- [6] M. A.-O. X. Zhu, F. A.-O. Chen, X. Chen, and Y. Yang, "The more the better? Effects of L1 tonal density and typology on the perception of non-native tones," no. 1932-6203 (Electronic).
- [7] K. Tsukada, and M. Kondo, "The Perception of Mandarin Lexical Tones by Native Speakers of Burmese," *Language and Speech*, vol. 62, no. 4, pp. 625-640, 2018.
- [8] A. L. Francis, V. Ciocca, L. Ma, and K. Fenn, "Perceptual learning of Cantonese lexical tones by tone and non-tone language speakers," *Journal of Phonetics*, vol. 36, no. 2, pp. 268-294, 2008.
- [9] Y.-C. Hao, "Second language acquisition of Mandarin Chinese tones by tonal and non-tonal language speakers," *Journal of Phonetics*, vol. 40, no. 2, pp. 269-279, 2012.
- [10] C. K. So, and C. T. Best, "Cross-language Perception of Non-native Tonal Contrasts: Effects of Native Phonological and Phonetic Influences," *Language and Speech*, vol. 53, no. 2, pp. 273-293, 2010.
- [11] J. T. Gandour, "Tone perception in Far Eastern languages," *Journal of Phonetics*, vol. 11, pp. 149-175, 1983.
- [12] E. Pelzl, J. Liu, and C. Qi, "Native language experience with tones influences both phonetic and lexical processes when acquiring a second tonal language," *Journal of Phonetics*, vol. 95, pp. 101197, 2022.
- [13] V. Schaefer, and I. Darcy, "Lexical function of pitch in the first language shapes cross-linguistic perception of Thai tones," vol. 5, no. 4, pp. 489-522, 2014.
- [14] Y.-S. Lee, D. A. Vakoch, and L. H. Wurm, "Tone perception in Cantonese and Mandarin: A cross-linguistic comparison," *Journal of Psycholinguistic Research*, vol. 25, no. 5, pp. 527-542, 1996.
- [15] C. T. Best, "A Direct Realist View of Cross-language Speech Perception," *Speech perception and linguistic experience: Issues in cross-language research*, W. Strange, ed., pp. 171-204, Baltimore, MD: York Press, 1995.
- [16] C. T. Best, and M. D. Tyler, "Nonnative and second-language speech perception: Commonalities and complementarities," *Language experience in second language speech learning*, M. J. Munro and O.-S. Bohn, eds., pp. 13-34, Amsterdam: John Benjamins Publishing, 2007.
- [17] J. E. Flege, "Second language speech learning: Theory, findings, and problems," *Speech perception and linguistic experience: Issues in cross-language research*, W. Strange, ed., Timonium, MD: York press, 1995.
- [18] Y. R. Chao, "A system of tone letters," *Le Maître Phonétique*, vol. 45, pp. 24-27, 1930.
- [19] Y.-Y. Fok-Chan, *A Perceptual Study of Tones in Cantonese* p.^pp. 191, Hong Kong: University of Hong Kong Press, 1974.
- [20] X. Li, and M. Xiang, *An introduction to Chinese Dialectology*, Beijing: Peking University Press, 2015.
- [21] J. Da, "A corpus-based study of character and bigram frequencies in Chinese e-texts and its implications for Chinese language instruction," *Proceedings of 4th International Conference on New Technologies in Teaching and Learning Chinese*, P. Zhang, T. Xie and J. Xu, eds., pp. 501-511, Beijing: The Tsinghua University Press, 2004.
- [22] P. Boersma, and D. Weenink, "Praat: doing phonetics by computer [Computer program]. Version 6.4.01," 2023.
- [23] D. Bates, M. Mächler, B. M. Bolker, and S. C. Walker, "Fitting linear mixed-effects models using lme4," *Journal of Statistical Software*, vol. 67, no. 1, pp. 1-48, 2015.
- [24] R Core Team, "R: A Language and Environment for Statistical Computing," R Foundation for Statistical Computing, 2023.
- [25] K. Schneider, G. Dogil, and B. Möbius, "Reaction time and decision difficulty in the perception of intonation." pp. 2221-2224.