

1 **Perceptual clustering of high-pitched vowels in Chinese Yue Opera**

2 Yixin Zhang,<sup>1</sup> Francis Nolan,<sup>1</sup> and Daniel Friedrichs<sup>1,2,3</sup>

3 <sup>1</sup> *Phonetics Laboratory, Faculty of Modern and Medieval Languages and Linguistics, University of Cambridge,*  
4 *Sidgwick Avenue, Cambridge CB3 9DA, United Kingdom*

5 <sup>2</sup> *Department of Speech, Hearing and Phonetic Sciences, UCL, 2 Wakefield Street, London WC1N 1PF, United*  
6 *Kingdom*

7 <sup>3</sup> *Department of Computational Linguistics, University of Zurich, Andreasstrasse 15, 8050 Zurich, Switzerland*

8 **Abstract**

9 Numerous studies on Western Opera singing have shown that listeners' vowel identification  
10 performance decreases with an increasing fundamental frequency ( $f_0$ ). This study explores the  
11 intelligibility of high-pitched vowels in Yue Opera, the largest dialectal opera in China. Six long  
12 vowels (/i y e a o u/) were recorded by a professional female singer at ten  $f_0$ s between 220 and  
13 932 Hz, of which 700-ms nuclei with flat  $f_0$  contours and resonance trajectories were extracted as  
14 stimuli. In a within-subject design, sixteen phonetically trained listeners responded on a free-  
15 choice vowel quadrilateral (task 1) and in a two-alternative forced-choice task (task 2) to indicate  
16 which vowel was presented. Results show that vowels cluster in the perceptual space into three  
17 groups (/i y e/, /u o/, /a/) above 521 Hz and that listeners could identify vowels between but not  
18 within groups with high accuracy up to at least 932 Hz. Multidimensional scaling (MDS) of  
19 simulated auditory excitation patterns reveals highly differentiable spectral shapes between  
20 groups. These findings put into question whether previous results on Western Opera could be  
21 generalized to other forms of opera singing.

22 **Key Words:** Vowel Intelligibility, High-pitched Singing, Chinese Yue Opera

## 23 I. INTRODUCTION

24 It is often assumed that high-pitched singing is difficult to understand due to the loss of  
25 vowel intelligibility. The exploration of the loss has a rich history in Western Opera singing, and  
26 some literature has well summarized the relevant studies (e.g., Sundberg, 2013). As early as in  
27 1885, von Helmholtz described an observation that the timbre of the vowel /u/ shifted towards /o/  
28 when  $f_0$  in a male voice exceeded roughly 175 Hz (i.e., the musical note F3). A relatively recent  
29 study by Hollien et al. (2000) has confirmed this and shown that the identification of the sung  
30 vowels /i/ and /u/ shifted towards categories with  $F_1$  just above the  $f_0$  in the sung stimulus vowels.  
31 Similar observations were reported in several other studies (for an overview, see Sundberg, 2013,  
32 p. 87), which all indicate that above a certain absolute  $f_0$  of approximately 523 Hz (i.e., the  
33 musical note C5) listeners' identification performance for all vowels but /a/ and /ɑ/ (which have  
34 the highest  $F_1$  in normal speech) would successively decrease towards chance-level.

35 It is widely assumed among researchers from the field of singing that the aforementioned  
36 reduction in vocalic intelligibility is due to the sparse sampling of the vocal tract transfer  
37 function at high fundamental frequencies ( $f_0$ ), which leads to a poor specification of the formants.  
38 As a soprano's vocal range reaches musical notes corresponding to  $f_0$  around 1 kHz (e.g.,  
39 soprano C = 1046 Hz), the wide spacing of the harmonics makes it unlikely that typical formant  
40 frequency patterns can be found in the acoustic signal. This is particularly true for close vowels  
41 such as /i/ and /u/ that usually exhibit relatively low first formants ( $F_1$ ), which would be  
42 exceeded by such high  $f_0$ .

43 However, studies outside Western Opera singing reported a non-uniform relationship  
44 between vowel intelligibility and  $f_0$ . For example, Smith and Scott (1980) found that the vowels  
45 /i i e æ/ were identifiable (70% correct) up to an  $f_0$  of 880 Hz when they were produced in

46 isolation by a soprano in a non-singing style with a raised larynx (i.e., shortened vocal tract).  
47 When asked to produce the same vowels in her typical soprano singing style at the corresponding  
48 musical note A5, the identification score dropped to 4%. The results of a relatively recent study  
49 by Nolan and Sykes (2015) have put these findings into question as it was shown that the vowels  
50 /i ε a ɔ u ə/ (produced in CV context with an initial lateral) all were perceived as or close to /a/  
51 at an  $f_0$  of 880 Hz (A5), although the soprano was asked to produce them in a non-singing style.  
52 On the contrary, Maurer and Landis (1996) demonstrated that the isolated vowels /i a u o/ (but  
53 not /e/) could be identified accurately by listeners between 497 and 873 Hz when they were  
54 produced by untrained children, women, and men at individually chosen  $f_0$ .

55 The contradicting results of these studies may be due to the uncontrolled secondary cues to  
56 vowel category perception (e.g., vowel duration, formant frequency movements, and co-  
57 articulation in the consonantal environment; for more information on this, see, for example,  
58 Strange et al., 1976, Lehiste and Peterson, 1961). Nevertheless, some studies that used excised  
59 vowels with a single duration and quasi-flat  $f_0$  contours and resonance trajectories still reported  
60 satisfying identification performance of the participants.

61 For example, Friedrichs et al. (2015a) found that the phonological function of the isolated  
62 steady-state vowels /i y e ø ε a u o/ can be maintained at  $f_0$ s up to 880 Hz when they were tested  
63 in a listening test with only two response options. In a follow-up study investigating the  
64 influence of talker variability, Friedrichs et al. (2017) found that the cardinal vowels (point  
65 vowels) /i a u/ remained identifiable even up to 1046 Hz when they were tested in isolation and  
66 multiple response options were provided. In the same experiment, it was shown that listeners'  
67 identification performance decreased significantly for /y ε/ and dropped to chance for /e ø o/  
68 within the range of 523–1046 Hz. Based on the analyses of auditory excitation pattern

69 simulations, the authors proposed that the overall spectral shape of the cardinal vowels /i a u/  
70 may be utilized by listeners as acoustic landmarks that aid vowel perception at high  $f_0$ . This  
71 assumption is supported by several studies that indicated that gross spectral shapes as  
72 represented by, for example, Mel Frequency Cepstral Coefficients (MFCCs) (Davis and  
73 Mermelstein, 1980) carry superior acoustic cues to vowel category identification than formants  
74 (e.g., Ito et al., 2001, and Zahorian and Jagharghi, 1993).

75 As many studies suggest that vowel identification is possible even when no typical formant  
76 patterns can be found in the acoustic signal (for a comprehensive overview, see Maurer 2016), it  
77 seems plausible that the reduction in vocalic intelligibility, especially the bias towards open  
78 vowels in high-pitched Western Opera, may largely be due to its special singing style. Joliveau et  
79 al. (2004) have demonstrated that Western Opera singers shift their first resonance (by opening  
80 their jaws and lips), and hence  $F_1$ , to the vicinity of  $f_0$  when they are singing at high pitches to  
81 gain vocal power. This so-called resonance or formant tuning may be beneficial when  
82 performing in large auditoria without microphones. However, such adjustments made to the  
83 articulation inevitably lead to changes in the acoustic patterns, which may explain the previously  
84 described migration of vowel category perception to those with higher  $F_1$ . Therefore, it is not  
85 only the listeners who are ‘mishearing’ but also the singers who are ‘mispronouncing’ the vowels  
86 in Western Opera.

87 Western Opera are not the only musical drama that became popular before the wide use of  
88 microphones, which may have played an important role in the evolution of contemporary singing  
89 styles. Various styles of musical drama also exist in China, in which the characters are played by  
90 specially trained singers. Similar to their Western counterparts, these Chinese Opera singers need  
91 to sing loudly while achieving a certain aesthetical norm. More importantly, many Chinese

92 people also find Chinese Operas hard to understand. The naïve audience often attributes the  
93 unintelligibility of Chinese Operas to their special music styles, slow rhythm, and stylised  
94 languages, but whether changes in vowel qualities with changing  $f_0$  contribute to this has rarely  
95 been explored.

96 The only study on vowel intelligibility in Chinese Operas was carried out by Maurer et al.  
97 (2014), who examined the identifiability of vowels in Cantonese Opera singing by phonetically  
98 trained Cantonese speakers. The results showed high identification scores (>80% correct  
99 responses) for the vowels /i a ɔ u/ (but not /y æ/) in consonant-vowel (CV) or consonant-vowel-  
100 consonant (CVC) context. It is worth mentioning that the stimuli used in this study were  
101 extracted from a DVD of a famous female Cantonese Opera singer. Therefore, the vowels were  
102 not separated from the melody, and they could only roughly control the  $f_0$  levels. Some of these  
103 vowels might have co-occurred with musical notes that reflect the lyrics' lexical tones (for more  
104 information on tone and melody, see Wee, 2007 and You, 2006). In this way, the melody  
105 associated with the nine tones in Cantonese may narrow the lexical set and contribute to the  
106 identification by the native speakers. Thus, whether and why there is also a decrease in vocalic  
107 intelligibility as  $f_0$  increases in Chinese Operas requires empirical investigation under stricter  
108 conditions, namely, using vowels produced in isolation at strictly controlled  $f_0$ s. The influence of  
109 the melody, tone as well as some secondary cues should be carefully controlled.

110 Another Chinese Opera style, which has not been studied in this context yet is Yue Opera.  
111 Unlike Western Opera and other Chinese Operas, it has the unique feature that all its characters,  
112 including all gender and ages, are played by females. This obviously requires a vast amount of  
113 control over phonation and articulation. The language used in Yue Opera is a stylized language  
114 specific to the use on stage. This language is based on the Wu dialect spoken in Shengzhou but

115 also influenced by Mandarin during the development of Yue Opera (Qiu, 1995). Wu and  
116 Mandarin (and many other Chinese dialects) share the same logographic writing systems and  
117 similar syllabic structures. In both dialects, each character corresponds to a single morpheme and  
118 a syllable in the form of C<sup>G</sup>VC (C: consonant, G: glide, V: vowel or diphthong), and the onset  
119 and coda consonants are optional; namely, open syllables with a vowel or a diphthong only are  
120 allowed. However, the same character usually has different pronunciations in Wu and Mandarin  
121 dialects, and these two dialects have different vocalic, consonantal, and tonal inventories.  
122 According to previous studies (You, 2006; Huang, 2000; Qiu, 1999), the stage language of Yue  
123 Opera includes 13 groups of rhymes (i.e. *Dazhe*, ‘big rhymes’). The rhymes within the same  
124 group are considered as rhyming with each other, though they do not necessarily contain the  
125 same segments (e.g. /a/, /ia/, /ua/, /aʔ/, /iaʔ/, /uaʔ/ all belong to the same *Dazhe*). These 13 groups  
126 involve 20 long, short, nasalized, or dentalised vowels, with eight of them being long vowels (i.e.,  
127 /i/ /y/ /e/ /u/ /o/ /ɔ/ /a/ /ɤ/). It is noteworthy that the male and female characters tend to realize  
128 some vowels slightly different from each other in Yue Opera, and the mid-close vowel is more  
129 commonly realized as /e/ by female characters but /ɛ/ when a male character is played (personal  
130 communication with Shuyang Sheng, the invited singer, and Weitao Mao, a well-known male  
131 character player and the vice-chancellor of the China Theatre Association).

132         In the present study, we recorded a professional female Yue Opera singer producing  
133 seven long isolated vowels (/i y e a ɔ o u/) in her singing style at ten *f*<sub>0</sub>s between 220 Hz and 932  
134 Hz by presenting her the corresponding morphemic characters containing the vowels only. The  
135 vowel /ɤ/ was not recorded because no morpheme corresponds to an open syllable with /ɤ/,  
136 namely, the very few morphemic characters containing /ɤ/ all have onset consonants that  
137 interfere with vowel quality. We conducted listening tests to compare the results with those from

138 previous studies on Western Opera singing. To investigate the spectral properties underlying the  
139 listeners' identification process at high pitches, multidimensional scaling (MDS) was employed  
140 to geometrically model the changes in the perceptual space and simple versions of excitations  
141 patterns were analyzed that the vowels would be expected to generate in the auditory periphery.

## 142 **II. METHODS**

### 143 **A. Participants**

144 Fifteen phonetically trained listeners participated in the perceptual experiments (7 females, 8  
145 males; mean age = 24.6, standard deviation = 3.5). They were all students at the University of  
146 Cambridge, and none of them reported any hearing impairments when asked before the  
147 experiment.

### 148 **B. Stimuli and Apparatus**

149 A professional female Yue Opera singer (age = 35) who received special training since  
150 school age was recorded in a noise-controlled room at the Phonetics Laboratory of the University  
151 of Cambridge using a MixPre-6 recorder and a Sennheiser M64 microphone with a K6 battery  
152 module. The sampling frequency of the recordings was 44100 Hz. She was asked to produce the  
153 vowels /i y u e o a ə/ in the Yue Opera style at ten  $f_0$ s corresponding to musical notes between A4  
154 to B<sup>b</sup>6 without lexical tones (i.e., 220, 350, 440, 521, 659, 740, 784, 831, 880, and 932 Hz).  
155 Piano notes were presented as reference sounds to the singer via Sony MDR-Z7M2 headphones  
156 before each vowel production. She was asked to produce long and monotone vowels as  
157 accurately as possible while keeping a constant distance from the microphone of approximately  
158 30 cm. The recordings were done twice to elicit more accurate stimuli, once by vowel (i.e.,  
159 recording each vowel at all  $f_0$ s before moving on to the next vowel) and once by  $f_0$  (i.e., recording  
160 all vowels at one  $f_0$  before moving on to the next  $f_0$ ). The vowel recordings with the most

161 accurate  $f_0$  realization were selected as stimuli. For reference purposes, spoken versions of the  
162 vowels in Yue Opera style were recorded at an  $f_0$  that was comfortable for the singer (mean  $f_0$  =  
163 376.8 Hz, standard error = 25.93 Hz).

164 As the singer was unfamiliar with the international phonetic alphabet, she was presented  
165 with logographic characters corresponding to open syllables containing the target vowels. For  
166 each vowel, three different characters were presented to the singer to ensure the correct  
167 elicitation of each vowel. The characters used during the recordings were taken from Huang  
168 (2000), and the singer confirmed that the three characters in each group share the same vowel.

169 After the recording session, it was found that the singer diphthongized the vowel /ɔ/ into  
170 /ou/ throughout almost all  $f_0$ s. This change may be due to the influence of Mandarin, the common  
171 language the singer used in conversational speech, in which the /ɔ/-carrying syllables are realized  
172 as /ou/. This diphthongization makes it impossible to investigate the categorical perception of /ɔ/  
173 as a single vowel in the two perceptual tasks, so that /ɔ/ was dropped from the subsequent  
174 experiment and analysis. Only the recordings of the six long vowels /i y u e o a/ were used.

175 For each stimulus, 700-ms sound segments were extracted from the vowel centers. The  
176 excised sounds showed relatively flat  $f_0$  contours with a maximum deviation from the target  $f_0$  of  
177 4 %. The sounds were normalized in Praat (Borsma & Weenink, 2021) to 75 dB SPL, and the  
178 onsets and offsets of the sounds were faded over 5ms by amplitude modulating the waveform  
179 with raised cosines. During the experiment, the output level was adjusted by listeners  
180 individually to a comfortable listening level.

### 181 C. Procedure

182 The perceptual experiment involved a guided transcription task and a two-alternative  
183 forced-choice task conducted successively through E-prime 2.0 (Psychology Software Tools,



184 Pittsburgh, PA). The guided transcription task was chosen to investigate possible gradual  
185 changes in the vocalic intelligibility at different  $f_0$ s, while the subsequent two-alternative forced-  
186 choice task allowed a more refined exploration of the categorical perception of different vowels.  
187 The participants could take a break as long as they wanted between the two tasks.

188 In the guided transcription task, the six vowels were presented at ten  $f_0$ s in a pseudo-  
189 randomized order, resulting in 60 trials (6 vowels  $\times$  10  $f_0$ s). In each trial, the participants were  
190 presented with a figure representing the perceptual vowel space (including reference vowels, see  
191 FIG. 1) after receiving a vowel as an auditory stimulus. The perceptual space was presented as  
192 the vowel quadrilateral, in which the position of the vowels reflected a two-formant space. For  
193 instance, front rounded vowels were shown retracted from fully front.

194 The participants were asked to click at any point on the figure to indicate where they  
195 thought the vowel in the stimulus belonged to in the perceptual space. After the click, the screen  
196 would refresh automatically, signaling the start of the new trial, and the participants would hear  
197 the next stimulus simultaneously. The coordinates of their clicks were recorded. There was no  
198 time limit.

199 In the two-alternative forced-choice experiment, 300 trials were involved (6 intended  
200 vowels  $\times$  10  $f_0$ s  $\times$  5 noise vowels). In each trial, the participants were presented with an auditory  
201 stimulus and saw a screen that contained two horizontally arranged vowels out of the six, one of  
202 the two being the vowel intended by the singer. The left-right order of the vowel pairs, as well as  
203 the order of the auditory stimuli, was pseudo-randomized. The participants were asked to  
204 indicate whether it was the vowel on the right or the left they had heard by pressing two keys on  
205 the computer keyboard that were labeled beforehand by the investigator as 'right' or 'left'. After

206 the participants made their choice, they would hear the next stimulus automatically. There was  
207 no time limit, and the participants could only listen to a stimulus once.

#### 208 **D. Perceptual Data Analysis**

209 To analyze the results of the guided transcription, we indexed the change in vowel quality  
210 by the distance between the coordinates of the participants' clicks and the coordinates of the  
211 intended reference vowels on the diagram of perceptual space (henceforth *Perceptual Distance*).  
212 Here, *Perceptual Distance* is not used to index whether the participants made a correct or  
213 incorrect response, but the perceptual changes, which might also reflect the potential changes in  
214 the singer's articulatory strategy.

215 We constructed several linear mixed effects (LME) models in R (R core team, 2020)  
216 using lmer in lme4 (Bates et al., 2015). We selected the optimal fixed structure by using stepwise  
217 comparisons from the most complex effect to the simplest and the random effects by the smallest  
218 Akaike Information Criterion (AIC). The final model has *Perceptual Distance* as the dependent  
219 variable,  $f_0$  and *Intended Vowel* as the fixed effects and *Participant* as the random effect.

220 Following Friedrichs' design (2015a), the participants' responses in the two-alternative  
221 forced-choice task were analyzed with the bias-free non-parametric sensitivity measure  $A'$   
222 according to Signal Detection Theory (Tanner and Swets, 1954; Stanislaw and Todorov, 1999;  
223 Pallier, 2002) in R (R Core Team, 2020). Signal Detection Theory applies to the situation in  
224 which participants are asked to determine which one of the two categories (i.e., which one of a  
225 vowel pair in our case) a stimulus belongs to. The task generates two measures of behavioral  
226 performance: the hit rate and the false alarm rate. In the present study, the response option of the  
227 lower  $F_1$  (i.e., the closer vowel) was arbitrarily assigned to the signal (signal vowel), the other to  
228 the noise (noise vowel). Then, a *hit* (H) referred to when "the signal vowel was presented and

229 chosen”, a *miss* to when “the signal vowel was presented but not chosen”, a *false alarm* (F) to  
 230 when “the noise vowel was presented but not chosen” and a *correct rejection* to when “the noise  
 231 vowel was presented and chosen”. In studies using Signal Detection Theory, H and F are  
 232 transformed into indices of sensitivity and bias based on statistical models like A' and d' (Pollack  
 233 and Norman, 1964; Smith, 1995; Zhang and Mueller, 2005). Here, A' rather than d' was used  
 234 because it is a non-parametric measure that can deal with situations when *hit* or *false alarm* rates  
 235 are 0 or 1. In such instances, d', the z-score difference between the signal and noise distribution  
 236 (=Z(H) - Z(F)), is either -infinite or +infinite (Zhang and Mueller, 2005). A' was calculated  
 237 using the following formula (1) (Zhang and Mueller, 2005: 207):

$$238 \quad (1) A' = \begin{cases} 0.75 + \frac{H-F}{4} - F(1-H) & \text{when } F \leq 0.5 \leq H \\ 0.75 + \frac{H-F}{4} - \frac{F}{4H} & \text{when } F \leq H < 0.5 \\ 0.75 + \frac{H-F}{4} - \frac{1-H}{4(1-F)} & \text{when } 0.5 < F \leq H \end{cases}$$

239 A' ranges between 0 and 1, 1 indicating maximum performance and 0.5 indicating chance  
 240 performance. The participants' response bias was indexed by B''<sub>D</sub>, which correlates to the slope  
 241 of the receiver operating characteristic function at the point of observation. B''<sub>D</sub> is calculated as  
 242 described in formula (2) (Pallier, 2002) and ranges from -1 (maximum bias to the noise vowel)  
 243 and 1 (maximum bias to the signal vowel).

$$244 \quad (2) B''_D = \frac{(1-H) \times (1-F) - H \times F}{(1-H) \times (1-F) + H \times F}$$

245 Since a very high A' leads to meaningless B''<sub>D</sub> as it is based on a small number of misses  
 246 and false alarms (Stanislaw and Todorov, 1999; Zhang and Mueller, 2005), we only calculated  
 247 B''<sub>D</sub> of the vowel pairs with A' values smaller than 0.7. We pooled over the participants (N = 15)

248 to calculate  $A'$  for each *Intended Vowel Pair* at each  $f_0$  as each vowel was only presented once to  
249 each participant.

### 250 **E. Acoustic Analyses**

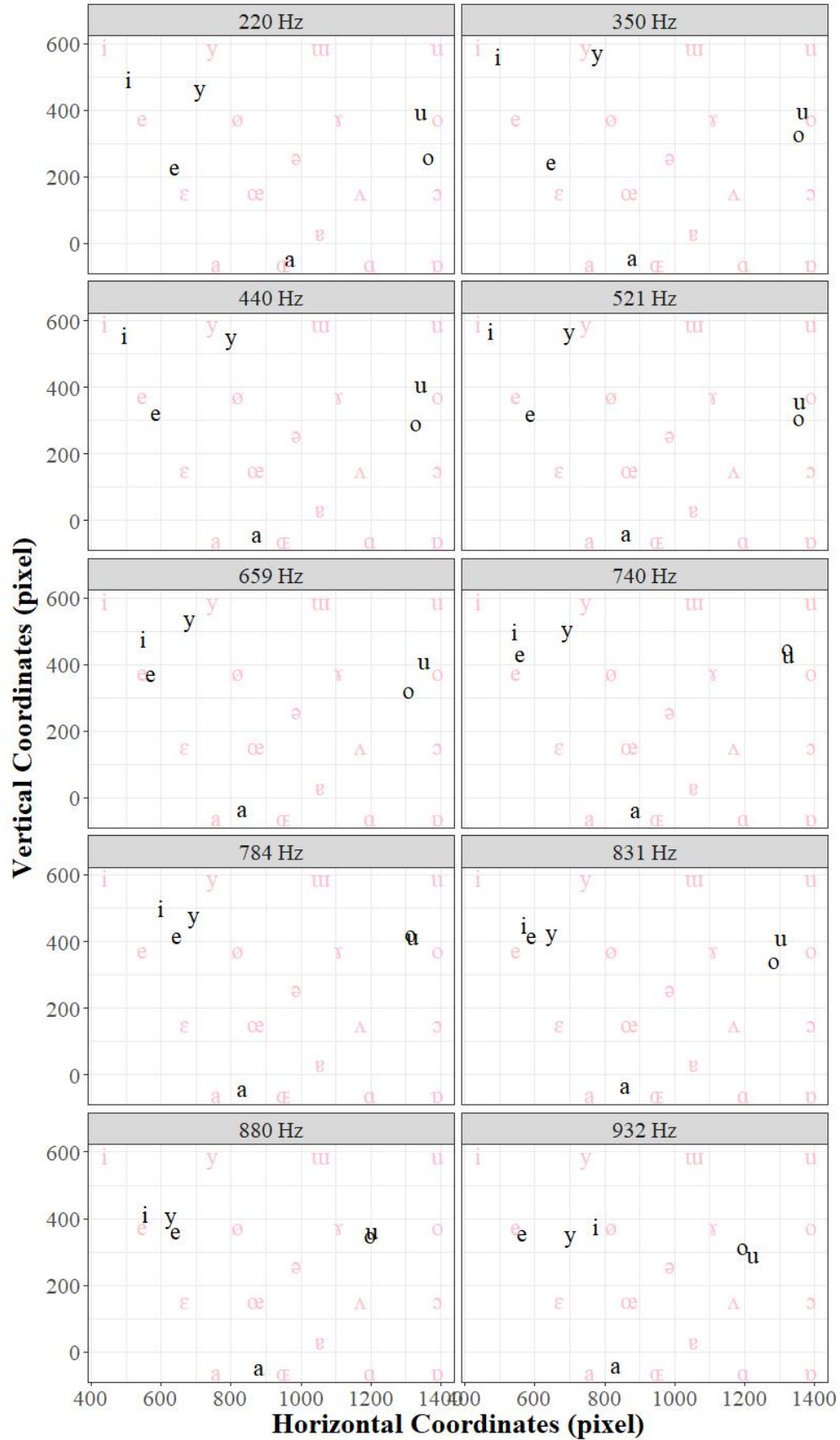
251 Acoustic analyses were conducted to help to understand the perceptual results. Simple  
252 simulated auditory excitation patterns of the vowel stimuli were computed using a 200-channel  
253 linear gammatone filter bank. The bandwidths and centre frequencies were calculated according  
254 to the ERB formulae given by Glasberg and Moore (1990). For each filter channel, the rms level  
255 of the output wave was calculated and converted to dB. To account for the transmission  
256 properties of the middle ear, a frequency weighting based on measurements made by Puria et al.  
257 (1997) was applied.

258 Classical multidimensional scaling (MDS) analysis (Shepard, 1962a, b) of the simulated  
259 auditory excitation patterns was further employed to geometrically model vowel changes at  
260 higher  $f_0$ s in the auditory perceptual space. MDS has been shown in previous studies (e.g.,  
261 Iverson & Kuhl, 1995; Kewley-Port & Atal, 1989) to be a good technique to illustrate the  
262 perceptual similarity of vowels. Each vowel at each  $f_0$  was assigned to a point in a two-  
263 dimensional geometric space with distances in the MDS space linearly related to spectral  
264 distance. Hence, MDS can map the correspondence between perceptual and acoustic properties  
265 and show acoustic differences between and among phonetic categories across the different  $f_0$ s.

## 266 **III. Results**

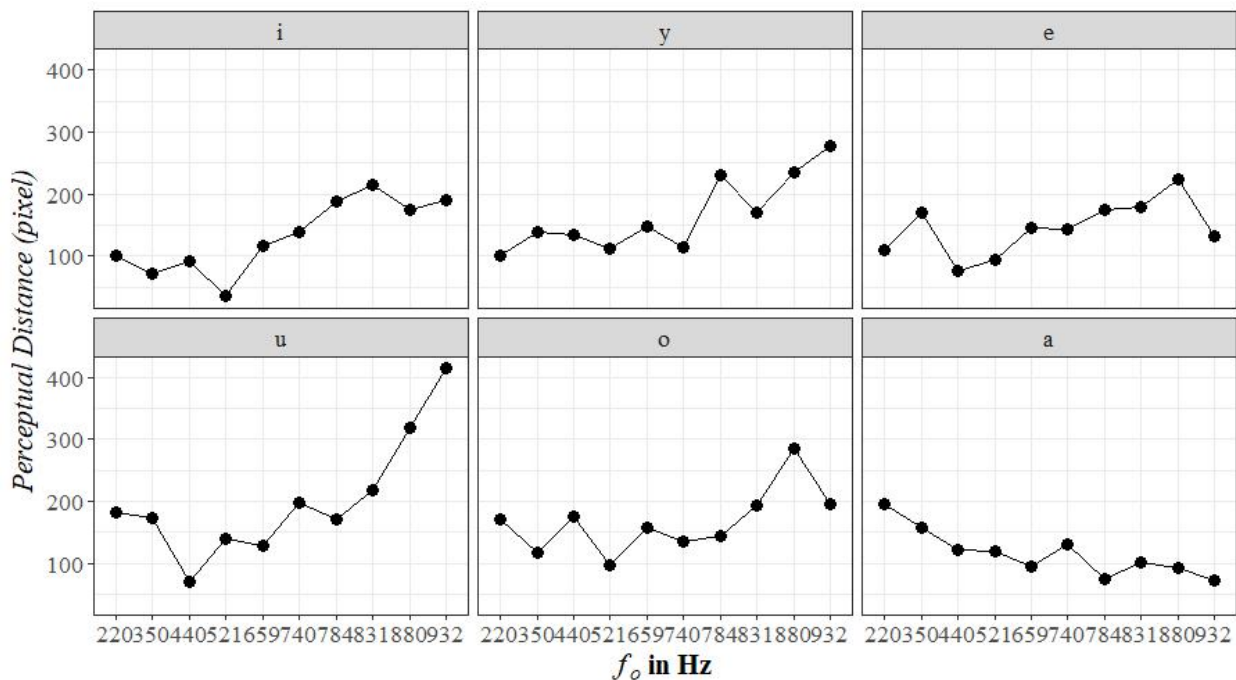
### 267 **A. Perceptual Experiments**

268 *Guided Transcription Task.* FIG. 1 shows that the basic shape of the transcribed vowel  
269 quadrilateral was maintained at all  $f_0$ s, but the high front vowels /i y e/ as well as the high back  
270 vowels /u o/ started to cluster together from 740 Hz.



272 FIG. 1. (Color Required). Results of the guided transcription task for all  $f_0$ s. Transcribed vowels  
 273 are plotted in black at the averaged coordinates of the clicks. The reference vowels are shown in  
 274 light red (i.e., the vowel quadrilateral shown to the participants). Note that the scales on the x and  
 275 y-axes do not represent frequencies but numeric coordinates on a 20-inch screen (1600 × 900  
 276 pixels).

277 The clustering involved /i y u/ moving towards the categories with higher  $F_1$ , namely,  
 278 towards [e ø o], and /e o/ moving up towards [i u] (henceforth, we describe phonemic vowels in  
 279 Yue Opera in ‘/\_/’ and the vowels perceived by the listeners in ‘[ ]’). From 831 Hz, the  
 280 perceived categories of high vowels /i y u/ all shifted further towards the vowels with the next  
 281 higher  $F_1$ . At the highest  $f_0$ , 932 Hz, /i/ was almost perceived as [e], and the perceived category of  
 282 /u/ was close to [Λ]. A closer examination of the *Perceptual Distance* (the distance between the  
 283 average vowel placements in the guided transcription task and the relevant reference vowel on  
 284 the quadrilateral) revealed that all the vowels except /a/ increased in mean *Perceptual Distance*  
 285 from an  $f_0$  of 521 Hz (FIG. 2).



286

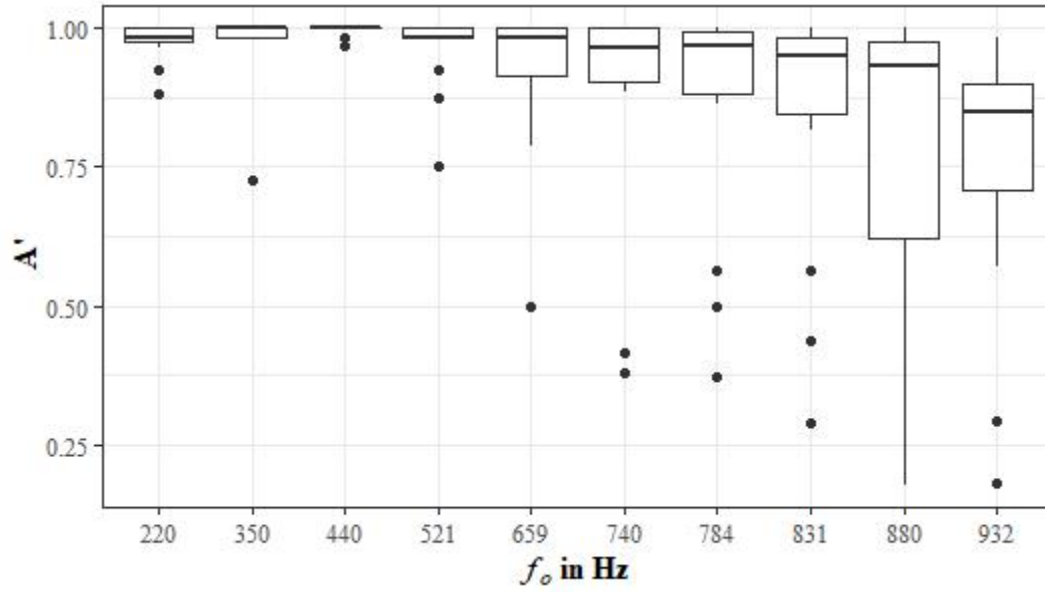
287 FIG. 2. Average *Perceptual Distance* for each vowel at all  $f_0$ s. Note that the scale on the y-axes  
 288 represents the distance in pixels between clicks and reference vowels on a 20-inch screen (1600  
 289  $\times$  900 pixels).

290 The smallest average *Perceptual Distance* was often not found at the lowest  $f_0$  (220 Hz),  
 291 but at the next higher  $f_0$ s, which correspond to the singer’s speaking  $f_0$  range. LME further  
 292 revealed highly significant effects of *Intended Vowel*,  $f_0$ , as well as their interactions (TABLE I).  
 293 The pairwise comparison (Tukey test) confirms that the differences in *Perceptual Distance* are  
 294 significant ( $ps < 0.001$ ) between the high  $f_0$ s (880 and 932 Hz) and the relatively low  $f_0$ s (440 and  
 295 521 Hz).

296 TABLE I. Results of the linear mixed-effects model on *Perceptual Distance* (Significance levels  
 297 \* =.05, \*\* = .01, \*\*\* = .001)

Final Model	Perceptual Distance $\sim f_0 + \text{Intended Vowel} + f_0: \text{Intended Vowel} + (1 \backslash \text{Participant})$			
	SS	df	F	p
$f_0$	115.518	9	7.973	<.0001***
Intended Vowel	59.782	5	7.43	<.0001***
$f_0: \text{Intended Vowel}$	131.85	45	1.82	<.0001***

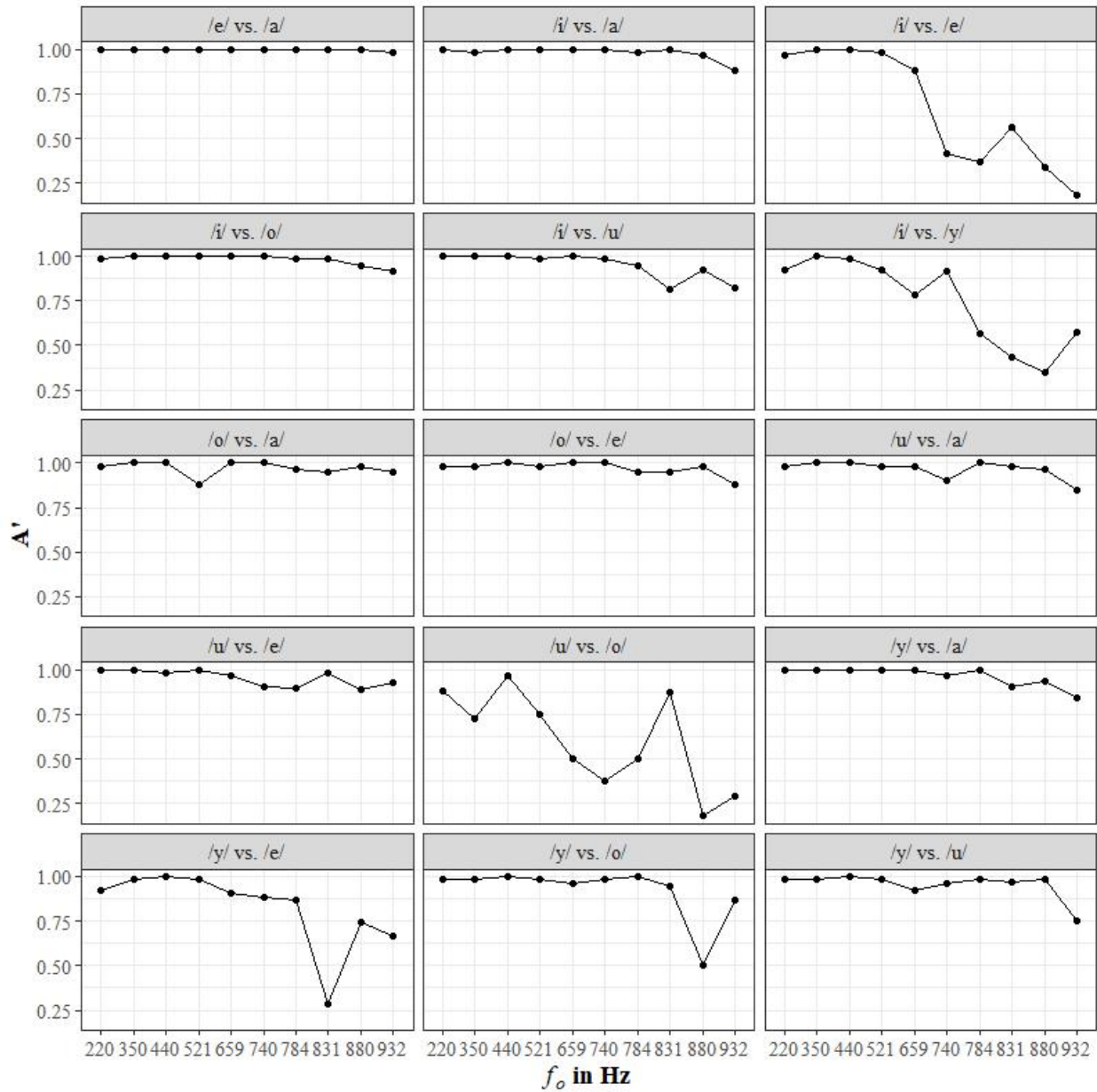
298  
 299 *Two-alternative forced-choice task.* A high identification accuracy was found throughout  
 300 all  $f_0$ s up to 932 Hz with median A' above 0.75 (FIG. 3).



301  
 302 FIG. 3. Box plots showing the distributions of  $A'$  (y-axis) for all vowel pairs that were tested at  
 303 ten  $f_0$ s between 220 and 932 Hz (x-axis).  $A'$  of 0.5 represents chance level, and  $A'$  of 1 represents  
 304 maximum performance.

305 Vowel pairs involving the low vowel /a/ or pairs composed of front and back low vowels  
 306 showed a stable and high identification accuracy across  $f_0$ s up to 831 Hz (FIG. 4). At this  $f_0$ ,  $A'$   
 307 for these pairs ranged roughly between 0.75 and 1, except for /y-o/. In contrast, from an  $f_0$  of 659  
 308 and 740 Hz, respectively,  $A'$  values for the pairs /u-o/ and /i-e/ dropped to chance level. The  
 309 same observation was made for /i-y/ from 784 Hz upwards and for /y-e/ at 831 Hz before  
 310 showing higher  $A'$  again at the two highest  $f_0$ s.





311

312

FIG. 4.  $A'$  (y-axis) for each of the vowel pair contrasts at the ten investigated  $f_o$ s (x-axis).  $A'$  of

313

0.5 represents chance level, and  $A'$  of 1 represents maximum performance.

314

Listener bias calculation is not meaningful when  $A'$  is high as it is only based on a small

315

number of misses or false alarms (Stanislaw and Todorov, 1999). We therefore only calculated

316

$B''_D$  for the vowel pairs /i-e/, /i-y/, /u-o/, /y-e/, and /y-o/ at the highest  $f_o$ s from 659 Hz as  $A'$  was

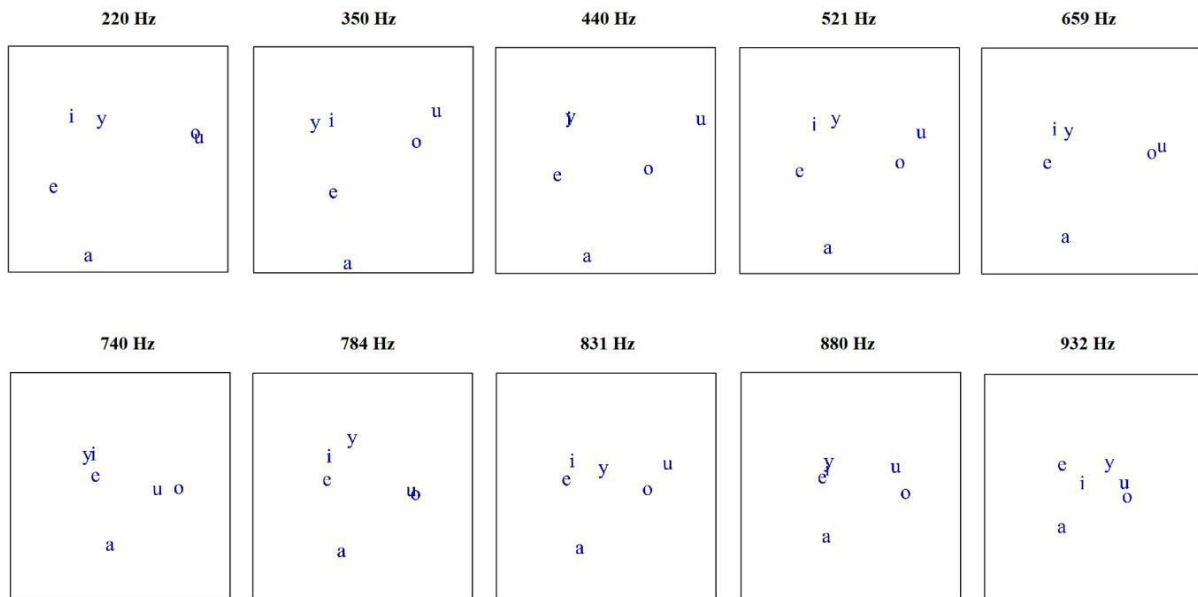
317

smaller than 0.7 in these cases. No consistent bias was found for /i-e/, /i-y/, and /u-o/ as results

318 revealed both positive and negative B''D values with a high absolute value (e.g., -0.833 and  
 319 0.894). While no bias was found for the pair /y-e/ at any of the  $f_0$ s, a strong bias towards /o/ was  
 320 found for the pair /y-o/ at an  $f_0$  of 880 Hz, but at no other frequency. In other words, no  
 321 consistent bias towards low vowels with higher  $F_{1S}$ s could be observed.

322 **B. Acoustics-derived auditory simulation**

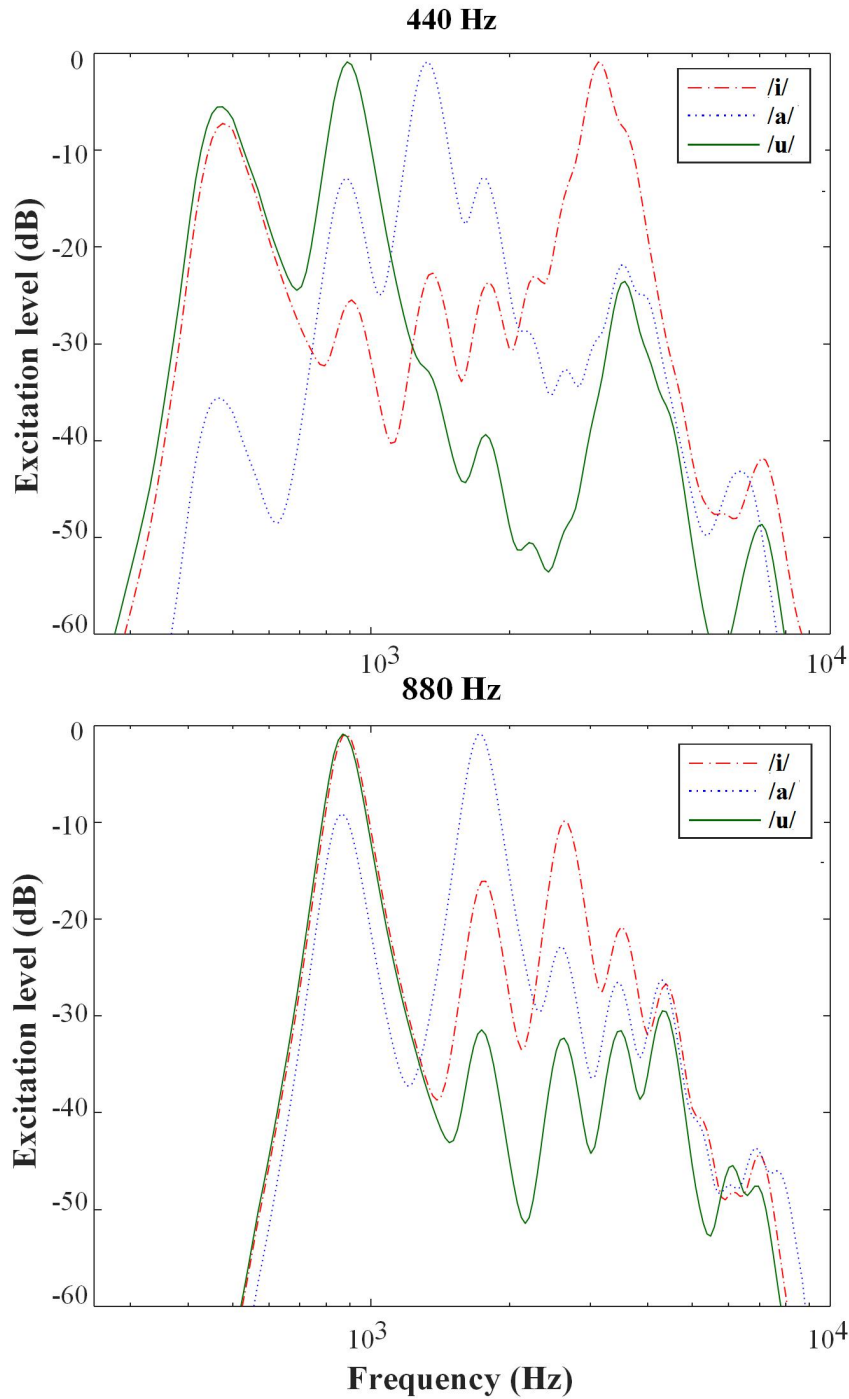
323 The results of the MDS analysis (FIG. 5) show the distribution of the stimulus vowels in  
 324 a two-dimensional space derived from the spectral similarity of the simulated auditory excitation  
 325 patterns. The spectral distances between the vowels at each  $f_0$  resembled to a high degree the  
 326 perceptual results. Above 521 Hz, high front vowels (/y e/) started to cluster around /i/ while /u/  
 327 and /o/ started to cluster together and /a/ remained clearly separated from all other vowels. At the  
 328 highest  $f_0$ s, 880 and 932 Hz, the shape of the vowel quadrilateral was considerably less clear than  
 329 in the perceptual space derived from the guided transcription.



330  
 331 FIG. 5. MDS plots showing the auditory perceptual distance between the vowels used in this  
 332 study throughout the  $f_0$ s between 220 and 932 Hz. The differences between the vowels were

333 derived from simple versions of excitation patterns that they would be expected to generate in  
334 the auditory periphery.

335         A closer examination of the individual excitation patterns showed that despite the severe  
336 under-sampling of the vocal tract transfer function at very high  $f_s$ s, the vowels /i u a/ still  
337 exhibited distinctive features up to at least 880 Hz (FIG. 6).



338

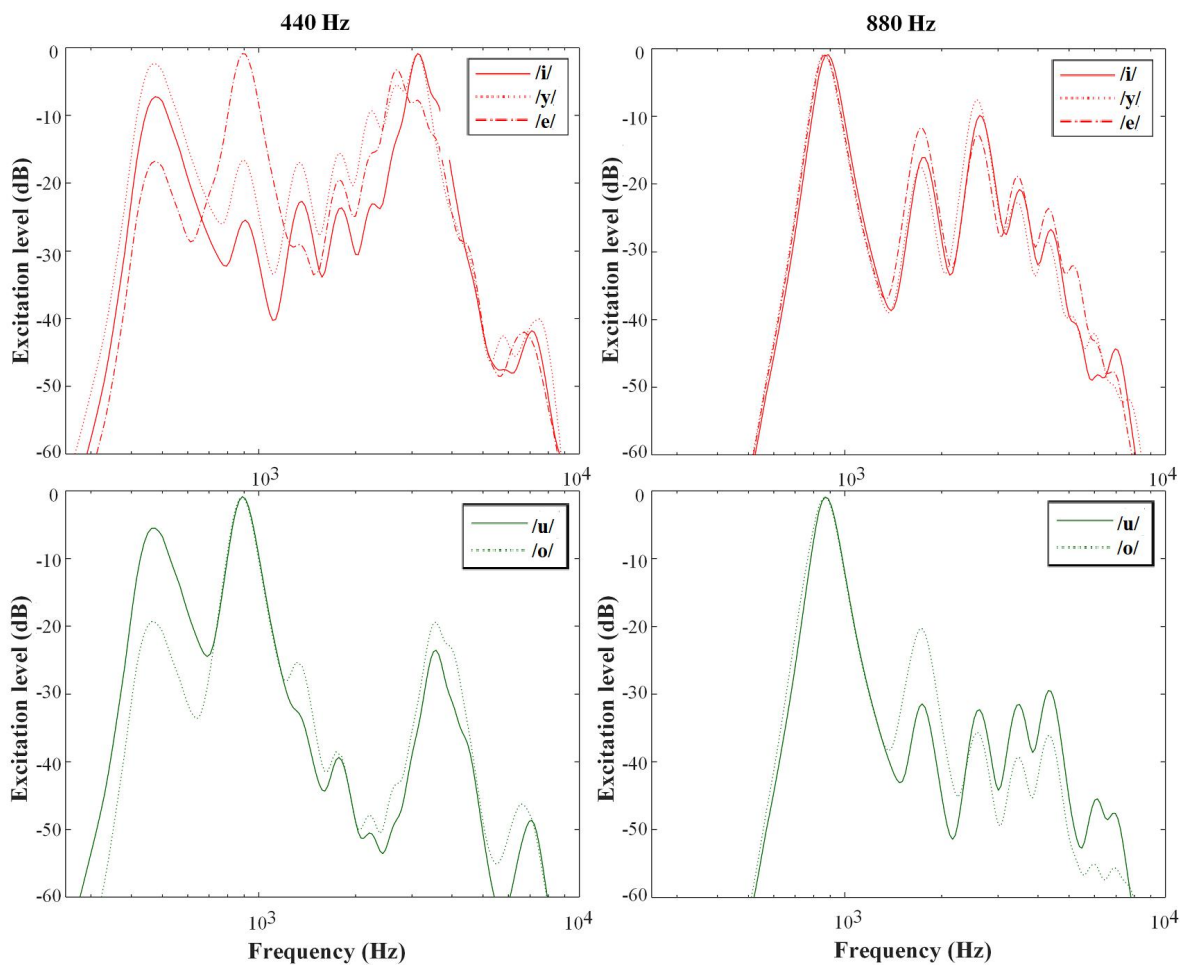
339 FIG. 6 (Color Required) Simulated auditory excitation patterns of the isolated vowels /i a u/ used

340 in this study at an octave interval with  $f_o$ s of 440 and 880 Hz. The excitation patterns reveal

341 highly differentiable spectral representations at both  $f_o$ s. At the higher  $f_o$ , the overall excitation

342 level in the frequency region above about 1.5 kHz can easily be distinguished. (The information  
343 in this figure may not be properly conveyed in black and white.)

344 The excitation patterns also showed a high degree of correspondence to the previously  
345 found confusion patterns in both perceptual tasks, namely, as  $f_0$  increased, the front high vowels  
346 tended to cluster together as well as the back vowels /u o/. For instance, the excitation patterns of  
347 /i y e/ and /u o/ at 880 Hz showed high within-group similarities (FIG. 7).



348  
349 FIG. 7 (Color Required) Simulated auditory excitation patterns of the vowel groups /i y e/ and /u  
350 o/ used in this study. Both groups were found to cluster in the perceptual space at higher  $f_0$ s. The

351 excitation patterns for 880 Hz reveal similar overall spectral shapes. (The information in this  
352 figure may not be properly conveyed in black and white.)

#### 353 **IV. Discussion**

354 The results of the present study reveal a perceptual clustering phenomenon of the high  
355 front vowels /i y e/ and the high back vowels /u o/ at high  $f_0$ s, with both clusters being highly  
356 differentiable from each other and from /a/. There was a considerable reduction in vowel  
357 distinctiveness in Yue Opera style singing within these clusters at higher pitches. However, the  
358 findings do not support the previously made assumption that listeners' identification would bias  
359 towards open vowels like [a] which many studies using vowels from Western Opera singers have  
360 suggested before (e.g., Nolan & Sykes, 2016; for an overview, see also Sundberg, 2013).

361 Furthermore, high vowels with low  $F_1$  such as /i/ and /u/ were not always the first to lose  
362 their intelligibility, as found in some studies (Hollien et al., 2000, Howie and Delattre, 1962), but  
363 could even remain identifiable up to the highest  $f_0$ . The findings of the present study might, in  
364 fact, explain the results of Smith and Scott's research (1980) who reported good identification  
365 accuracy (70 % correct) for the vowels /i ɪ e æ/ when they were presented in a non-operatic style  
366 and in isolation at  $f_0$ s around 880 Hz (i.e., A5). It seems likely that listeners could distinguish  
367 well between the two vowel pairs /i ɪ/ and /e æ/ but not always the vowels within the pair.

368 The guided transcription task used in this study revealed that the perceptual space  
369 resembles the basic shape of the vowel quadrilateral up to high registers in Yue Opera singing.  
370 However, the high front vowels /i y e/ started to cluster from 740 Hz and the high back vowels /u  
371 o/ from 659 Hz. Furthermore, towards the highest  $f_0$  investigated, the perceived categories of /i y/  
372 shifted towards the lower categories [e, ø] and /u o/ towards [ɤ, ʌ]. The latter shift of /u o/ not

373 only started from lower  $f_0$  but was also the most extensive observed in terms of perceptual  
374 distance. In contrast, the perceived category of /a/ remained accurate and stable across all the  $f_0$ s.

375         The major confusions in the two-alternative forced-choice task were found between  
376 vowel pairs within a cluster, namely, /i-e/, /i-y/, /y-e/, and /u-o/ at 659 Hz and above. No  
377 significant confusions were observed between clusters, that is, for vowels from different clusters.  
378 As it is likely that no typical formant frequency distribution could be found in such vowel pairs  
379 at very high  $f_0$ s, it seems likely that the overall spectral shape may carry enough acoustic  
380 information to distinguish between clusters, but not vowels within clusters.

381         Calculations of auditory excitation patterns and MDS analyses revealed apparent spectral  
382 differences between clusters and thus supported this hypothesis. Closer examination of the  
383 auditory excitation patterns revealed that the vowels constituting the observed clusters retain  
384 distinct spectral shapes, which kept them distinguishable from each other and /a/ throughout the  
385  $f_0$  range investigated. However, at higher  $f_0$ s, the vowels within each cluster exhibited very  
386 similar spectral shapes to one another, which may explain the decrease in listeners' identification  
387 performance in the two-alternative choice task for vowels within a cluster. These results indicate  
388 that highly differentiable overall spectral shapes (e.g., those representing [i a u]) can be used by  
389 listeners as acoustic landmarks to maintain some degree of vowel category perception at very  
390 high pitches. The calculations of the excitation patterns used in this study revealed distinct  
391 excitation levels in the frequency region above roughly 1.5 kHz for the vowels /i a u/, but highly  
392 similar levels for vowels within the clusters. Therefore, the present findings support the view that  
393 models of vowel perception based on formant peak patterns cannot provide such a full account of  
394 vowel perception as theories based on overall spectral shape (for a comprehensive review of  
395 several overall-spectral-shape models, see Kiefte et al., 2013).

396 This is also supported by the MDS analyses, which have shown that a triangular  
397 distribution of the six Yue Opera vowels could be observed up to about 880 Hz, though on a  
398 gradually reduced scale as an increasing  $f_0$  brought compression of the vowel space in both the  
399 dimensions of the MDS plots. The distance between /i/ and /u/ (and therefore *vowel frontness*)  
400 decreased as well as that between /i/ and /a/ (and therefore *vowel height*) (see FIG. 6). In contrast  
401 to Friedrichs et al. (2016), who observed an expansion in the front-back distinction when the  
402 vowel height dimension collapsed towards higher  $f_0$  of a Western Musical Theatre singer, the  
403 perceptual space containing the vowels produced by our Yue Opera singer did not show such  
404 compensation. This may either be due to the singer's personal habit or could also be because Yue  
405 Opera singers employ other mechanisms to protect the distinctiveness of sung vowels at high  
406 pitches – for instance, the association with melody and the tones embedded in the melody. A  
407 kinetic pitch on a vowel sweeps the transfer function with any available harmonic and therefore  
408 better reveals it than a static harmonic. More importantly, as previously mentioned, lexical tones  
409 embedded in the melody of an actual performance will also help lexical access by narrowing the  
410 candidates. Chinese theatre composition requires the composers and lyricists to follow the rule  
411 that the melody associated with each syllable should not conflict with the lexical tone in the  
412 beginning part, only allowing limited modification of the tone contour (Wee, 2007; Zhang, 1980:  
413 91). Although in Mandarin popular song composition, this rule may not be followed strictly (for  
414 an overview on correspondence between lexical tone and sung melody, see Schellenberg and  
415 Gick, 2020: Table 1), theatre composition is much stricter on tone-melody harmony (Zhou and  
416 You, 1997: 190).

417 A mismatch between the results of the two perceptual tasks used in the present study is  
418 worth special attention. In the guided transcription, the two vowels that typically exhibited the



419 lowest  $F_1$ , /i/ and /u/, are on average placed near [e] and [o] respectively on the response  
420 quadrilateral, but the vowels with typically the next higher  $F_1$ , /e/ and /o/, were also located in  
421 this perceptual vicinity (i.e., perceived correctly and placed near the most relevant reference  
422 vowels [e] and [o]) at 932 Hz. However, in the two-alternative forced-choice, no bias towards /e/  
423 or /o/ was found for these vowels at high  $f_0$ s. Instead, /e/ was often identified as /i/ at 880 Hz.  
424 When  $f_0$  exceeded 784 Hz, participants showed a bias towards [u] rather than [o], which would  
425 be more consistent with the findings of the transcription task. The acoustic analyses cannot fully  
426 explain these mismatched bias patterns. It may be that the participants were sensitive to the  
427 changes in vowel quality, but the categorical perception of vowels may not necessarily  
428 correspond completely to the perceived quality but was influenced as well by other factors like  
429 the task. Previous studies have indicated that different tasks do affect the listener's identification  
430 performance. For instance, participants performed better when they were presented with more  
431 meaningful response options (e.g., written words containing the target vowel vs. vowel letters),  
432 fewer response options, and a lower degree of talker variability (Friedrichs et al., 2017, 2015a,  
433 2015b).

434         The results presented here, especially the identifiability of the high vowels (i.e., those  
435 with typically low  $F_1$ ) at high  $f_0$ , and the much higher  $f_0$  at which identifiability started to decline  
436 compared to the studies on Western Opera (for an overview, see Sundberg, 2013), may also  
437 partly be driven by the features of Yue Opera. As there are no male singers in traditional Yue  
438 Opera, male, female, and even child characters in a single performance are all played by female  
439 singers. In order to distinguish between the gender and age of the different characters they are  
440 portraying, Yue Opera singers employ style-specific aesthetic and articulatory adjustments. For  
441 example, singers typically portray female characters with a more reduced mouth opening than

442 male characters. This makes resonance tuning, as described in studies on Western Opera singing,  
443 very unlikely because it requires articulatory actions to increase mouth opening (i.e., opening the  
444 jaw, widening the lips). It seems plausible that tongue height and advancement and anatomical  
445 dimensions such as those of the pharynx might play a role in distinguishing between characters  
446 and maintaining intelligibility. To investigate this further and fully understand the correlation  
447 between character gender and vowel realisation, experiments with more Yue Opera singers  
448 performing in different gender and age groups are required (for further discussion on *the*  
449 *influence of gender and age on vowel quality*, see Maurer et al., 2015). This investigation may  
450 even be expanded to other Chinese Operas, which involve males playing females. Further  
451 research in this area may also be helpful to test whether vowel clustering can solely explain the  
452 relatively high intelligibility at high  $f_0$ s or whether other factors contribute to this.

## 453 **V. Conclusion**

454 The present study on Yue Opera demonstrated that vowels clustered in the perceptual  
455 space into three groups (/i y e/, /u o/, /a/) at high  $f_0$  above about 521 Hz, and that listeners were  
456 able to identify vowels between but not within groups with high accuracy up to 932 Hz. The  
457 results, therefore, show that previous findings on vowel intelligibility in Western Opera may be  
458 style-specific and cannot be generalized to other forms of opera singing. The findings presented  
459 here furthermore support the view that the overall spectral shape provides a more robust cue than  
460 formant peak patterns for the perception of the high-pitched vowels. Further studies on  
461 articulatory strategies in high-pitched Yue Opera singing may be useful to fully understand the  
462 underlying mechanisms resulting in the perceptual clustering of vowels at high  $f_0$ s.

463

## 464 **ACKNOWLEDGEMENTS**

465 Thanks to Shuyang Sheng, the National-level Yue Opera performer in Shanghai Yue  
466 Opera House, who was recorded for this study. Further thanks go to Nick Clark, whose software  
467 was used to perform the gammatone filtering, as well as Stuart Rosen and Paul Iverson for their  
468 help with the auditory excitation patterns and MDS, respectively. The first author would also like  
469 to thank the CHINA Scholarship COUNCIL (CSC) and Cambridge Trust for their doctoral  
470 scholarship. The last author was supported by the Forschungskredit of the University of Zurich,  
471 Grant No. FK-18-077, and the Swiss National Science Foundation (SNSF), Grants No.  
472 P400PG\_180693 and P2ZHP1\_168375.

473

#### 474 REFERENCES (BIBLIOGRAPHIC)

475 Bates D.M., Mächler M., Bolker B., and Walker S. (2015). “Fitting Linear Mixed-Effects

476 Models Using lme4.” *Journal of Statistical Software*, **67**(1), 1–48.

477 doi:10.18637/jss.v067.i01.

478 Boersma, P., and Weenink, D. (2021). “Praat: Doing phonetics by computer [computer program]

479 (version 6.0.39),” <http://www.praat.org/> (Last viewed January 22, 2021).

480 Davis, S., and Mermelstein, P. (1980). “Comparison of parametric representations for

481 monosyllabic word recognition in continuously spoken sentences,” *IEEE transactions on*

482 acoustics, speech, and signal processing, **28**(4), 357–366.

483 E-Prime (Version 2.0.8.22) [Computer software]. Pittsburgh, PA: Psychology Software Tools.

484 Friedrichs, D., Maurer, D., Rosen, S., and Dellwo, V. (2017). “Vowel recognition at fundamental

485 frequencies up to 1 kHz reveals point vowels as acoustic landmarks,” *J. Acoust. Soc. Am.*

486 **142**(2), 1025–1033.

487 Friedrichs, D., Rosen, S., Iverson, P., Maurer, D., and Dellwo, V. (2016). “Mapping vowel  
488 categories at high fundamental frequencies using multidimensional saling of cochlea-  
489 scaled spectra,” *J. Acoust. Soc. Am.* **140** (1), 3219.

490 Friedrichs, D., Maurer, D., and Dellwo, V. (2015a). “The phonological function of vowels is  
491 maintained at fundamental frequencies up to 880 Hz,” *J. Acoust. Soc. Am.* **138**(1),  
492 EL36–EL42.

493 Friedrichs, D., Maurer, D., Suter, H., and Dellwo, V. (2015b). “Vowel identification at high  
494 fundamental frequencies in minimal pairs,” in *Proceedings of the 18th International  
495 Congress on Phonetic Science*, paper number 0438, pp. 1–5.

496 Glasberg, B. R., and Moore, B. C. (1990). “Derivation of auditory filter shapes from notched-  
497 noise data,” *Hearing research*, **47**(1-2), 103–138.

498 Hollien, H., Mendes-Schwartz, A. P., and Nielsen, K. (2000). “Perceptual confusions of high-  
499 pitched sung vowels,” *Journal of Voice*, **14**(2), 287–298.

500 Howie, J., and Delattre, P. (1962). “An experimental study of the effect of pitch on the  
501 intelligibility of vowels,” *Natl. Assoc. Teachers Singing Bull.* **18**(4), 6–9.

502 Huang, W. (2000). *Phonology of Yue Opera* (Doctoral dissertation, Fudan University). [黄玮.  
503 (2000). 越剧音韵研究. (博士论文, 复旦大学)].

504 Ito, M., Tsuchida, J., and Yano, M. (2001). “On the effectiveness of whole spectral shape for  
505 vowel perception,” *J. Acoust. Soc. Am.* **110**(2), 1141–1149.

506 Iverson, P., and Kuhl, P. (1995). “Mapping the perceptual magnet effect for speech using signal  
507 detection theory and multidimensional scaling,” *J. Acoust. Soc. Am.* **97**(1), 553–562.

508 Joliveau, E., Smith, J., and Wolfe, J. (2004). “Tuning of vocal tract resonance by sopranos,”  
509 *Nature*, **427**(6970), 116.

510 Kewley-Port, D., and Atal, B. S. (1989). "Perceptual differences between vowels in a limited  
511 phonetic space," *J. Acoust. Soc. Am.* **85**, 1726–1740.

512 Lehiste, I., and Peterson, G. E. (1961). "Transitions, glides, and diphthongs," *J. Acoust. Soc. Am.*  
513 **33**(3), 268–277.

514 Maurer, D. (2016). "Acoustics of the Vowel – Preliminaries," (Peter Lang AG, International  
515 Academic Publishers, Bern, Switzerland).

516 Maurer, D., Suter, H., Friedrichs, D., and Dellwo, V. (2015). "Gender and age differences in  
517 vowel-related formant patterns: What happens if men, women, and children produce  
518 vowels on different and on similar F0?," *J. Acoust. Soc. Am.* **137**(4), 2416.

519 Maurer, D., Mok, P., Friedrichs, D., and Dellwo, V. (2014). "Intelligibility of high-pitched vowel  
520 sounds in the singing and speaking of a female Cantonese Opera singer," in *15th Annual  
521 Conference of International Speech Communication Association*, 2132–2133.

522 Maurer, D., and Landis, T. (1996). "Intelligibility and spectral differences in high-pitched  
523 vowels," *Folia phoniatrica et logopaedica*, **48**(1), 1–10.

524 McKight, P. E., and Najab, J. (2010). "Kruskal-wallis test," *The corsini encyclopedia of  
525 psychology*, **1**.

526 Nolan, F., and Sykes, H. (2015). "Vowel and consonant identification at high pitch: the acoustics  
527 of soprano unintelligibility," *Proceedings of the 18th International Congress of Phonetic  
528 Sciences, Paper number 14*, 1–5.

529 Pallier, C. (2002). "Computing discriminability and bias with the R software," URL:  
530 <http://www.pallier.org/pdfs/aprime.pdf> (Last viewed October 13, 2020)

531 Pollack, I., & Norman, D. A. (1964). A non-parametric analysis of recognition  
532 experiments. *Psychonomic science*, *1*(1), 125-126.

533 Puria, S., Peake, W. T., & Rosowski, J. J. (1997). "Sound-pressure measurements in the cochlear  
534 vestibule of human-cadaver ears," *J. Acoust. Soc. Am.* **101**(5), 2754–2770.

535 Qiu, D. (1999). The Rhyme List of Yue Opera. *Theatre Lyrics* (4), 4.[裘达人. (1999). 越剧唱词  
536 音韵表. *戏文*(4), 4.]

537 Qiu, D. (1995). On the compilation of the Rhyme List of Yue Opera. *Theatre Lyrics* (1), 2.[裘达  
538 人. (1995). 关于修订《越剧唱词音韵表》的琐见. *戏文*(1), 2.]

539 R Core Team. (2020). "R: A language and environment for statistical computing [computer  
540 software] (version 3.1.3.)," R Foundation for Statistical Computing, Vienna, Austria,  
541 <https://www.R-project.org/>.

542 Schellenberg, M., & Gick, B. (2020). Microtonal variation in sung Cantonese. *Phonetica*, *77*(2),  
543 83-106.

544 Shepard, R. N. (1962a). "The analysis of proximities: Multidimensional scaling with an  
545 unknown distance function. I.," *Psychometrika* **27**, 125–140.

546 Shepard, R. N. (1962b). "The analysis of proximities: Multidimensional scaling with an  
547 unknown distance function. II.," *Psychometrika* **27**, 219–246.

548 Smith, L. A., and Scott, B. L. (1980). "Increasing the intelligibility of sung vowels," *J. Acoust.*  
549 *Soc. Am.* **67**(5), 1795–1797.

550 Smith, W.D. (1995). Clarification of sensitivity measure A'. *Journal of Mathematical*  
551 *Psychology* *39*, 82–89

552 Stanislaw, H., and Todorov, N. (1999). "Calculation of signal detection theory measures," *Behav.*  
553 *Res. Methods, Instrum., Comput.* **31**, 137–149.

554 Strange, W., Verbrugge, R. R., Shankweiler, D. P., and Edman, T. R. (1976). "Consonant  
555 environment specifies vowel identity," *J. Acoust. Soc. Am.* **60**, 213–224.

- 556 Sundberg, J., Lã, F. M., and Gill, B. P. (2013). "Formant tuning strategies in professional male  
557 opera singers," *Journal of Voice*, **27**(3), 278–288.
- 558 Sundberg, J. (2012). "Perception of singing," in *Psychology of Music*, 3rd ed., edited by D.  
559 Deutsch (Academic Press, London), pp. 69–106.
- 560 Tanner, W.P., Jr., & Swets, J.A. (1954). A decision-making theory of visual detection.  
561 *Psychological Review* *61*, 401–409.
- 562 Wee, L. H. (2007). "Unraveling the relation between Mandarin tones and musical melody,"  
563 *Journal of Chinese Linguistics*, **35**(1), 128–144.
- 564 You, R. (2006). *Phonology of Chinese Operas*. The Commercial Press. [游汝杰. (2006). *地方戏*  
565 *曲音韵研究*. 商务印书馆.]
- 566 Zahorian, S. A., & Jagharghi, A. J. (1993). "Spectral-shape features versus formants as acoustic  
567 correlates for vowels," *J. Acoust. Soc. Am.* **94**(4), 1966–1982.
- 568 Zhang, G. (1980). *The art of Chinese Opera*. China Drama Press. [张庚. (1980). *戏曲艺术论*. 中  
569 国戏剧出版社.]
- 570 Zhang, J., & Mueller, S. T. (2005). A note on ROC analysis and non-parametric estimate of  
571 sensitivity. *Psychometrika*, **70**(1), 203–212. doi:10.1007/s11336-003-1119-8
- 572 Zhou, Z., & You, R. (1997). *Dialects and Chinese Culture*. Shanghai People's Publishing House.  
573 [周振鹤&游汝杰. (1997). *方言和中国文化*. 上海人民出版社.]