

# The effects of high versus low talker variability and individual aptitude on phonetic training of Mandarin lexical tones

Hanyu Dong <sup>Corresp.</sup> <sup>1</sup>, Meghan Clayards <sup>2</sup>, Helen Brown <sup>3</sup>, Elizabeth Wonnacott <sup>Corresp.</sup> <sup>1</sup>

<sup>1</sup> Division of Psychology and Language Sciences, University College London, University of London, London, United Kingdom

<sup>2</sup> Department of Linguistics, School of Communications Sciences and Disorders, McGill University, Montreal, QC, Canada

<sup>3</sup> Department of Psychology, Nottingham Trent University, Nottingham, United Kingdom

Corresponding Authors: Hanyu Dong, Elizabeth Wonnacott  
Email address: hanyu.dong.10@ucl.ac.uk, e.wonnacott@ucl.ac.uk

High variability training has been found more effective than low variability training in learning various non-native phonetic contrasts. However, little research has considered whether this applies to the learning of tone contrasts. The only two relevant studies suggested that the effect of high variability training depends on the perceptual aptitude of participants (Perrachione, Lee, Ha, & Wong, 2011; Sadakata & McQueen, 2014). The present study extends these findings by examining the interaction between individual aptitude and input variability using natural, meaningful L2 input (both previous studies used pseudowords). Sixty English speakers took part in an eight session phonetic training paradigm. They were assigned to high/low/high-blocking variability training groups and learned real Mandarin tones and words. Individual aptitude was measured following previous work. Learning was measured using one discrimination task, one identification task and two production tasks. All tasks assessed the generalisation of learning. Overall, all groups improved in both production and perception of tones which transferred to novel voices and items, demonstrating the effectiveness of training despite the increased complexity compared with previous research. Although the low variability group exhibited an advantage with the training stimuli, there was no evidence that the different variability training led to different performance in any of the tests of generalisation. Moreover, although aptitude significantly predicted performance in discrimination, identification and training tasks, no interaction between individual aptitude and variability was revealed. We discuss these results in light of previous findings.

1

2 **The effects of high versus low talker variability and individual aptitude on phonetic**  
3 **training of Mandarin lexical tones**

4 Hanyu Dong<sup>1</sup>, Meghan Clayards<sup>2</sup>, Helen Brown<sup>3</sup>, & Elizabeth Wonnacott<sup>1</sup>

5 *<sup>1</sup>Division of Psychology and Language Sciences, University College London, London, UK*

6 *<sup>2</sup>Department of Linguistics, School of Communications Sciences and Disorders, McGill*  
7 *University, Montreal, QC, Canada*

8 *<sup>3</sup> Department of Psychology, Nottingham Trent University, Nottingham, UK*

9

10 Correspondence concerning this article should be addressed to Elizabeth Wonnacott, Division of  
11 Psychology and Language Sciences, Chandler House, 2 Wakefield Street, London, WC1N 1PF.  
12 Email: [e.wonnacott@ucl.ac.uk](mailto:e.wonnacott@ucl.ac.uk)

13 Declarations of interest: None

**14 Abstract**

15 High variability training has been found more effective than low variability training in learning  
16 various non-native phonetic contrasts. However, little research has considered whether this  
17 applies to the learning of tone contrasts. The only two relevant studies suggested that the effect  
18 of high variability training depends on the perceptual aptitude of participants (Perrachione, Lee,  
19 Ha, & Wong, 2011; Sadakata & McQueen, 2014). The present study extends these findings by  
20 examining the interaction between individual aptitude and input variability using natural,  
21 meaningful L2 input (both previous studies used pseudowords). Sixty English speakers took part  
22 in an eight session phonetic training paradigm. They were assigned to high/low/high-blocking  
23 variability training groups and learned real Mandarin tones and words. Individual aptitude was  
24 measured following previous work. Learning was measured using one discrimination task, one  
25 identification task and two production tasks. All tasks assessed the generalisation of learning.  
26 Overall, all groups improved in both production and perception of tones which transferred to  
27 novel voices and items, demonstrating the effectiveness of training despite the increased  
28 complexity compared with previous research. Although the low variability group exhibited an  
29 advantage with the training stimuli, there was no evidence that the different variability training  
30 led to different performance in any of the tests of generalisation. Moreover, although aptitude  
31 significantly predicted performance in discrimination, identification and training tasks, no  
32 interaction between individual aptitude and variability was revealed. We discuss these results in  
33 light of previous findings.

34 *Keywords: Phonetic training; L2 phonetic contrasts; Lexical tone learning*

35

## 36 **1 Introduction**

37 One challenging aspect of learning a second language (L2) is learning to accurately  
38 perceive non-native phonetic categories. This task is particular difficulty where the L2 contains  
39 the same acoustic properties as the first language (L1), but used differently (Bygate, Swain, &  
40 Skehan, 2013), suggesting that it is challenging to adjust existing acoustic properties in the L1 to  
41 learn new L2 categories. This challenge is compounded by the fact that speech is highly variable  
42 in the natural linguistic environment. Variability comes not only from the phonetic context but  
43 also from differences between speakers. Thus, learners must learn to distinguish the new L2  
44 categories despite all the variability present in the learning input. There is evidence that native  
45 listeners can process this variability in speech faster and more accurately than non-native  
46 listeners (Bradlow & Pisoni, 1999), indicating that it is indeed a challenge for L2 learners.  
47 Despite this, it has been suggested that input variability may be beneficial for second language  
48 learning and generalization (Barcroft & Sommers, 2005; Lively, Logan & Pisonni, 1993).  
49 However recent evidence suggests that the ability to benefit from variability may depend on  
50 individual learner aptitude (Perrachione, Lee, Ha, & Wong, 2011; Sadakata & McQueen, 2014),  
51 at least in the learning of lexical tones i.e. the distinctive pitch patterns carried by the syllable of  
52 a word which, in certain languages, distinguish meaningful lexical contrasts. The current paper  
53 further explores how and when variability supports or impedes learning of new L2 phonetic  
54 categories, focusing on English learners of Mandarin tone contrasts.

### 55 **1.1 High Variability L2 Phonetic Training for Non-Tonal Contrasts**

56 A substantial body of literature has explored whether phonetic training can be used to  
57 improve identification and discrimination of non-native phonetic contrasts in L2 learners. An  
58 early study by Strange and Dittman (1984) attempted to train Japanese speakers on the English

59 /r/- /l/ distinction, a phoneme contrasts that does not exist in Japanese. This training study used a  
60 discrimination task in which participants made same–different judgments about stimuli from a  
61 synthetic *rock-lock* continuum, receiving immediate trial-by-trial feedback. Participants were  
62 given a variety of discrimination and identification tasks pre- and post-training. The key result  
63 was that although performance increased both for trained items on the synthesized *rock-lock*  
64 continuum, and for novel items on a synthesized *rake-lake* continuum, participants failed to  
65 show any improvement for naturally produced minimal pair speech tokens. Later research  
66 suggested that a key factor which prevented generalization to natural speech tokens was a lack of  
67 variability in the training materials: Variability was present in the form of the ambiguous  
68 intermediate stimuli along the continuum, however, there was a single phonetic context and a  
69 single (synthesized) speaker. Logan, Lively, and Pisoni (1991) also trained Japanese learners on  
70 the English /r/-/l/ contrast, but included multiple natural exemplars (67 minimal pairs, where the  
71 target speech sounds appeared in different phonetic contexts) and multiple speakers (four males  
72 and two females). Their pre- and post- training tests involved novel and trained words spoken by  
73 both trained and novel speakers. In contrast to Strange and Dittman, they found that participants  
74 successfully generalized to both new speakers and new words. This was the first study to indicate  
75 the importance of variability within the training material. A follow up study by Lively, Logan,  
76 and Pisoni (1993) provided further evidence for this by contrasting a condition with *high*  
77 *variability* input with one with *low variability* input in which the stimuli were spoken by a single  
78 speaker (although still exemplified in multiple phonetic environments). Participants in the low  
79 variability group improved during the training sessions but failed to generalise this learning to  
80 new speakers.

81           Following Logan et al. (1993) the use of high variability training materials has become  
82 standard in L2 phonetic training – the so called “*high variability phonetic training*” (HVPT)  
83 methodology. This methodology has been successfully extended to training a variety of contrasts  
84 in various languages such as learning of the English /u:/-/ʊ/ distinction by Catalan/Spanish  
85 bilinguals (Aliaga-García & Mora, 2009), learning of the English /i:/-/ɪ/ contrasts by native  
86 Greek speakers (Lengeris & Hazan, 2010; Giannakopoulou, Uther & Ylinen, 2013), and learning  
87 of the English /w/-/v/ distinction by native German speakers (Iverson, Ekanayake, Hamann,  
88 Sennema, & Evans, 2008).

89           There is also some evidence that this type of perceptual training benefits production in  
90 addition to perception. Bradlow, Pisoni, Akahane-Yamada, and Tohkura (1997) found that  
91 production of the /r/-/l/ contrast improved in Japanese speakers following HVPT, with this  
92 improvement being retained even after three months. Similar improvement on the production of  
93 American English mid to low vowels by Japanese’s speakers following HVPT was also reported  
94 by Lambacher, Martens, Kakehi, Marasinghe, and Molholt (2005). However, the evidence here  
95 is mixed: a recent study (Alshangiti & Evans, 2014) employed HVPT to train Arabic learners on  
96 non-native English vowel contrasts and found no improvements in production, although  
97 participants receiving additional explicit production training did show some limited  
98 improvement.

99           The finding that variability boosts generalization is intuitively sensible: Experience of  
100 variation allows the formation of generalized representations that include only phonetically  
101 relevant cues and exclude irrelevant speaker identity cues. However it is notable that the seminal  
102 experiments of Logan and colleagues had a small sample (the tests of generalization were  
103 administered to only three of the participants in Logan et al. 1991), and since this work,

104 relatively few studies have explicitly tested the benefit of high variability training by directly  
105 comparing high variability and low variability training conditions. Clopper and Pisoni (2004)  
106 found a benefit of high variability, although this focused on dialect categorization rather than L2  
107 phonetic learning. They tested participants' ability to categorize dialects following exposure to  
108 high variability training (three speakers per dialect) compared with low variability training (one  
109 speaker per dialect), finding better generalization after high variability training. Sadakata and  
110 McQueen (2013) trained native Dutch speakers with geminate and singleton variants of the  
111 Japanese fricative /s/. Participants were trained with either a limited set of words recorded by a  
112 single speaker (low-variability) or with a more variable set of words recorded by multiple  
113 speakers (high-variability). Critically, the total amount of exposure to the contrast was held  
114 constant across conditions such that each item in the low-variability condition was repeated more  
115 frequently than each item in the high-variability condition. Both types of training led to increases  
116 in both the identification and discrimination of the novel contrast, including generalization to  
117 untrained fricatives and speakers, however for the identification task the improvement was  
118 greater following high variability training.

119 More recently, Giannakopoulou, Brown, Clayards, and Wonnacott (2017) compared  
120 matched high variability (four speakers) and low variability (one speaker) training for adult and  
121 child (8 year old) native Greek speakers who were trained on the English /i:/-/i/ contrast. In  
122 contrast to the results of Logan et al. (1993), this study did *not* show a benefit for high variability  
123 compared to low variability training in either age group, even for generalization items. However,  
124 for adult participants, it is unclear the extent to which this was due to ceiling effects. Two other  
125 previous studies which specifically manipulated variability during learning of novel phonetic  
126 categories are those by Perrachione, Lee, Ha, & Wong (2011) and Sadakata and McQueen

127 (2014) which both looked at the learning of lexical tone. We discuss these studies in more detail  
128 in the following section.

129 Finally, there is also evidence of a benefit of high variability training in L2 vocabulary  
130 learning: With more varied training materials, (either multiple speakers or multiple voice quality  
131 types) participants show greater learning in both production and reception tests (Barcroft &  
132 Sommers, 2005, 2014; Sommers & Barcroft, 2007, 2011).

### 133 **1.2 Phonetic Training of L2 Lexical Tones**

134 Each of the phonetic training studies discussed above involved training a *segmental*  
135 contrast (consonantal or vocalic). Another type of phonological contrast which exists in some  
136 natural languages is lexical tone, whereby the pitch contour is used to distinguish lexical or  
137 grammatical meanings (Yip, 2002). For example, Mandarin Chinese has four lexical tones; level-  
138 tone (Tone 1), rising-tone (Tone 2), dipping tone (Tone 3) and falling-tone (Tone 4). These pitch  
139 contours combine with syllables to distinguish meanings. For instance, the syllable *ba* combines  
140 with the four tones to mean: eight (*bā*, Tone 1), pluck (*bá*, Tone 2), grasp (*bǎ*, Tone 3) and father  
141 (*bà*, Tone 4). Each of these words thus forms a minimal pair with each of the others. Note that  
142 while languages such as English use pitch information extensively for intonation – such as  
143 forming a question or for emphasis – they do not use pitch information lexically, causing  
144 difficulties for learners of Mandarin as an L2.

145 The first study examining lexical tone training was conducted by Wang, Spence,  
146 Jongman, and Sereno (1999). A similar paradigm to that used by Logan et al. (1991) was  
147 adopted using four speakers for training. Training consisted of a two-alternative forced choice  
148 (2AFC) task in which participants heard a syllable whilst viewing two standard diacritic  
149 representations (i.e., →, ↗, v, ↘, which are iconic in nature). They were asked to pick out the



150 picture of the arrow that corresponded to the tone and received feedback. At test, participants  
151 chose which tone they had heard out of a choice of all four (4AFC task). There were also two  
152 generalisation tasks, one with 60 new words produced by one of the training speakers, and the  
153 other with an additional 60 new words produced by a new speaker. Training materials were all  
154 real monosyllabic Mandarin words that varied in the consonants, vowels and syllable structure.  
155 Native speakers of American English showed significant improvement in the accuracy of tone  
156 identification after eight sessions of high variability training over two weeks and this generalized  
157 to both new words and new speakers.

158         In a follow up study, Wang, Jongman and Sereno (2003) used the same training paradigm  
159 to test whether learning transferred to production. They recruited participants taking Mandarin  
160 courses and asked them to read through a list of 80 Mandarin words written in Pinyin (an  
161 alphabetic transcription) before and after training. These production were rated by 82 native  
162 Mandarin speakers blind to whether each recording was from pre- or post-test. They found  
163 improvements in production, although these were mainly seen in pitch height rather than pitch  
164 contour.

165         These studies suggested that as with segmental phoneme contrasts, high variability  
166 training could also facilitate the learning of tone contrasts. However, Wang and colleagues  
167 (1999, 2003) used only HVPT. Following the results of Logan et al. (1991, 1993) there is an  
168 interest in exploring whether high variability training has an advantage over low variability  
169 training. The first study to investigate this for the training of lexical tone was conducted by  
170 Perrachione et al. (2011). They trained native American English speakers with no previous  
171 knowledge of Mandarin (or any other tonal language), using English monosyllabic pseudowords  
172 combined with Mandarin tones 1 2, and 4. The training task used either low variability (one

173 speaker) or high variability (four speaker) input. The pseudowords were associated with concrete  
174 objects displayed in pictures. During the training, participants matched the sound they heard with  
175 one of three pictures presented, where the three words associated with these pictures were  
176 minimal trios that differed only in tone. They received feedback on a trial-by-trial basis.  
177 Learning was tested using a version of the training task with new talkers (and with feedback  
178 removed). Importantly, Perrachione et al. (2011) were also interested in the role of individual  
179 differences in learning. Therefore, in addition to the key tests of the training materials, they also  
180 determined participants' baseline ability to perceive the tone contrasts using a *Pitch Contour*  
181 *Perception Test (PCPT)*. In this task, participants heard a vowel produced with either Mandarin  
182 tone 1, 2 or 4 whilst viewing pictures of the three standard diacritics, and were asked to pick out  
183 the picture of the arrow that corresponded to the tone. Based on performance in this task before  
184 training, the researchers grouped participants into high and low aptitude groups. The key finding  
185 of this study was that while the low variability group outperformed the high variability group  
186 during training (presumably due to accommodation to a repeated speaker through the task), there  
187 were no differences between the high and low variability groups during test, even though test  
188 items involved novel speakers and thus probed generalization. Critically however, there was an  
189 interaction between individuals' aptitude categorization (as defined by the PCPT) and the type of  
190 variability training: Only the participants with high aptitude benefited from high variability  
191 training, while those with low aptitude actually benefited more from low variability training.

192 Another training study by Sadakata and McQueen (2014) also explored the relationship  
193 between input variability and individual aptitude in lexical tone training, though using rather  
194 different training and testing materials. They trained native Dutch speakers (with no prior  
195 knowledge of Mandarin or any other tonal language) using naturally produced bisyllabic

196 Mandarin pseudowords. The two syllables in each word either had Tone 2 followed by Tone 1,  
197 or Tone 3 followed by Tone 1, and each tone pair was randomly assigned one of two numeric  
198 labels (1, 2 - so for example for one participant Tone 2-Tone 1 was labelled “1”, Tone 3-Tone 1  
199 was labelled “2”). During the training task, participants were asked to identify the tone pair type  
200 of each stimulus by choosing the correct numeric label (e.g. hear /pasa/ with Tone 2-Tone 1,  
201 correct response is 2). Thus, in contrast to the study by Perrachione et al. (2011), participants did  
202 not need to learn the meaning of each word. Input variability was manipulated, with three levels  
203 (low/medium/high). In contrast to the work by Perrachione et al., where the high variability and  
204 low variability conditions differed only in terms of the number of speakers, in this study  
205 variability was increased both by including more speakers and more items (pseudowords). The  
206 test session used a similar design to the training sessions but included a 3AFC test (to prevent  
207 ceiling effect, a new untrained tone pair [Tone 1 – Tone 1], was included alongside the trained  
208 contrasts and assigned a new numeric label (“3”).

209         As in the study by Perrachione et al. (2011), Sadakata and McQueen (2014) also tested  
210 individual aptitude but with a different method. They employed a categorization task using  
211 stimuli from a six step Tone 2 to Tone 3 continua (created using natural productions of the two  
212 tones with the Mandarin vowel /a/ as endpoints and linearly interpolating between these  
213 endpoints). Participants were asked to identify if the sound they heard was more like Tone 2 or  
214 Tone 3 and a categorization slope was obtained for each participant, providing a measure of their  
215 ability to discriminate this contrast (which is generally found to be the most challenging tone  
216 contrast for L2 learners of Mandarin). Participants were grouped according to their slopes, and as  
217 in Perrachione et al., this grouping was entered as a factor in the analyses of the main test of  
218 learning. The results were similar to those of Perrachione et al.: there was no group level benefit

219 of high variability training but instead an interaction between individual aptitude and variability  
220 condition, which was due to the fact that only participants with high aptitude benefited from high  
221 variability training, while those with lower aptitude actually benefitted more from low variability  
222 training. There was also no interaction between aptitude and variability condition in the tests of  
223 generalization to new speakers or items.

224         The results of these studies thus provide mutually corroborating evidence – using  
225 somewhat different training and testing methods - that the ability to learn from high variability  
226 input is dependent on learner aptitude. Perrachione et al. (2011) suggest that one reason why low  
227 aptitude participants may struggle with multi-speaker input is that the speakers were intermixed  
228 during training: This requires trial-by-trial adaption to each speaker, which was not required in  
229 the corresponding single speaker low variability conditions. This may place a burden on learners  
230 (see Nusbaum & Morin, 1992; Mattys & Wiget, 2011 for evidence that intermixed multi-speaker  
231 stimuli are difficult even for L1 processing and that this interacts with constraints on working  
232 memory and attention). To test this, Perrachione et al. included a second experiment in which  
233 items from each speaker were presented in separate blocks (as is more common in high  
234 variability phonetic training). This improved performance with trained items compared with  
235 unblocked training for low aptitude learners only, confirming the hypothesis that switching  
236 between speakers interferes with learning for low aptitude learners. On the other hand, Sadaka  
237 and McQueen (2014) employed a blocked presentation in their high variability condition, so that  
238 trial-by-trial inconsistency *cannot* explain the greater difficulty of low aptitude learners in this  
239 study.

240

### 241 **1.3     The Current Study**

242           The finding that learning from multiple voices is more or less effective for different  
243 groups of learners may have implications for those interested in designing training tools for  
244 educational purposes. The fact that the effect has been found using quite different methods is  
245 encouraging. Here we further probe this finding in a new paradigm in which naive participants  
246 are trained using natural, meaningful stimuli from Mandarin Chinese. The current study serves as  
247 a partial replication and extension of the two previous studies by Perrachione et al. (2011) and  
248 Sadakata and McQueen (2014).

249           There are three important points to note with regards to our methodology. First, we  
250 trained participants on real Mandarin words produced by native speakers. This stands in contrast  
251 to previous studies which have trained participants only on pseudowords: Perrachione et al.  
252 (2011) used Mandarin tones with English pseudowords, whilst Sadakata and McQueen (2014)  
253 used Mandarin pseudowords. Second, while previous studies have trained participants on only  
254 three of the four tones, we trained participants on all four Mandarin tones (six tone contrasts)  
255 given that learners of Mandarin will need to learn the complete set. Thirdly, we embedded tone  
256 learning in a vocabulary learning task. This contrasts with the procedure used by Sadakata and  
257 McQueen, where participants were trained to map tonal categories onto (arbitrary) numbers, as  
258 well as with other HVPT studies in which participants were trained to map phonetic categories to  
259 orthographic categories (e.g. “r” and “l”, Logan et al. 1993). However the procedure is in line  
260 with that used by Perrachione et al. (described above), where participants were trained to  
261 associate pseudowords containing tonal information with pictures of common objects such as  
262 table, bus, or phone. Learning both tones and lexical items simultaneously more closely  
263 resembles real world L2 learning situations.

264           The key manipulation in the current study was the amount and type of variability that  
265 occurred during training. Following Perrachione et al. (2011), we compared training given to  
266 different groups of learners: low variability training (one speaker), high variability training (four  
267 speakers intermixed within each training session) and high variability blocking training (four  
268 speakers each presented in separate blocks). We predicted that the difficulty of high variability  
269 input for lower aptitude participants would be greater in the unblocked condition, thus potentially  
270 increasing the likelihood of seeing the predicted interaction between variability and learner  
271 aptitude. On the other hand, blocked input is more usual of HVPT (e.g. Logan et al. 1991;  
272 Iverson, Hazan & Bannister, 2005) and may increase the possibility of seeing any benefits of  
273 speaker variability on generalization.

274           We used two perceptual tasks designed to tap individual aptitude. These were adapted  
275 from those used Perrachione et al. (2011) and Sadakata and McQueen (2014). However, while  
276 the previous studies grouped participants into one of two categories (high aptitude *vs.* low  
277 aptitude) based on the aptitude score, in current study they were used as continuous measures  
278 (allowing us to avoid assigning an arbitrary “cut off” for high *vs.* low aptitude groups, and the  
279 loss of information which occurs when an underlying continuous variable is turned into a binary  
280 measure). Note that the statistical approach used in this paper (logistic and linear mixed effect  
281 models) allowed us to include continuous predictors and look at their interactions with other  
282 factors.

283           We also included several measures of learning. The three interval oddity task required  
284 participants to pick out the “different word” after hearing three words spoken aloud. The three  
285 words were minimal triplets but with only two tone used (e.g. *bā*, Tone 1; *bā*, Tone 1; *bà*, Tone  
286 4). Both speaker novelty and item novelty were manipulated. The word repetition task, in which

287 participants repeated spoken Mandarin words, provided a test of production which could be  
288 conducted both pre and post-test. Item novelty was again manipulated. In the post-test session  
289 only, we included two additional tests: a picture identification test and a picture naming task. The  
290 picture identification test was similar in form to the training session (2AFC picture  
291 identification), however new speakers were used in order to test speaker generalization. The  
292 picture naming task required participants to name the pictures used in training in Mandarin. Note  
293 that last two tasks test both the ability to perceive/produce the tone distinctions in Mandarin, but  
294 also to link these to meaning, potentially tapping more directly in to mechanisms relevant to  
295 word learning.

296 In sum, the following experiment assessed whether individuals' aptitude would interact  
297 with high/low variability training. It used real Mandarin stimuli with all four Mandarin tones  
298 embedded in a vocabulary learning task, and included tests of both perception and production.  
299

## 300 **2 Method**

### 301 **2.1 Participants**

302 Sixty adults recruited from UCL Psychology Subject Pool participated in the experiment,  
303 twenty in each of the three conditions (low variability, high variability, high variability blocking).  
304 Participant information is summarised in *Table 1*. There was no difference between these groups  
305 in age,  $F(2, 57) = 1.95, p = .15$ . Participants had no known hearing, speech, or language  
306 impairments. Written consent was obtained from participants prior to the first session. Each  
307 participant was paid £45 at the end of the study.

308 All participants except three were native speakers of English. Of these three, one participant  
309 (low variability condition) was a native bilingual of English and Hindi, one participant (high

310 variability condition) was a native French speaker, and one participant (high variability condition)  
311 was a native Finnish speaker. Critically none had any prior experience of Mandarin Chinese or  
312 any other tonal language. On average, participants learned 2.4 ( $SD = 0.8$ ) languages and the  
313 average age for starting to learn the first L2 was 12.6 ( $SD = 1.3$ ).

314

## 315 2.2 *Stimuli*

### 316 2.2.1 *Stimuli used in Training and in the Picture Identification, Three Interval Oddity, Word* 317 *Repetition and Picture Naming Tests*

318 These stimuli consisted of 36 minimal pairs of Mandarin words (6 minimal pairs for each  
319 of the six tone contrasts for each of the four Mandarin tones). The words in each pair contained  
320 the same phonemes, differing only in tones (e.g. *māo*, Tone 1 [*cat*] vs. *mào*, Tone 4 [*hat*]). The  
321 words were chosen to be picturable and to start with a wide range of phonemes (see Appendix  
322 A). In order to examine generalization across items, half of the word pairs (3 per tone contrast)  
323 were designated "trained" words and used in both training and testing: the other half were  
324 designated "untrained" words and were encountered only at test.

325 The full set of 72 Mandarin words was recorded by two groups of native Mandarin  
326 speakers using a Sony PCM-M10 handheld digital audio recorder. The first group was made up  
327 of three female speakers and two male speakers, (F1, F2, F3, M1, M2). These stimuli were used  
328 in the training, word repetition and picture identification tasks. The second group consisted of  
329 three new female speakers and two new male speakers (FN1, FN2, FN3, MN1, MN2). These  
330 stimuli were used in the Three interval oddity task (making all new speakers in that task). Table  
331 2 summarises how speakers were assigned to each task.



332 In the low variability condition only one speaker (Trained voice 1) was used in training,  
333 and this same speaker was also used as the test voice in the Word Repetition test and for trained  
334 test items in the Picture Identification test. In the high variability condition, four speakers were  
335 used in training. Only one of these speakers (Trained voice 1) was used in the Word Repetition  
336 test and for trained items in the Picture Identification test (the same speaker across both tests). In  
337 both conditions, a further speaker (New voice 1) was assigned to the untrained test items in the  
338 Picture Identification test. The assignment of speakers was rotated across participants, resulting  
339 in 5 counterbalanced versions of each condition (see Table 2). This ensured that any difference  
340 found between the low and high variability conditions, and between trained and new voices,  
341 were not due to idiosyncratic difference between voices. There was no counterbalancing of  
342 speaker in other tasks.

343 All words were edited into separate sound files, and peak amplitude was normalised  
344 using Audacity (Audacity team, 2015, <http://audacity.sourceforge.net/>). Any background noise  
345 was also removed. All recordings were perceptually natural and highly distinguishable as judged  
346 by native Chinese speakers. Clipart pictures of the 72 words were selected from free online  
347 clipart databases.

#### 348 2.2.2 *Stimuli used in the Aptitude Tests:*

349 Pitch Contour Perception Test: Six Mandarin vowels (/a/, /o/, /e/, /i/, /u/, /y/) were  
350 repeated in the four Mandarin tones by two male and two female native Mandarin speakers  
351 (MN1, MN2, FN1, FN2 from taker set 2) making 96 stimuli in total. Stimuli were identical  
352 across conditions and participants.

353 Categorization of Synthesized Tonal Continua: Natural endpoints were chosen from a  
354 native Mandarin male speaker producing the word 'wan' with both Tone 2 and Tone 3. A neutral

355 vowel was also recorded by a native male English speaker producing the ‘father vowel’ /a/. This  
356 vowel was edited slightly to remove portions containing creaky voice at the end.

357         The three syllables (*wan* [Tone 2], *wan* [Tone 3], /a/) were then manipulated in Praat  
358 (Boersma & Weenink, 2015). All three syllables were normalized to be approximately 260 ms  
359 long using the PSOLA method. The neutral vowel was manipulated to have a flat pitch (148 Hz)  
360 and a flat intensity contour (75dB). The pitch contours of the two natural endpoints were  
361 extracted and a 6-step pitch continuum (Step 1: Tone 2, Step 6: Tone3) was generated by linearly  
362 interpolating between the endpoints. These six pitch contours were then each superimposed on a  
363 copy of the neutral vowel using the PSOLA method. Stimuli were identical across participants  
364 and conditions.

365

### 366 **2.3 Procedure**

367         The experiment involved three stages (see *Figure 2.3*): Pre-test (session 1), training  
368 (sessions 2-7), and post-test (session 8). Participants were required to complete all eight sessions  
369 within two weeks, with the constraint of one session per day at most. The majority of sessions  
370 took place in a quiet, soundproof testing room in Chandler House, UCL. The remaining sessions  
371 took place in a quiet room in a student house.

372         Participants were given a brief introduction about the aim of the study and told that they  
373 were going to learn some Mandarin tones and words. They were explicitly told that Mandarin  
374 has four tones (flat, rising, dipping and falling) and that the tonal differences were used to  
375 distinguish meanings. The experiment ran on a on a Dell Alienware 14R laptop with a 14-inch  
376 screen. The experiment software was built using a custom-built software package developed at  
377 the University of Rochester.

378           The specific instructions for each task were displayed on- screen before the task started.  
379   After each task, participants had the opportunity to take a 1-minute break. The tasks completed  
380   in each session are listed in *Figure 2.3* and described in more detail below. Note that the PCPT  
381   and CSTC were carried out at the beginning of the first session as they provided the measure of  
382   individual aptitude prior to exposure to any Mandarin stimuli. There was no time limit for  
383   making responses in any of the tasks. Participants wore a pair of HD 201 Sennheiser headphones  
384   throughout the experiment.

385

### 386 *2.3.1 The Pitch Contour Perception Test*

387           This test was based on the work of Wong and Perrachione (2007). Participants heard a  
388   tone (e.g. /a/ [Tone 1]), while viewing pictures of four arrows indicating the different pitch  
389   contours on the screen. Participants clicked on the arrow that they thought matched the tone  
390   heard. No feedback was provided. There were 96 stimuli in total (4 speakers \* 4 tones \* 6  
391   vowels). Participants completed this task twice, at both pre- and post-test. The main purpose of  
392   this task was to provide a measure of individual differences in tone perception prior to training,  
393   following Perrachione et al. (2011). Although Perrachione et al. only conducted this task at pre-  
394   test, for consistency with the CSTC (described below) we also repeated the test at post-test and  
395   conducted analyses to identify whether performance on this task was itself improved as a result  
396   of training (see Section 3.3.2).

397

### 398 *2.3.2 Categorization of Synthesized Tonal Continua*

399           This test was based on Sadakata and McQueen (2014). Participants first practiced  
400   listening to Tone 2 and Tone 3. They heard the tone while viewing the corresponding picture of  
401   an arrow. Each tone was repeated 10 times. Then, for each test trial, participants were asked to

402 decide if the sound they heard was closer to Tone 2 or Tone 3 by clicking on the corresponding  
403 arrow. No feedback was provided. The speech continua consisted of 6 steps (Step 1: Tone 2,  
404 Step 6: Tone 3). Each of the six steps was repeated 10 times per block. Participants completed  
405 two blocks, with an optional 1 minute break in the middle, resulting in 120 trials in total. The  
406 main purpose of this task was to provide a measure of individual differences in tone perception  
407 prior to training, following Sadakata and McQueen (2014). In line with their procedure,  
408 participants completed the task both before and after training and we conducted analyses to  
409 explore whether there was improvement from pre to post-test (see Section 3.2.1).

410

### 411 2.3.3 *Three Interval Oddity Test*

412 This task required subjects to identify the “different” stimulus from a choice of three  
413 Mandarin words. Each of the three words within a trial was spoken by a different speaker. Four  
414 speakers were used (3 female, 1 male). All speakers were untrained (i.e., not used during  
415 training; see *Table 2*). Each trial used one of the 36 minimal pairs from the main stimuli set (18  
416 trained pairs, 18 untrained pairs). Preliminary work suggested that trials differed in difficulty  
417 depending on whether the “different” stimulus was spoken by the single male speaker, or one of  
418 the three female speakers. We therefore ensured that there were equal numbers of the following  
419 trial types: (i) “Neutral” - all three words were spoken by female speakers (ii) “Easy” - the  
420 “different” word was spoken by a male speaker and the other two were spoken by female  
421 speakers; (iii) “Hard” - the “different” word was spoken by a female speaker and the other two  
422 were spoken by one male speaker and one female speaker. Each of the words in the minimal pair  
423 was used once as the target (“different”) word, making 72 trials in total.

424 During the task, three frogs were displayed on the screen. Participants heard three words  
425 (played with ISIs of 200ms) and indicated which word was the odd one out by clicking on the  
426 appropriate frog, which could be in any of the three positions. They could not make their  
427 response until after all three words had been heard, at which point a red box containing the  
428 instruction “click on the frog that said the different word” appeared at the bottom of the screen.  
429 No feedback was given after each trial. Participants completed this task twice – once in the pre-  
430 test, and once in the post-test (see *Figure 2.3*).

431

#### 432 2.3.4 *Word Repetition Test*

433 All seventy-two Mandarin words from the main stimuli set were presented one at a time  
434 in a randomised order. They were always spoken by the same speaker and this speaker was also  
435 used in their training stimuli (Training voice 1; see *Table 2*). After each word, two seconds of  
436 white noise was played. Participants were instructed to listen carefully to the word and then to  
437 repeat the word aloud after the white noise. The white noise was included to make sure that  
438 participants had to encode the stimulus they were repeating, rather than relying on the  
439 phonological loop, which would be pure imitation (Flege, Takagi & Mann, 1995). Verbal  
440 responses were digitally recorded and were later transcribed and rated by native speakers of  
441 Mandarin (see Section 3.3.1.1). This task was completed once in the pre-test and once in the  
442 post-test.

443

#### 444 2.3.5 *English Introduction Task*

445 This task was included in case the meaning of some pictures were ambiguous (not all  
446 items were concrete nouns – e.g. “to paint”). Participants saw each of the 36 pictures from the

447 training set presented once each in random order and heard the corresponding English word. No  
448 response was recorded. Participants completed this task only once, at the end of the pre-test  
449 session.

450

### 451 2.3.6 Training Task

452 Participants completed the training task in Session 2-7. On each trial, participants heard a  
453 Mandarin word and selected one of two candidate pictures displayed on the computer screen.  
454 The two picture always belonged to the same minimal pair (see *Figure 2.3.6*). After selecting a  
455 picture, the participant was informed whether their answer was correct (a green happy face  
456 appeared) or incorrect (a red sad face appeared). If the correct choice was made, a picture of a  
457 coin also appeared in a box on the left-hand side of the screen, with the aim of motivating  
458 participants to try to earn more coins in each subsequent session of training. After that,  
459 everything but the correct picture was removed from the screen and the participant heard the  
460 correct word again. In the lower right corner of the screen a trial indicator of X/288 was  
461 displayed where X indicated the number of trials completed. This tool helped participants to  
462 keep track of their performance (see *Figure 2.3.6*).

463 There were 18 picture/word pairs used. Each word was used as the target word four  
464 times. Thus, each picture pair appeared eight times, resulting in 288 trials in total per session.  
465 Participants were assigned to one of the following condition: low variability, high variability and  
466 high variability blocking (with the assignment of speakers counterbalanced – see *Table 2*). Each  
467 session lasted for approximately 30 minutes.

468 In the low variability condition, only *one* speaker was used. In the high variability  
469 condition, *four* speakers were used. For these two condition, all 288 trials were randomized so

470 there was no fixed order of speaker. For each participant, each of their six training sessions was  
471 identical. In the high variability blocking condition, the stimuli were the same as those in the  
472 high variability condition. However, from Day 1 to Day 4 of training (i.e., Session 2-5), only one  
473 speaker was involved on each day's training session, with the trained speaker that was used in  
474 the test tasks (e.g. F1 for Version 1) always occurring on Day 3 (i.e., Session 4). On Days 5 and  
475 6 of training (i.e., Sessions 6 and 7), participants heard all four speakers, each in a separate  
476 block, each word was repeated twice in each voice on these days. The trained speaker used in the  
477 test tasks always occurred in the third block. After each block, the number of coins they had  
478 earned so far was displayed on the screen. For each participant, the structure of the training task  
479 was identical on Days 5 and 6.

480

#### 481 2.3.7 *Picture Identification Test*

482 This task was the same as the training task with the following changes. Firstly, each word  
483 was only repeated twice, once by a trained speaker (Trained voice 1) and once by an untrained  
484 speaker (New voice 1), making 72 trials in total. Secondly, no feedback was given. This task was  
485 completed only in the post-test.

#### 486 2.3.8 *Picture Naming Test*

487 All 36 pictures from the training words were presented in a randomised order.  
488 Participants were instructed to try to name the picture using the appropriate Mandarin word.  
489 Verbal responses were recorded and were later transcribed and rated by native Mandarin  
490 speakers (see Section 3.5.2). This task was completed only in the post-test.

### 491 2.3.9 Questionnaires

492 Participants completed a language background questionnaire after the experiment.  
493 Participants were asked to list all the places they had lived for more than 3 months and any  
494 languages that they had learned. For each language the participant was asked to state: (a) how  
495 long they learned the language for and their starting age; (b) to rate their own current proficiency  
496 of the language.

497

## 498 3 Results and Discussion

### 499 3.1 Statistical Approach

500 Three different sets of analyses are reported. First, we conducted the analysis on two  
501 individual measures: CSTC (Section 3.2.1) and PCPT (Section 3.2.3). The primary aim of these  
502 analyses was to ensure that the three groups did not differ at pre-test, however we also looked for  
503 possible differences at post-test. Second, separate analyses are reported on: data from the tests  
504 administered pre- and post- training (i.e. word repetition task (Section 3.3.1) and Three Interval  
505 oddity task (Section 3.3.2), the data collected during training (Section 3.4) and the data from the  
506 two tasks administered only at post-test (i.e. the picture identification task (Section 3.5.1) and  
507 picture naming task (Section 3.5.2). These analyses, explore the effects of our experimentally  
508 manipulated conditions on the various measures of Mandarin tone learning. Third, analyses were  
509 conducted exploring the role of aptitude in each of these tasks (Section 3.6). Specifically, we  
510 wanted to see whether aptitude interacted with *variability-condition* in predicting the benefits of  
511 training, in line with the predictions of previous research (Perrachione et al., 2011; Sadakata &  
512 McQueen, 2014).



513 Except where stated, analyses used logistic mixed effect models (LMEs; Baayen,  
514 Davidson, & Bates, 2008; Jaeger, 2008; Quené & van den Bergh, 2008) using the package lme4  
515 (Bates, Maechler, & Bolker, 2013) for the R computing environment (R Development Core  
516 Team, 2010). LMEs allow binary data to be analysed with logistic models rather than as  
517 proportions, as recommended by Jaeger (2008). In each of the analyses, the factor *variability-*  
518 *condition* has three levels (low variability [LV], high variability [HV], and high variability  
519 blocking [HVB]) which we coded into two contrasts with LV as the baseline (LV versus HV, LV  
520 versus HVB). An exception to this is the training data, where a model containing all three  
521 conditions would not converge and we took a different approach, as described in Section 3.4. We  
522 also included the interactions between these contrasts and the other factors. We used centred  
523 coding which ensued that other effects were evaluated as averaged over all three levels of  
524 *variability-condition* (rather than the reference level of LV<sup>1</sup>). Similarly, in the Three Interval  
525 Oddity, we included a *trial-type* factor (to control for the fact that participants were likely to find  
526 some trial types easier than others) – this had three levels ((i) “Neutral” - all three words were  
527 spoken by female speakers (ii) “Easy” - the “different” word was spoken by the one male  
528 speaker (iii) “Hard” - the “different” word was spoken by one of the two female speakers) and  
529 for this we included contrasts with neutral (“neutral versus easy” and “neutral versus hard”)  
530 again using centered coding. In order to perform the analysis comparing pre- and post-test  
531 performance, *test-session* was coded as a factor with two levels (pre-test/post-test) with “pre-  
532 test” set as the reference level. This allowed us to look at the (accidental) possible differences  
533 between the experimental conditions at the pre-test stage, as well as whether post-test  
534 performance differed from this baseline. All other predictors, including both discrete factor

---

<sup>1</sup> This differs from the default coding of contrasts in the lme4 package. It was achieved by replacing the three-way factor “condition” with two centred dummy variables and using the main fixed effects from the output of this model.

535 codings with two levels (*item-novelty* in the Word Repetition and Three Interval Oddity tasks,  
536 and *voice-novelty* in the Picture Identification task) and numeric predictors (*training-session*) in  
537 the Training data analyses and the individual difference measures in the models reported in  
538 Section 3.7), were centred to reduce the effects of collinearity between main effects and  
539 interactions, and in order that main effects were evaluated as the average effects over all levels of  
540 the other predictors (rather than at a specified reference level for each factor). We automatically  
541 put experimentally manipulated variables and all of their interactions into the model, without  
542 using model selection (except for “*trial-type*” in the Three Interval Oddity task which works as a  
543 control factor and for this factor we only used its main effect and the interaction with *test-*  
544 *session*). However, we did not inspect the models for all main effects and interactions. Instead,  
545 we report statistics which were necessary to look for accidental differences at pre-test, and those  
546 related to our hypotheses. We aimed to examine whether the training improves participants’  
547 performance on both new items and new voices and whether such improvement was modulated  
548 by their individual aptitudes. Participant is included as a random effect and a full random slope  
549 structure was used (i.e., by-subject slopes for all experimentally manipulated within-subject  
550 effects (*test-session*, *voice-novelty*, *item-novelty*) and interactions, as recommended by Barr,  
551 Levy, Scheepers, and Tily, 2013. In some cases the models did not converge and in those cases  
552 correlations between random slopes were removed. Models converged with Bound Optimization  
553 by Quadratic Approximation (BOBYQA optimization; Powell, 2009). R scripts showing full  
554 model details can be found here:

555 [https://osf.io/wdh8a/?view\\_only=d1557462138447ffbafaf7a59662df8](https://osf.io/wdh8a/?view_only=d1557462138447ffbafaf7a59662df8).

556

### 557 3.2 *Individual Aptitude Tasks*

558 3.2.1 *Categorisation of Synthesized Tonal Continua*

559 We estimated individual's performance on the CSTC task following Sadakata and  
560 McQueen (2014). We used the Logistic Curve Fit function in SPSS to calculate a slope  
561 coefficient for each participant (Joanisse, Manis, Keating & Seidenberg, 2000). The slope  
562 (standardized  $\beta$ ) indicates individual differences in tone perception. The smaller the slope, the  
563 better the performance. According to Sadakata and McQueen, the data of participants with a  
564 slope measure greater than 1.2 were removed from the analysis. Using this threshold 43 out of 60  
565 participants failed the threshold. This is consistent with the observation that most of the  
566 participants were not able to consistently categorize the endpoints of the continua, indicating that  
567 this was not a good test of aptitude. We do not report further analyses with this aptitude variable  
568 however they can be found in the supplemental materials  
569 ([https://osf.io/wdh8a/?view\\_only=d1557462138447ffbaafaf7a59662df8](https://osf.io/wdh8a/?view_only=d1557462138447ffbaafaf7a59662df8)).

570 3.2.2 *The Pitch Contour Perception Test*

571 The predicted variable was whether a correct response was given (1/0) on each trial. The  
572 predictors were the contrasts between conditions (LV versus HV; LV versus HVB) and *test-*  
573 *session* (pre-test, post-test). (Note - average accuracy in each condition is also included in the  
574 table of participant details; Table 1, section 2.1). There was no significant difference between the  
575 LV and HV groups at pre-test ( $\beta = -0.35$ ,  $SE = 0.26$ ,  $z = -1.38$ ,  $p = 0.17$ ) or between the LV and  
576 HVB groups ( $\beta = 0.17$ ,  $SE = 0.26$ ,  $z = 0.66$ ,  $p = 0.51$ ) on this measure. Participants showed  
577 significant improvement after training ( $\beta = 0.21$ ,  $SE = 0.05$ ,  $z = 4.13$ ,  $p < 0.001$ ).

578 In sum, for this measure of perceptual ability our three participant groups did not differ in  
579 their performance and the groups showed equivalent improvement from pre- to post- test. Given

580 that this measure is affected by training, we used participants scores at pre-test as our measure of  
581 individual differences in the analyses reported in Section 3.6.

582

### 583 **3.3 Tests Administered Pre- and Post- Training**

#### 584 *3.3.1 Word Repetition*

##### 585 *3.3.1.1 Coding and inter-rater reliability analyses*

586 The same methods were used for both production tests – i.e. the Word Repetition test  
587 (pre- and post-) and the Picture naming task (post-test only). The files were combined into a  
588 single set, along with the 360 stimuli which were used in the experiment (and which were  
589 produced by native Mandarin speakers). The latter items were included in order to examine  
590 whether the raters were reliable. All stimuli were rated by two raters: Rater 1 was the first author  
591 and Rater 2 was recruited from the UCL MA Linguistics program and was naïve to the purposes  
592 of the experiment. Raters were presented with recordings in blocks in a random sequence (blind  
593 to test-type, condition, whether the stimulus was from pre-test or post-test and whether it was  
594 produced by a participant or was one of the experimental stimuli). For each item, raters were  
595 asked to (i) identify the tone, (ii) give a rating quantifying how native-like they thought the  
596 pronunciation was compared (1-7 with 1 as not recognizable and 7 as native speaker level), and  
597 (iii) transcribe the pinyin (segmental pronunciation) produced by the participants.

598 Three measurements were taken from the production tasks: mean accuracy of tone  
599 identification (Tone accuracy), mean tone rating (Tone rating) and mean accuracy of production  
600 of the pinyin (derived by coding each production as correct (1= the entire string is correct) or  
601 incorrect (0 = at least one error in the pinyin)). As a first test of rater reliability, performance  
602 with the native speaker stimuli was examined– these were near ceiling: Rater 1: Tone accuracy =

603 98%, Tone rating = 6.7, Pinyin accuracy = 80%; Rater 2: Tone accuracy = 87%, Tone rating =  
604 6.5, Pinyin accuracy = 80%).

605 Furthermore, for the remaining data (i.e. the experimental data) inter-rater reliability was  
606 examined for both measures for the two production tasks. For the binary measures (Tone  
607 accuracy and Pinyin accuracy), kappa statistics were calculated using the “fmsb” package in R  
608 (Cohen, 2014). For the word repetition data, for Tone accuracy  $kappa = 0.43$  (“moderate  
609 agreement”), and for Pinyin accuracy  $kappa = 0.33$  (“fair agreement”; Landis & Koch, 1977). For  
610 the Picture Identification test, for Tone accuracy  $kappa = 0.68$  (“substantial agreement”) and for  
611 Pinyin accuracy  $kappa = 0.54$  (“moderate agreement”); For the Tone rating, the package “irr” in  
612 R was used to access the intra-class correlation (McGraw & Wong, 1996) based on an average-  
613 measures, consistency, two-way mixed-effects model. For Word Repetition,  $ICC = 0.28$  and for  
614 Picture Identification  $ICC = 0.44$ ; according to Cicchetti (1994), values less than .40 are regarded  
615 as “poor”. Given this, we do not include analyses with Tone Rating as the dependent variable  
616 (though these data are included in the data set  
617 [https://osf.io/wdh8a/?view\\_only=d1557462138447ffbaafaf7a59662df8](https://osf.io/wdh8a/?view_only=d1557462138447ffbaafaf7a59662df8)). All of the analyses  
618 presented in Sections 3.3.1 and 3.5.2 were based on Rater 2 (the naive rater).

619

### 620 3.3.1.2 Tone accuracy

621 The predicted variable was whether a correct response was given (1/0) on each trial (as  
622 identified by the coder). The predictors were *test-session* (pre-test, post-test), *variability-*  
623 *condition* (LV versus HV, LV versus HVB) and *item-novelty* (trained, untrained). The mean  
624 accuracy, split by test session and training condition, is shown in *Figure 3.3.1.2*.

625 At pre-test, there was no significant difference between the LV and the HV group ( $\beta = -$   
626  $0.09$ ,  $SE = 0.20$ ,  $z = -0.46$ ,  $p = .65$ ) nor between the LV and the HVB group ( $\beta = 0.05$ ,  $SE = 0.20$ ,  
627  $z = 0.27$ ,  $p = .79$ ), suggesting the groups started at a similar level. There was also no difference  
628 between trained and untrained words at pre-test ( $\beta = 0.06$ ,  $SE = 0.11$ ,  $z = 0.51$ ,  $p = 0.61$ ).

629 Across the three groups, participants' performance increased significantly after training  
630 ( $M_{pre} = 0.70$ ,  $SD_{pre} = 0.14$ ,  $M_{post} = 0.76$ ,  $SD_{post} = 0.14$ ,  $\beta = 0.37$ ,  $SE = 0.13$ ,  $z = 2.90$ ,  $p <$   
631  $.01$ ). There was no significant difference in the improvement for trained and untrained items  
632 (word-type by test-session interaction:  $\beta = 0.08$ ,  $SE = 0.16$ ,  $z = 0.49$ ,  $p = .63$ ). The interactions  
633 between the variability contrasts and test-session were not significant (LV versus HV:  $\beta = -0.20$ ,  
634  $SE = 0.31$ ,  $z = -0.65$ ,  $p = .52$ ; LV versus HVB:  $\beta = -0.31$ ,  $SE = 0.31$ ,  $z = -0.99$ ,  $p = .32$ ), and they  
635 were not qualified by any higher level interactions with *item-novelty* (LV versus HV:  $\beta = 0.01$ ,  
636  $SE = 0.38$ ,  $z = 0.02$ ,  $p = .99$ ; LV versus HVB:  $\beta = -0.30$ ,  $SE = 0.38$ ,  $z = -0.79$ ,  $p = .44$ ).

### 637 3.3.1.3 Pinyin accuracy

638 The predicted variable was whether the participants produced the correct string of  
639 phonemes (1/0) in each trial (as determined by the rater). The predictors were *test-session* (pre-  
640 test, post-test), *variability-condition* (LV versus HV, LV versus HVB) and *item-novelty* (trained,  
641 untrained). Mean pinyin accuracy is displayed in *Figure 3.3.1.3*.

642 At pre-test, there was no significant difference between the LV and the HV group ( $\beta = -$   
643  $0.05$ ,  $SE = 0.13$ ,  $z = -0.41$ ,  $p = .68$ ) nor between the LV and the HVB group ( $\beta = -0.08$ ,  $SE =$   
644  $0.13$ ,  $z = -0.60$ ,  $p = .55$ ), suggesting that the groups started at a similar level. However,  
645 participants did better on untrained words than trained words at pre-test ( $\beta = 0.25$ ,  $SE = 0.09$ ,  $z =$   
646  $2.82$ ,  $p < .01$ ), suggesting potential accidental differences in these items. Participants showed no  
647 improvement after training ( $M_{pre} = 0.54$ ,  $SD_{pre} = 0.13$ ,  $M_{post} = 0.55$ ,  $SD_{post} = 0.13$ ,  $\beta = 0.07$ ,

648  $SE = 0.09, z = 0.81, p = .42$ ). In addition, there was no evidence of different improvements for  
649 different variability conditions (*test-session* by LV versus HV:  $\beta = -0.02, SE = 0.22, z = -0.09, p$   
650  $= .93$ ; *test-session* by LV versus HVB:  $\beta = -0.27, SE = 0.22, z = -1.24, p = .22$ ) or any interaction  
651 between *variability condition, test-session* and *item-novelty* (LV versus HV:  $\beta = 0.07, SE = 0.31,$   
652  $z = 0.23, p = .82$ ; LV versus HVB:  $\beta = -0.41, SE = 0.31, z = -1.33, p = .18$ ).

653

### 654 3.3.2 Three Interval Oddity Task

655 The predicted variable was whether a correct response was given (1/0) on each trial. The  
656 predictors were *test-session* (pre-test, post-test), *variability-condition* (LV versus HV, LV versus  
657 HVB), *trial-type* (neutral versus easy, neutral versus hard) and *item-novelty* (trained item,  
658 untrained item). The mean accuracy is displayed in *Figure 3.3.2*.

659 At pre-test, there was no significant difference between the LV and the HV group ( $\beta = -$   
660  $0.002, SE = 0.14, z = -0.01, p = .99$ ) nor between the LV and the HVB group ( $\beta = 0.12, SE =$   
661  $0.14, z = 0.86, p = .39$ ), suggesting the groups started at a similar level. However, performance  
662 with the items classified as “untrained” was significantly greater at pre-test ( $\beta = -0.31, SE = 0.06,$   
663  $z = -4.95, p < 0.01$ ), suggesting accidental differences between items. As expected, at pre-test  
664 participants performed significantly better on “easy” trials (where the target speaker had a  
665 different gender) than “neutral” trials (where all three speakers had the same gender),  $\beta = 0.40,$   
666  $SE = 0.08, z = 5.09, p < 0.01$ ; and “neutral” trials were marginally easier than “hard” trials  
667 (where one of the foil speakers had the odd gender out),  $\beta = -0.14, SE = 0.08, z = -1.81, p = 0.07$ .

668 Overall, participants’ performance increased significantly after training ( $M_{pre} = 0.59,$   
669  $SD_{pre} = 0.21, M_{post} = 0.66, SD_{post} = 0.19, \beta = 0.31, SE = 0.05, z = 6.54, p < .001$ ). Critically,  
670 there was no reliable interaction between *test-session* and *item-novelty* ( $\beta = 0.14, SE = 0.09, z =$

671 1.49,  $p = .14$ ), suggesting no evidence that training had a greater effect for trained words than for  
672 novel words. There was also no interaction with *test-session* for either the contrast between the  
673 LV versus the HV conditions ( $\beta = -0.01$ ,  $SE = 0.12$ ,  $z = -0.12$ ,  $p = .90$ ) or the contrast between  
674 the LV versus the HVB conditions ( $\beta = 0.01$ ,  $SE = 0.12$ ,  $z = 0.11$ ,  $p = .91$ ) and no higher-level  
675 interactions. This suggests that the extent to which participants improved on this task between  
676 pre and post-test did not differ across *variability-conditions* or *item-novelty*.

677         Although not part of our key predictions, we also looked to see if there was evidence that  
678 participants improved more with the easier or harder trials. In fact, the interaction between *test-*  
679 *session* and the contrast between “easy” and “neutral” was significant ( $\beta = -0.27$ ,  $SE = 0.11$ ,  $z = -$   
680  $2.39$ ,  $p = .02$ ) while the contrast between “neutral” and “hard” was not ( $\beta = 0.12$ ,  $SE = 0.11$ ,  $z = -$   
681  $1.06$ ,  $p = .29$ ). This was due to the fact that there was improvement for “neutral” (Mpre = 0.57,  
682 SDpre = 0.14, Mpost = 0.65, SDpost = 0.15) and “hard” trials (Mpre = 0.54, SDpre = 0.16,  
683 Mpost = 0.65, SDpost = 0.15) but not for “easy” trials (Mpre = 0.66, SDpre = 0.16, Mpost =  
684 0.68, SDpost = 0.15).

685

### 686 3.3.3 Summary of Tests administered Pre-and Post-Training

687         The analysis of Word Repetition and Three Interval Oddity data showed that participants’  
688 performance increased significantly after training (except for Pinyin accuracy in Word  
689 Repetition) for both tasks. For the Pinyin accuracy measure in Word Repetition, and for the  
690 Three Interval Oddity task, there was a main effect of *item-novelty* at pre-test, suggesting that  
691 items designated to be “untrained” were accidentally easier than those designated as “trained”,  
692 but no interaction with *test-session* suggesting that training did not differentially affect  
693 improvement with trained and untrained items. However, the critical finding was that there was



694 no interaction between *test-session* and *variability-condition*, or between *test-session*, *item-*  
695 *novelty* and *variability-condition*, providing no evidence that the variability manipulation  
696 affected the extent of improvement in these tests.

697

### 698 3.4 Training Data

699 Here, a model containing data from all three conditions did not converge; however two  
700 separate models, one including the LV and HV conditions, and the other the LV and HVB  
701 conditions (with condition as a factor with two levels), did converge. In each case the predicted  
702 variable was whether a correct response was given (1/0) on each trial. The predictors were the  
703 numeric factor *training-session* (1→6) and the factor *variability-condition* which had two levels  
704 (model 1: LV versus HV; model 2, LV versus HVB). The mean accuracy is displayed in *Figure*  
705 3.4.

706

707 In both models, there was an effect of *training-session* (model 1:  $\beta = 0.49$ ,  $SE = 0.04$ ,  $z =$   
708  $11.52$ ,  $p < .001$ ; model 2:  $\beta = 0.53$ ,  $SE = 0.04$ ,  $z = 12.17$ ,  $p < .001$ ): Participants' performance  
709 increased significantly with training-sessions. Overall, the LV group performed better than both  
710 the HV group ( $\beta = -0.79$ ,  $SE = 0.16$ ,  $z = -5.03$ ,  $p < .001$ ) and the HVB group ( $\beta = -0.83$ ,  $SE =$   
711  $0.32$ ,  $z = -2.61$ ,  $p < .01$ ). However the LV versus HV contrast was also modulated by an  
712 interaction with *test-session* ( $\beta = -0.19$ ,  $SE = 0.04$ ,  $z = -4.59$ ,  $p < .001$ ), as was the LV versus  
713 HVB contrast ( $\beta = -0.35$ ,  $SE = 0.08$ ,  $z = -4.33$ ,  $p < .001$ ). From *Figure 3.4* it can be seen that the  
714 LV and the HVB group did not differ in the first session (i.e. where they get identical input) but  
715 the difference gradually increased over the next sessions. For the LV and the HV group, they

716 differed starting from the first session and this difference continued to increase throughout  
717 training.

718

### 719 3.4.1 Summary of training data

720 The analysis of training data revealed significant improvements for all three groups. The  
721 LV group performed better than the other two groups due to repetitive exposure to just one  
722 speaker throughout the six sessions. In the first session, the difference between the LV and the  
723 HVB groups was not significant. However, the difference between conditions increased over  
724 time for both LV-HVB and LV-HV contrasts.

725

## 726 3.5 Tests Administered at Post-Test Only

### 727 3.5.1 Picture Identification

728 The coding and reliability analyses for this data is described in section 3.3.1.1. The  
729 predicted variable was whether a correct response was given (1/0) on each trial. The predictors  
730 were the factor *voice-novelty* (Trained voice, New voice) and the factor *variability-condition*  
731 which had two contrasts (LV versus HV, LV versus HVB). The mean accuracy is displayed in  
732 *Figure 3.5.1.1*.

733 There was a main effect of *voice-novelty* ( $\beta = 1.07, SE = 0.16, z = 6.53, p < .001$ )  
734 reflecting higher performance in trials with trained voices. Participants in the LV group  
735 performed better than those in the HV group ( $\beta = -0.71, SE = 0.32, z = -2.23, p = .03$ ) but there  
736 was no significant difference between the LV and the HVB group ( $\beta = -0.14, SE = 0.32, z = -$   
737  $0.44, p = .66$ ). There was a significant interaction between *voice-novelty* and both the LV-HV  
738 contrast ( $\beta = -1.19, SE = 0.35, z = -3.43, p < .01$ ) and the LV-HVB contrast ( $\beta = -1.11, SE =$

739 0.36,  $z = -3.08$ ,  $p < .01$ ). Breaking this down by condition: for each condition there was  
740 significantly better performance with trained than new voices (LV:  $\beta = 1.83$ ,  $SE = 0.29$ ,  $z = 6.42$ ,  
741  $p < 0.001$ ; HV:  $\beta = 0.64$ ,  $SE = 0.23$ ,  $z = 2.86$ ,  $p < 0.01$ ; HVB:  $\beta = 0.73$ ,  $SE = 0.26$ ,  $z = 2.82$ ,  $p <$   
742  $0.01$ ). Breaking it down by voice-novelty: For new voices, neither of the contrasts between  
743 conditions was significant (LV versus HV:  $\beta = -0.12$ ,  $SE = 0.26$ ,  $z = -0.45$ ,  $p = 0.65$ ; LV versus  
744 HVB  $\beta = 0.41$ ,  $SE = 0.27$ ,  $z = 1.51$ ,  $p = 0.13$ ). For trained items, there was significantly higher  
745 performance in the LV than HV condition, but no difference between the LV and HVB  
746 conditions (LV versus HV:  $\beta = -1.30$ ,  $SE = 0.44$ ,  $z = -2.97$ ,  $p < 0.01$ ; LV versus HVB:  $\beta = -0.70$ ,  
747  $SE = 0.45$ ,  $z = -1.55$ ,  $p = 0.12$ ).

748

### 749 3.5.2 Picture Naming

750 These data used the same two measures as the Word Repetition data (see section 3.4.1),  
751 i.e. (i) tone identification accuracy and (ii) pinyin accuracy analysed with two logistic mixed  
752 effect models. There was only one predictor, *variability-condition* (LV versus HV, LV versus  
753 HVB) for both models. The descriptive statistics are displayed in *Figure 3.5.2*.

754 For tone accuracy, participants in the LV group performed showed no significant  
755 difference compared with the HV group ( $\beta = -0.33$   $SE = 0.22$ ,  $z = -1.54$ ,  $p = 0.12$ ) and the HVB  
756 group ( $\beta = -0.24$ ,  $SE = 0.22$ ,  $z = -1.13$ ,  $p = .26$ ). There was also no significant difference between  
757 groups in pinyin accuracy (LV versus HV:  $\beta = 0.09$ ,  $SE = 0.25$ ,  $z = 0.35$ ,  $p = 0.73$ ; LV versus  
758 HVB:  $\beta = -0.04$ ,  $SE = 0.25$ ,  $z = -0.17$ ,  $p = 0.86$ ).

759

### 760 3.5.3 *Summary of Tests Administered at Post-Test Only*

761 In sum, the analysis of the Picture Identification results suggests that on average,  
762 participants had higher accuracy on trained voice trials, demonstrating greater ease in identifying  
763 the words which had been trained repeatedly with the same speaker. The interaction between  
764 *voice-novelty* and *variability-condition* suggests that exclusive training on a single speaker in the  
765 LV condition boosted performance specifically for that speaker. Critically, there is no evidence  
766 for greater performance with untrained items in either the HV or the HVB condition, in contrast  
767 to the hypotheses. For Picture Naming, no significant result was found.

768

### 769 3.6 *Analyses with Individual Aptitude*

770 In order to look at the effect of learner aptitude and the interaction between this factor  
771 and variability condition, we first calculated the mean accuracy at pre-test on the PCPT for each  
772 participant. This was used as a continuous predictor (*aptitude*) and added to each of the models  
773 reported above. In addition we added the interaction between this factor and key experimental  
774 factors (see *Table 3*). Based on Perrachione et al. (2011) and Sadakata and McQueen (2014),  
775 high variability should benefit high aptitude participants only, while low variability would  
776 benefit low aptitude participants only. In our design, we used a continuous measure of individual  
777 ability rather than a binary division of high and low variability. We therefore predicted a stronger  
778 positive correlation between *aptitude* and amount of learning in the high variability condition  
779 than in the low variability condition. In the models for the pre- and post-test data (i.e. Three  
780 Interval Oddity and Word repetition) this could show up as a three way interaction between  
781 *condition*, *test-session* and *aptitude*. This interaction could possibly be modulated by *item-*  
782 *novelty* (four way interaction), since variability is thought to be key for generalization. In the

783 tests only administered post training, we looked for an interaction between *aptitude* and  
784 *condition* (since we have no measure at pre-test, and since there was no novelty manipulation  
785 here).

786 Each model reported in *Table 3* contained all the fixed and random effects included in the  
787 original models (although in some cases we had to remove correlations between slopes due to  
788 problems with convergence). For each of the new models we first confirmed that adding in the  
789 new effects and interactions with the individual measures did not change any of the previously  
790 reported patterns of significance for the experimental effects (see script  
791 [https://osf.io/wdh8a/?view\\_only=d1557462138447ffbaafaf7a59662df8](https://osf.io/wdh8a/?view_only=d1557462138447ffbaafaf7a59662df8))<sup>2</sup>.

792 The results are shown in *Table 3*. *Aptitude* can be seen to contribute to the model for  
793 training, the Three Interval Oddity task and the pinyin accuracy measure in the Word Repetition  
794 and the Picture Identification task; however there was no interaction between *aptitude* and any  
795 other factor. Thus there was no evidence that this measure of aptitude correlated with  
796 participants ability to benefit from training (no interaction with *test-session*), nor – critically for  
797 our hypothesis - did this differ by training condition (no interaction with *condition* or with *test-*  
798 *session* by *condition*).

799 Although the analyses use a continuous measure of PCPT for the purposes of  
800 visualization, *Figure 3.6.1* and *Figure 3.6.2* uses the mean accuracy for participants split into  
801 aptitude groups using a median split based on their PCPT score. *Figure 3.6.1* demonstrated the

---

<sup>2</sup> Note that models did not include all of the interactions between aptitude and each of the fixed effects in the original model, due to problems of convergence. Therefore the effects reported in Table 3 are the full set of additional fixed effects included in the new version of the model. For training data, recall that in section 3.4 we could not fit a converging model to the data from all three conditions, and instead presented two models – one for the LV+HV data, one for the LV+HVB data. We therefore attempted to include the effects of *aptitude* in each of these models; however neither model would converge if interactions with training-session were included and so these were removed. In the second model it was also necessary to remove the random slope for training session to achieve convergence.

802 results for the Three Interval Oddity task and Training task. The post-test data from the Picture  
803 Naming and Picture Identification tasks are shown in *Figure 3.6.2*. The production task, Word  
804 Repetition is shown in *Figure 3.6.3*. The results of the main effect of aptitude and its interaction  
805 with other predictors are summarised in *Table 3*.

806 In sum, participants with higher aptitude measure were better at the tasks, but there is no  
807 evidence either that this affected their improvement due to training, or, critically, their ability to  
808 benefit from the different variability exposure sets.

809

#### 810 **4 Discussion**

811 The current study investigated the effect of different types of phonetic training on English  
812 speakers learning of novel Mandarin words and tones. To our knowledge, this is the first study to  
813 train naive participants on all four Mandarin tones, using real language stimuli embedded in a  
814 word learning task. Learning was examined using a range of perception and production tasks.  
815 Following previous literature, we compared three training conditions: low variability (single  
816 speaker), high variability (four speakers, presented intermixed) and high variability blocked (four  
817 speakers, presented in blocks). We also administered tests designed to tap individual aptitude in  
818 the perception of pitch contrasts, adapted from the previous literature. The results indicated that  
819 participants' performance increased during training and that training also led to improved  
820 performance on pre- to post- tests of discrimination and production, with evidence of  
821 generalization to new voices and items. Participants also showed some ability to recall trained  
822 words – including their tones – in a naming task administered at post-test. However the only  
823 place where we saw any effect of the variability manipulation was in the training task (and with  
824 trained items in the picture identification task, which was highly similar to training), where the

825 *low* variability group outperformed both high variability groups. Critically, we found no  
826 evidence in any of our tests that high variability input benefitted learning or generalization, nor  
827 did we find any evidence of an interaction between individual aptitude and the ability to benefit  
828 from high variability training. In the following discussion, we first consider the findings from  
829 each task in turn before turning to a more general discussion of our findings concerning the  
830 benefit of high variability input.

831

#### 832 **4.1 Training and Picture Identification Tasks**

833 The training task employed in this study was a 2AFC task, where participants had to  
834 identify the correct meaning of a Mandarin word based on its tone. The results from training  
835 indicate that participants performed better in the single speaker LV training than in either the  
836 multiple speaker HV or HVB groups. This difference was present from the first session for the  
837 LV-HV contrast, and from the second session for the LV-HVB contrast (i.e. the first session  
838 where the two conditions differ). Greater difficulty with multiple speaker input is line with the  
839 findings of Perachione et al. (2011), although the differences did not emerge so rapidly in that  
840 study, possibly due to there being fewer trials per session). Intuitively, repeated exposure to the  
841 single speaker in the LV condition allows for greater adaption to speaker specific cues, whereas  
842 in the HV condition participants have to adapt to multiple speakers. This is particularly difficult  
843 in the unblocked HV condition, where trial-by-trial adaption is needed, which is effortful for  
844 participants (Magnuson & Nusbaum, 2007). Importantly, however, for all three groups, their  
845 performance gradually increased over each session. In combination with the fact that their  
846 performance on the other tasks increased after training, this indicates that the training task and  
847 materials were effective.

848 Critically, the *Picture Identification* test– a version of the training task without feedback  
849 which was administered post training – replicated this LV benefit for trained items, but  
850 demonstrated it did *not* extend to new *untrained speakers*. In fact, performance on *untrained*  
851 *speakers* was similar across conditions: participants performed more poorly than with the *trained*  
852 speaker, but were nonetheless above chance even with the untrained speaker. This indicates  
853 across-speaker generalization which did *not* depend on witnessing speaker variability in training,  
854 a point to which we return below.

#### 855 **4.2 Three Interval Oddity Task**

856 Our key test of perceptual discrimination was a three interval oddity task, where  
857 participants had to indicate the “different” word from a set of three. In each trial, the two foil  
858 words were the same word, and differed from the target word only in tone. Improvement in this  
859 task was significant but relatively modest (from 59% to 66%, following 8 training sessions),  
860 however there are many aspects which make this task more difficult than those used in previous  
861 studies. In particular, having each stimulus produced by a different speaker makes noting the  
862 similarity across tokens much harder, something we discovered in pilot work, where even before  
863 training participants were near ceiling with an equivalent task speakerin which the same speaker  
864 produced all three stimuli within a single trial. This is not a feature of any of the tests used in  
865 Perrachione et al. (2011) or Sadakata and McQueen (2014). In addition, we tested all four tone  
866 contrasts, including those involving Tone 3 (which Perrachione et al., 2011, did not include since  
867 it was considered perceptually the most confusable tone).

868 It is important to note that since all of the speakers in these test items were new,  
869 improvement in this test indicates generalization over speakers. Moreover, we did not see  
870 differences in the extent of improvement for *trained* versus *untrained* items, indicating that



871 improved tone discrimination is not item specific. Critically, this improvement following training  
872 occurred equally across the three variability conditions, indicating that input variability was not  
873 necessary for generalization, a point to which we return below.

874         Another result from this test was that we found evidence that some trial types were harder  
875 than others. Specifically, at pre-test, participants showed greatest performance for trials where  
876 one of the speakers was male and the other two were female, and the target “odd man” was the  
877 male speaker (“easy” trials). In contrast, they showed worst performance if there was one male  
878 and two female speakers, but the “odd man” was one of the female speakers (“hard” trials).  
879 Middle level performance was shown for trials where all three speakers were female (“neutral”  
880 trials). This is presumably due to participants relying on perceptual cues associated with speaker  
881 gender to do the task. Interestingly, our analyses showed that performance only increased for the  
882 trials where the odd man was not the lone male (the “neutral” and “hard” ones), and not for those  
883 where the male was the odd man. Given that participants are not near ceiling at pre-test (67%), it  
884 is perhaps surprising that their trained knowledge of the tone contrasts does not boost their  
885 performance. One possibility is although they are now better able to use tone cues, they are also  
886 *less* likely to use gender based cues, which they may now realize are less reliable, masking  
887 improvement based on tone for these particular test items.

888

### 889 **4.3 Production Tasks**

890         In this study, we used two production tasks: a word repetition task, administered pre and  
891 post training, and a picture naming task administered at post-test only. In the word repetition  
892 task, participants repeated a selection of Mandarin words produced by a native speaker, half of  
893 which would occur/had occurred in the training set, and half of which were untrained. We saw a

894 significant, though relatively modest improvement in participants' ability to reproduce the tone  
895 of the stimuli, such that it could be identified by a native speaker (from pre- to post- test: 70% to  
896 76%). This provides some evidence that purely perceptual training can influence production, in  
897 line with the findings of Bradlow and Pisoni (1999) and Zeromskaite (2014). Moreover the fact  
898 that participants showed a small but nevertheless significant increase in their ability to accurately  
899 repeat the segmental information (63% to 64% of words produced with correct segments)  
900 suggests that even though our training specifically targeted tone discrimination (which was all  
901 that was necessary to succeed in the training task) there was some more incidental learning of  
902 other aspects of the stimuli. As in the three interval oddity task, we again saw equivalent  
903 improvement for both trained and untrained items, and there was no difference in the extent of  
904 improvement in the different types of conditions, indicating that transfer did not rely on speaker-  
905 variability in the input.

906       Finally, in our picture naming vocabulary test participants were required to produce the  
907 trained words in response to pictures, without prompts. Participants showed some ability to  
908 recall both the segmental phonology and the tones, although unsurprisingly, accuracy here was  
909 considerably less than in the word repetition test for both (tone accuracy: 47% pinyin accuracy  
910 50%). Again we saw no differences between variability conditions, which is surprising given the  
911 substantial literature on vocabulary learning showing that there is a benefit of training with  
912 multiple speakers which can be tapped by naming tasks. We return to this point in the following  
913 section.

914

#### 915 **4.4    *The Role of High Variability Materials in Training and Generalization***

916 In the current study, across all of our different tests, we did not find either an overall  
917 benefit of exposure to high variability training materials, or any interaction between such a  
918 benefit and individual aptitude. We consider first the lack of *overall* variability benefit. This  
919 finding is in line with the lack of a main effect in the previous tone-training studies, yet it is at  
920 odds with some other phonetic training studies (Logan et al. 1991, 1993; Clopper & Pisoni,  
921 2004; Sadakata & McQueen 2013). This suggests the possibility that this overall variability  
922 benefit is restricted to segmental rather than tonal phonetic learning, at least for speakers of a  
923 non-tonal L1. It is harder to reconcile the lack of benefit for vocabulary learning in the picture  
924 naming task, given the findings of Barcroft, Sommers and colleagues (Barcroft & Sommers,  
925 2005, 2014; Sommers & Barcroft, 2007, 2011), particularly for our measure of segmental  
926 learning which is quite similar to that used in previous experiments, although the nature of our  
927 focused phonetic training is a possible explanation. However, it is also important to acknowledge  
928 the limitations of a null result: we have no evidence of an effect, but we also don't have evidence  
929 that there is *no* effect (see Dienes, 2008, for discussion of this distinction), and type 2 error is of  
930 course a possibility. On the other hand, at least for the phonetic training literature, while there is  
931 a longstanding *assumption* that speaker variability is important for generalization, as discussed in  
932 the introduction, the original test of this by Logan, Lively and Pisoni was extremely low powered  
933 considering they only tested three participants for the learning effect of generalisation. In  
934 addition, there are only a handful of published studies which have revisited this result (e.g.  
935 Lively et al., 1993; Lee, Perrachione, Dees & Wong, 2007; Gao, Low, Jin & Sweller, 2013). The  
936 current results suggest that there is need for further research to establish the extent to which the  
937 variability effect is replicable, and the extent to which it applies across different types of  
938 linguistic domains.

939           Turning to the lack of interaction with individual differences, the key question is why our  
940 result is different from that of Perrachione et al. (2011) and Sadakata and McQueen (2014).  
941 There are a variety of differences across the studies which could underpin the difference. Recall  
942 that although we set out to use similar methods to the previous studies, we were unable to use the  
943 data from our version of the Sadakata and McQueen test, due to too few participants meeting  
944 their inclusion criteria. The test which we did use is similar to that used by Perrachione and  
945 Wong, however our task is harder since it uses all six Mandarin vowels (whereas the original  
946 study used five, without /u/) and all of the Mandarin tones (where they used three, without Tone  
947 3). This change means that that we cannot easily contrast the range of participant scores in the  
948 two studies and it may be that the spread of ability of our participant is different from theirs. We  
949 also note that our statistical analyses are different from both of the previous studies in that they  
950 took their continuous aptitude measures and turned these into binary factors using a “cut off”,  
951 where as our statistical approach allows us to use them as continuous variables. However this  
952 should in principle make our approach more powerful than in previous studies. Moreover, we  
953 included a variety of both perception and production tasks. Thus, even if individual aptitude  
954 affects only specific aspects of learning or are only discernible in certain types of tests, we would  
955 have expected it to emerge in at least one of our tasks. Again, we have to acknowledge the  
956 possibility of type 2 error in our study, particularly since we know that interactions require  
957 greater samples than main effects to achieve the same power. On the other hand, type 1 error in  
958 the original studies is of course always possible.

959

#### 960 **4.5    *Future Directions***

961           As discussed above, it is difficult to draw strong conclusions from the null effects in the  
962 current work. An important limitation here is that – given the differences in materials and tasks  
963 compared with previous work - it is not clear what the size of the effects we should have  
964 expected. This makes it difficult to conduct a power analysis. It also precludes an informed  
965 Bayes factor analysis – which could potentially allow us to differentiate evidence for the null  
966 from evidence that is ambiguous (Dienes, 2008) – since this also requires a measure of the  
967 predicted effect size for each hypothesis<sup>3</sup>. We therefore suggest that it would be useful to  
968 implement a direct, high powered replication of these previous studies. We note that obtaining  
969 90% power would likely require a much larger sample than is standard in these types of studies.  
970 Given the time consuming nature of these multiple session training studies, we suggest that  
971 moving to online testing may be necessary to make this feasible (see Xie et al. 2018 for an  
972 example of an acoustic training study done over the web), or alternately multi-lab collaboration  
973 may be necessary.

974           Although direct replication will play a useful role in establishing these effects, we believe  
975 that ultimately it will also be important to develop a more nuanced approach to measuring the  
976 factors leading to different levels of aptitude both in tone learning, and in other types of phonetic  
977 learning. We note that here in addition to not seeing the predicted interaction with variability, we  
978 also didn't see interactions between aptitude and training session in any of our tasks, suggesting  
979 that our aptitude measure predicted baseline performance on the task and *not* the ability to  
980 improve due to training. In addition, the tasks used to measure “aptitude” are quite similar in  
981 nature to the training and test tasks, decreasing their explanatory value. Our ongoing work  
982 explores the combined predictive value of a range of measures including measures of attention,

---

<sup>3</sup> It is possible to inform the H1 using other parts of the same dataset (e.g. see Dienes 2018). However in the current work it was unclear how to do this, particularly for the interactions which are the key hypothesis.

983 working memory and musical ability. Identifying factors which are predictive of aptitude for  
984 tone learning has clear implications for teaching and the personalisation of teaching methods.

985

## 986 **5 Conclusion**

987 We trained naive participants on all four Mandarin tones, using real language stimuli  
988 embedded in a word learning task. We found improvements in both production and perception of  
989 tones which transferred to novel voices and items. We found that learning was greatest for  
990 training with a single voice but that training with a single voice versus four voices (whether  
991 intermixed or blocked) lead to equal amounts of generalization. Although learner aptitude  
992 predicted performance in most tasks, there was no evidence that different levels of aptitude lead  
993 to better or worse learning from different types of training input.

994 **References**

- 995 Aliaga-García, C., & Mora, J. C. (2009). Assessing the effects of phonetic training on L2 sound  
996 perception and production. In M. A. Watkins, A. S. Rauber, & B.O. Baptista (Eds.). *Recent*  
997 *Research in Second Language Phonetics/Phonology: Perception and Production* (pp. 2-  
998 31). Newcastle upon Tyne, UK: Cambridge Scholars Publishing.
- 999 Audacity Team. (2015). Audacity (Version 2.1.1). *Computer Program*. Retrieved May, 2015,  
1000 from <http://audacityteam.org/>
- 1001 Alshangiti, W., & Evans, B. G. (2014, May). Investigating the domain-specificity of phonetic  
1002 training for secondlanguage learning: Comparing the effects of production and perception  
1003 training on the acquisition of English vowels by Arabic learners of English. In *the*  
1004 *Proceedings of the International Seminar for Speech Production, Cologne, Germany*.
- 1005 Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed  
1006 random effects for subjects and items. *Journal of memory and language*, 59(4), 390-412.  
1007 <https://doi.org/10.1016/j.jml.2007.12.005>
- 1008 Bao, Z. (1999). *The Structure of Tone*. New York: Oxford University Press. ISBN 0-19-511880-  
1009 4.
- 1010 Barcroft, J., & Sommers, M. S. (2005). Effects of acoustic variability on second language  
1011 vocabulary learning. *Studies in Second Language Acquisition*, 27, 387-414.  
1012 <https://doi.org/10.1017/S0272263105050175>
- 1013 Barcroft, J., & Sommers, M. S. (2014). Effects of variability in fundamental frequency on L2  
1014 vocabulary learning: A comparison between learners who do and do not speak a tone  
1015 language. *Studies in Second Language Acquisition*, 36(3), 423-449.  
1016 <https://doi.org/10.1016/j.neuropsychologia.2006.11.015>

- 1017 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for  
1018 confirmatory hypothesis testing: Keep it maximal. *Journal of memory and language*, 68(3),  
1019 255-278. <https://doi.org/10.1016/j.jml.2012.11.001>
- 1020 Boersma, P., & Weenink, D. (2015). Praat: doing phonetics by computer [Computer program].  
1021 Version 5.4.14, retrieved 24 July 2015 from <http://www.praat.org/>
- 1022 Bradlow, A. R., & Pisoni, D. B. (1999). Recognition of spoken words by native and non-native  
1023 listeners: Talker-, listener-, and item-related factors. *The Journal of the Acoustical Society*  
1024 *of America*, 106, 2074-2085. <http://dx.doi.org/10.1121/1.427952>
- 1025 Bradlow, A. R., Pisoni, D. B., Akahane-Yamada, R., & Tohkura, Y. I. (1997). Training Japanese  
1026 listeners to identify English/r/and/l: IV. Some effects of perceptual learning on speech  
1027 production. *The Journal of the Acoustical Society of America*, 101, 2299-2310.  
1028 <https://doi.org/10.1121/1.418276>
- 1029 Bygate, M., Swain, M., & Skehan, P. (2013). *Researching pedagogic tasks: Second language*  
1030 *learning, teaching, and testing*. London UK: Routledge.
- 1031 Cicchetti, D. V. (1994). Guidelines, criteria, and rules of thumb for evaluating normed and  
1032 standardized assessment instruments in psychology. *Psychological assessment*, 6(4), 284.  
1033 <http://dx.doi.org/10.1037/1040-3590.6.4.284>
- 1034 Clopper, C. G., & Pisoni, D. B. (2004). Some acoustic cues for the perceptual categorization of  
1035 American English regional dialects. *Journal of phonetics*, 32(1), 111-140.  
1036 [https://doi.org/10.1016/S0095-4470\(03\)00009-3](https://doi.org/10.1016/S0095-4470(03)00009-3)
- 1037 Cohen, A. D. (2014). *Strategies in learning and using a second language*. London UK:  
1038 Routledge.
- 1039



- 1040 Dienes, Z. (2008). *Understanding psychology as a science: An introduction to scientific and*  
1041 *statistical inference*. Macmillan International Higher Education.
- 1042 Dienes, Z., & Mclatchie, N. (2018). Four reasons to prefer Bayesian analyses over significance  
1043 testing. *Psychonomic bulletin & review*, 25(1), 207-218. [https://doi.org/10.3758/s13423-](https://doi.org/10.3758/s13423-017-1266-z)  
1044 [017-1266-z](https://doi.org/10.3758/s13423-017-1266-z)
- 1045 Flege, J. E., Takagi, N., & Mann, V. (1995). Japanese adults can learn to produce  
1046 English/ɪ/and/ɪ/accurately. *Language and Speech*, 38, 25-55.  
1047 <https://doi.org/10.1177/002383099503800102>
- 1048 Gao, Y., Low, R., Jin, P., & Sweller, J. (2013). Effects of speaker variability on learning foreign-  
1049 accented English for EFL learners. *Journal of Educational Psychology*, 105, 649-665.  
1050 <http://dx.doi.org/10.1037/a0033024>
- 1051 Giannakopoulou, A., Brown, H., Clayards, M., & Wonnacott, E. (2017). High or low?  
1052 Comparing high and low-variability phonetic training in adult and child second language  
1053 learners. *PeerJ*, 5, e3209. DOI:[10.7717/peerj.3209](https://doi.org/10.7717/peerj.3209)  
1054
- 1055 Giannakopoulou, A., Uther, M., & Ylinen, S. (2013). Enhanced plasticity in spoken language  
1056 acquisition for child learners: Evidence from phonetic training studies in child and adult  
1057 learners of English. *Child Language Teaching and Therapy*, 29, 201-218.  
1058 <https://doi.org/10.1177/0265659012467473>
- 1059 Iverson, P., Ekanayake, D., Hamann, S., Sennema, A., & Evans, B. G. (2008). Category and  
1060 perceptual interference in second-language phoneme learning: An examination of  
1061 English/w/-/v/learning by Sinhala, German, and Dutch speakers. *Journal of Experimental*

- 1062 *Psychology: Human Perception and Performance*, 34, 1305. <https://doi.org/10.1037/0096->  
1063 [1523.34.5.1305](https://doi.org/10.1037/0096-1523.34.5.1305)
- 1064 Iverson, P., Hazan, V., & Bannister, K. (2005). Phonetic training with acoustic cue  
1065 manipulations: A comparison of methods for teaching English/r/-/l/to Japanese adults. *The*  
1066 *Journal of the Acoustical Society of America*, 118, 3267-3278.  
1067 <https://doi.org/10.1121/1.2062307>
- 1068 Jaeger, T. F. (2008). Categorical data analysis: Away from ANOVAs (transformation or not) and  
1069 towards logit mixed models. *Journal of memory and language*, 59(4), 434-446.  
1070 <https://doi.org/10.1016/j.jml.2007.11.007>
- 1071 Joanisse, M. F., Manis, F. R., Keating, P., & Seidenberg, M. S. (2000). Language deficits in  
1072 dyslexic children: Speech perception, phonology, and morphology. *Journal of*  
1073 *Experimental Child Psychology*, 77, 30-60. <https://doi.org/10.1006/jecp.1999.2553>
- 1074 Lambacher, S. G., Martens, W. L., Kakehi, K., Marasinghe, C. A., & Molholt, G. (2005). The  
1075 effects of identification training on the identification and production of American English  
1076 vowels by native speakers of Japanese. *Applied Psycholinguistics*, 26(2), 227-247.  
1077 <https://doi.org/10.1017/S0142716405050150>
- 1078 Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical  
1079 data. *Biometrics*, 159-174. DOI: 10.2307/2529310
- 1080 Lee, J., Perrachione, T. K., Dees, T. M., & Wong, P. C. (2007, August). Differential effects of  
1081 stimulus variability and learners' pre-existing pitch perception ability in lexical tone  
1082 learning by native English speakers. Paper presented at the *16th International Congress of*  
1083 *Phonetic Sciences* (pp. 1589-1592).

- 1084 Lengeris, A., & Hazan, V. (2010). The effect of native vowel processing ability and frequency  
1085 discrimination acuity on the phonetic training of English vowels for native speakers of  
1086 Greek. *The Journal of the Acoustical Society of America*, *128*, 3757-3768.  
1087 <https://doi.org/10.1121/1.3506351>
- 1088 Lively, S. E., Logan, J. S., & Pisoni, D. B. (1993). Training Japanese listeners to identify  
1089 English/r/and/l/. II: The role of phonetic environment and talker variability in learning new  
1090 perceptual categories. *The Journal of the Acoustical Society of America*, *94*, 1242-1255.  
1091 <https://doi.org/10.1121/1.408177>
- 1092 Logan, J. S., Lively, S. E., & Pisoni, D. B. (1991). Training Japanese listeners to identify  
1093 English/r/and/l/: A first report. *The Journal of the Acoustical Society of America*, *89*, 874-  
1094 886. <https://doi.org/10.1121/1.1894649>
- 1095 Magnuson, J. S., & Nusbaum, H. C. (2007). Acoustic differences, listener expectations, and the  
1096 perceptual accommodation of talker variability. *Journal of Experimental Psychology:*  
1097 *Human Perception and Performance*, *33*, 391. <https://doi.org/10.1037/0096-1523.33.2.391>
- 1098 Mattys, S. L., & Wiget, L. (2011). Effects of cognitive load on speech recognition. *Journal of*  
1099 *Memory and Language*, *65*(2), 145-160. <https://doi.org/10.1016/j.jml.2011.04.004>
- 1100 McGraw, K. O., & Wong, S. P. (1996). Forming inferences about some intraclass correlation  
1101 coefficients. *Psychological methods*, *1*(1), 30. <http://dx.doi.org/10.1037/1082-989X.1.1.30>
- 1102 Nusbaum, H. C., & Morin, T. M. (1992). Paying attention to differences among talkers. *Speech*  
1103 *perception, production and linguistic structure*, 113-134.
- 1104 Perrachione, T. K., Lee, J., Ha, L. Y., & Wong, P. C. (2011). Learning a novel phonological  
1105 contrast depends on interactions between individual differences and training paradigm

- 1106 design. *The Journal of the Acoustical Society of America*, 130, 461-472.
- 1107 <https://doi.org/10.1371/journal.pone.0089642>
- 1108 Powell, M. J. (2009). The BOBYQA algorithm for bound constrained optimization without  
1109 derivatives. *Cambridge NA Report NA2009/06*, University of Cambridge, Cambridge, 26-  
1110 46.
- 1111 Quené, H., & Van den Bergh, H. (2008). Examples of mixed-effects modeling with crossed  
1112 random effects and with binomial data. *Journal of Memory and Language*, 59(4), 413-425.  
1113 <https://doi.org/10.1016/j.jml.2008.02.002>
- 1114 R Development Core Team (2010). R: A Language and Environment for Statistical Computing,  
1115 Version R 3.3.2. Available at [www.r-project.org](http://www.r-project.org). Accessed September, 2017.
- 1116 Sadakata, M., & McQueen, J. M. (2013). High stimulus variability in nonnative speech learning  
1117 supports formation of abstract categories: Evidence from Japanese geminates. *The Journal*  
1118 *of the Acoustical Society of America*, 134, 1324-1335. <https://doi.org/10.1121/1.4812767>
- 1119 Sadakata, M., & McQueen, J. M. (2014). Individual aptitude in Mandarin lexical tone perception  
1120 predicts effectiveness of high-variability training. *Frontiers in Psychology*, 5, 1318.  
1121 <https://doi.org/10.3389/fpsyg.2014.01318>
- 1122 Sommers, M. S., & Barcroft, J. (2007). An integrated account of the effects of acoustic  
1123 variability in first language and second language: Evidence from amplitude, fundamental  
1124 frequency, and speaking rate variability. *Applied Psycholinguistics*, 28(2), 231-249.  
1125 <https://doi.org/10.1017/S0142716407070129>
- 1126 Sommers, M. S., & Barcroft, J. (2011). Indexical information, encoding difficulty, and second  
1127 language vocabulary learning. *Applied Psycholinguistics*, 32(2), 417-434.  
1128 <https://doi.org/10.1017/S0142716410000469>

1129

1130 Strange, W., & Dittmann, S. (1984). Effects of discrimination training on the perception of /r/-/l/

1131 by Japanese adults learning English. *Perception & Psychophysics*, 36, 131-145.

1132 <https://doi.org/10.3758/BF03202673>

1133 Wang, Y., Jongman, A., & Sereno, J. A. (2003). Acoustic and perceptual evaluation of Mandarin

1134 tone productions before and after perceptual training. *The Journal of the Acoustical Society*

1135 *of America*, 113, 1033-1043. <https://doi.org/10.1121/1.1531176>

1136 Wang, Y., Spence, M. M., Jongman, A., & Sereno, J. A. (1999). Training American listeners to

1137 perceive Mandarin tones. *The Journal of the Acoustical Society of America*, 106, 3649-

1138 3658. <https://doi.org/10.1121/1.428217>

1139

1140 Wong, P., & Perrachione, T. K. (2007). Learning pitch patterns in lexical identification by native

1141 English-speaking adults. *Applied Psycholinguistics*, 28, 565-585.

1142 <https://doi.org/10.1017/S0142716407070312>

1143 Xie, X., Weatherholtz, K., Bainton, L., Rowe, E., Burchill, Z., Liu, L., & Jaeger, T. F. (2018).

1144 Rapid adaptation to foreign-accented speech and its transfer to an unfamiliar talker. *The*

1145 *Journal of the Acoustical Society of America*, 143, 2013-2031.

1146 <https://doi.org/10.1121/1.5027410>

1147 Yip, M. (2002). *Tone. Cambridge textbooks in linguistics*. Cambridge: Cambridge University

1148 Press.

1149 Zeromskaite, I. (2014). The potential role of music in second language learning: A review

1150 article. *Journal of European Psychology Students*, 5, 78-88. <http://doi.org/10.5334/jeps.ci>

1151

**Table 1** (on next page)

*Age mean, age range, average number of language learned and mean starting age of learning the first L2 for participants in each condition.*

1

---

<b>Condition</b>	<b>Age Mean</b>	<b>Age Range</b>	<b>Languages Learned</b>	<b>Average Starting Age</b>
Low Variability	26.15	19-53	2.7	13.8
High Variability	25.65	19-47	2.5	12.2
High Variability Blocking	22.05	19-30	2.0	11.8

---

2

**Table 2** (on next page)

*Counterbalancing of voices for each task, training condition and version. LV = Low Variability; HV = High Variability; HVB = High Variability Blocking; PCPT = Pitch Contour Perception Test; CSTC = Categorisation of Synthesized Tonal Continua.*



1

Task	Condition	Voice				
		<i>Version 1</i>	<i>Version 2</i>	<i>Version 3</i>	<i>Version 4</i>	<i>Version 5</i>
Training	LV	F1	F2	F3	M1	M2
	HV &	F1	F2	F3	M1	M2
	HVB	F3	F1	M2	F1	F2
		M1	M1	F1	F2	F3
		M2	M2	F2	F3	M1
Word Repetition	All	F1	F2	F3	M1	M2
Picture Identification						
Trained Items	All	F1	F2	F3	M1	M2
New Items	All	F2	F3	M1	M2	F1
Three Interval Oddity	All	All versions: MN1, FN1, FN2, FN3				
PCPT	All	All versions: MN1, FN1, FN2, FN3				
CSTC	All	All versions: Synthesized voice				

2

3

**Table 3** (on next page)

*Statistics obtained when adding in participant aptitude (as measured by performance on the Pitch Contour Perception Test task at pre-test) into the models predicting performance on the test and training tasks.*

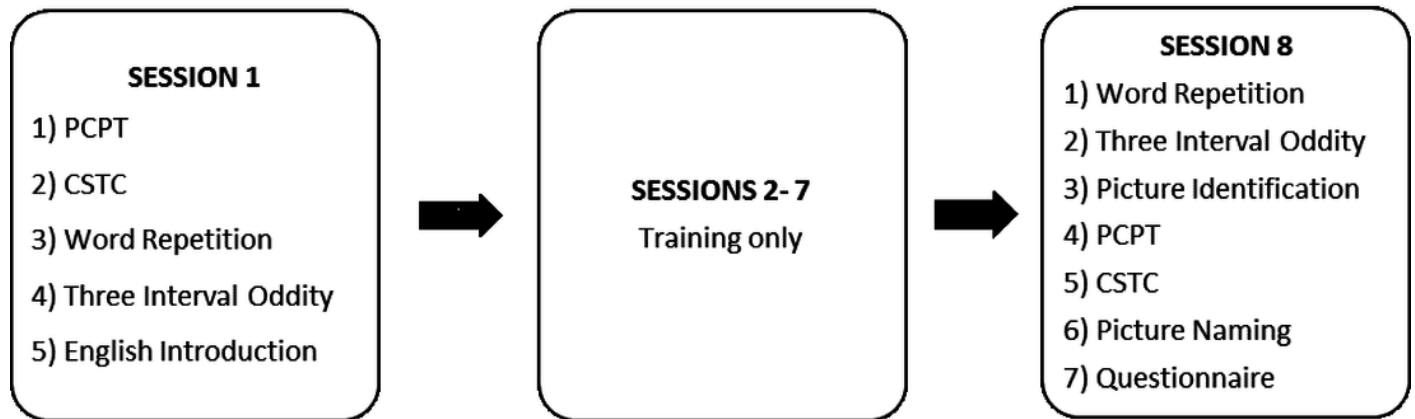
1

<b>Data Set</b>	<b>Coefficient Name</b>	<b>Statistics</b>
<i>Word Repetition:</i>	Aptitude	$\beta = 0.28$ , SE = 0.42, $z = 0.68$ , $p = .496$
<i>Tone Accuracy</i>	Aptitude by <i>Test-Session</i>	$\beta = -0.56$ , SE = 0.71, $z = -0.79$ , $p = .429$
<i>(Pre/Post)</i>	Aptitude by LV-HV Contrast by <i>Test-Session</i>	$\beta = 0.96$ , SE = 1.77, $z = 0.54$ , $p = .587$
	Aptitude by LV-HVB Contrast by <i>Test-Session</i>	$\beta = 0.11$ , SE = 1.51, $z = 0.07$ , $p = .941$
	Aptitude by LV-HV Contrast by <i>Test-Session</i> by	$\beta = -0.84$ , SE = 2.01, $z = -0.42$ , $p = .676$
	<i>Item-Novelty</i>	
	Aptitude by LV-HVB Contrast by <i>Test-Session</i> by	$\beta = 0.29$ , SE = 1.78, $z = 0.16$ , $p = .872$
	<i>Item-Novelty</i>	
<i>Word Repetition:</i>	<b>Aptitude</b>	<b><math>\beta = 0.62</math>, SE = 0.27, <math>z = 2.31</math>, <math>p = .021</math></b>
<i>Pinyin Accuracy</i>	Aptitude by <i>Test-Session</i>	$\beta = -0.28$ , SE = 0.51, $z = -0.56$ , $p = .576$
<i>(Pre/Post)</i>	Aptitude by LV-HV Contrast by <i>Test-Session</i>	$\beta = -0.07$ , SE = 1.28, $z = -0.05$ , $p = .958$
	Aptitude by LV-HVB Contrast by <i>Test-Session</i>	$\beta = -0.57$ , SE = 1.10, $z = -0.52$ , $p = .602$
	Aptitude by LV-HV Contrast by <i>Test-Session</i> by	$\beta = -1.70$ , SE = 1.74, $z = -0.98$ , $p = .328$
	<i>Item-Novelty</i>	
	Aptitude by LV-HVB Contrast by <i>Test-Session</i> by	$\beta = 0.21$ , SE = 1.55, $z = 0.14$ , $p = .892$
	<i>Item-Novelty</i>	
<i>Three Interval</i>	<b>Aptitude</b>	<b><math>\beta = 0.68</math>, SE = 0.31, <math>z = 2.19</math>, <math>p = .029</math></b>
<i>Oddity</i>	Aptitude by <i>Test-Session</i>	$\beta = 0.08$ , SE = 0.27, $z = 0.31$ , $p = .757$
<i>(Pre/Post)</i>	Aptitude by LV-HV Contrast by <i>Test-Session</i>	$\beta = 0.51$ , SE = 0.67, $z = 0.77$ , $p = .443$
	Aptitude by LV-HVB Contrast by <i>Test-Session</i>	$\beta = 0.48$ , SE = 0.58, $z = 0.83$ , $p = .409$
	Aptitude by LV-HV Contrast by <i>Test-Session</i> by	$\beta = 1.20$ , SE = 1.28, $z = 0.94$ , $p = .345$
	<i>Item-Novelty</i>	
	Aptitude by LV-HVB Contrast by <i>Test-Session</i> by	$\beta = -0.60$ , SE = 1.14, $z = -0.52$ , $p = .602$
	<i>Item-Novelty</i>	
<i>Training</i>	<b>Aptitude</b>	<b><math>\beta = 0.91</math>, SE = 0.31, <math>z = 2.93</math>, <math>p = .003</math></b>

<i>(Model including LV and HV conditions and LV only)</i>	Aptitude by LV-HV Contrast	$\beta = -0.43, SE = 0.33, z = -1.31, p = .192$
<i>Picture Identification (Post Only)</i>	<b>Aptitude</b>	<b><math>\beta = 1.48, SE = 0.75, z = 1.97, p = .049</math></b>
	Aptitude by Voice Novelty	$\beta = -0.28, SE = 0.86, z = -0.33, p = .744$
	Aptitude by LV-HV Contrast	$\beta = -0.23, SE = 1.85, z = -0.13, p = .899$
	Aptitude by LV-HVB Contrast	$\beta = 0.14, SE = 1.63, z = 0.09, p = .931$
	Aptitude by LV-HV Contrast by <i>Voice-Novelty</i>	$\beta = 3.47, SE = 2.07, z = 1.68, p = .094$
	Aptitude by LV-HVB Contrast by <i>Voice-Novelty</i>	$\beta = -1.07, SE = 1.82, z = -0.59, p = .558$
<i>Picture Naming: Tone Accuracy</i>	Aptitude	$\beta = 0.38, SE = 0.50, z = 0.75, p = .452$
	Aptitude by LV-HV Contrast	$\beta = -0.89, SE = 1.25, z = -0.71, p = .478$
	Aptitude by LV-HVB Contrast	$\beta = 0.11, SE = 1.09, z = 0.10, p = .921$
<i>Picture Naming: Pinyin Accuracy</i>	<b>Aptitude</b>	<b><math>\beta = -1.09, SE = 0.56, z = -1.93, p = .053</math></b>
	Aptitude by LV-HV Contrast	$\beta = 0.09, SE = 1.41, z = 0.06, p = .950$
	Aptitude by LV-HVB Contrast	$\beta = -0.10, SE = 1.23, z = -0.08, p = .939$

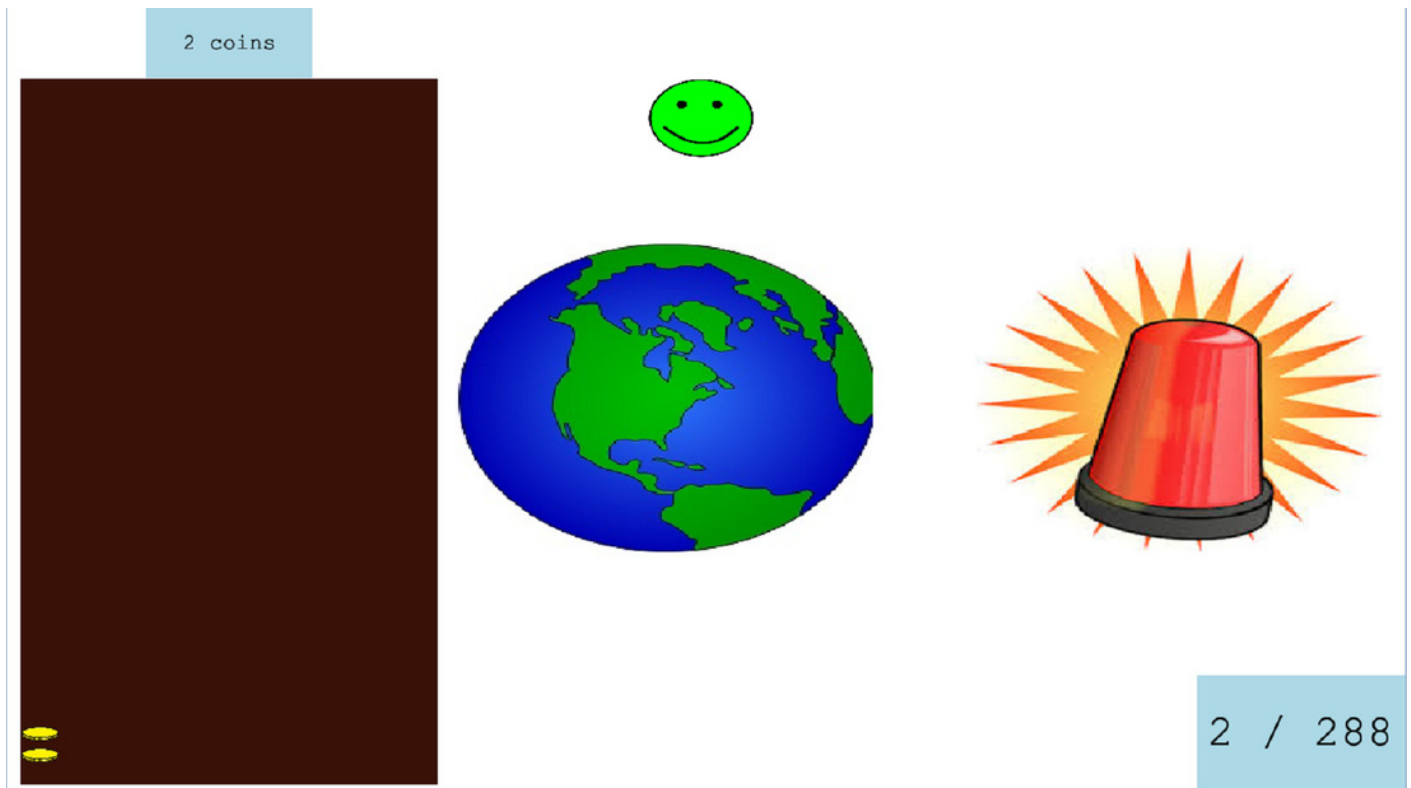
# Figure 1

Tasks completed in each of the eight sessions. (PCPT = Pitch Contour Perception Test; CSTC = Categorisation of Synthesized Tonal Continua).



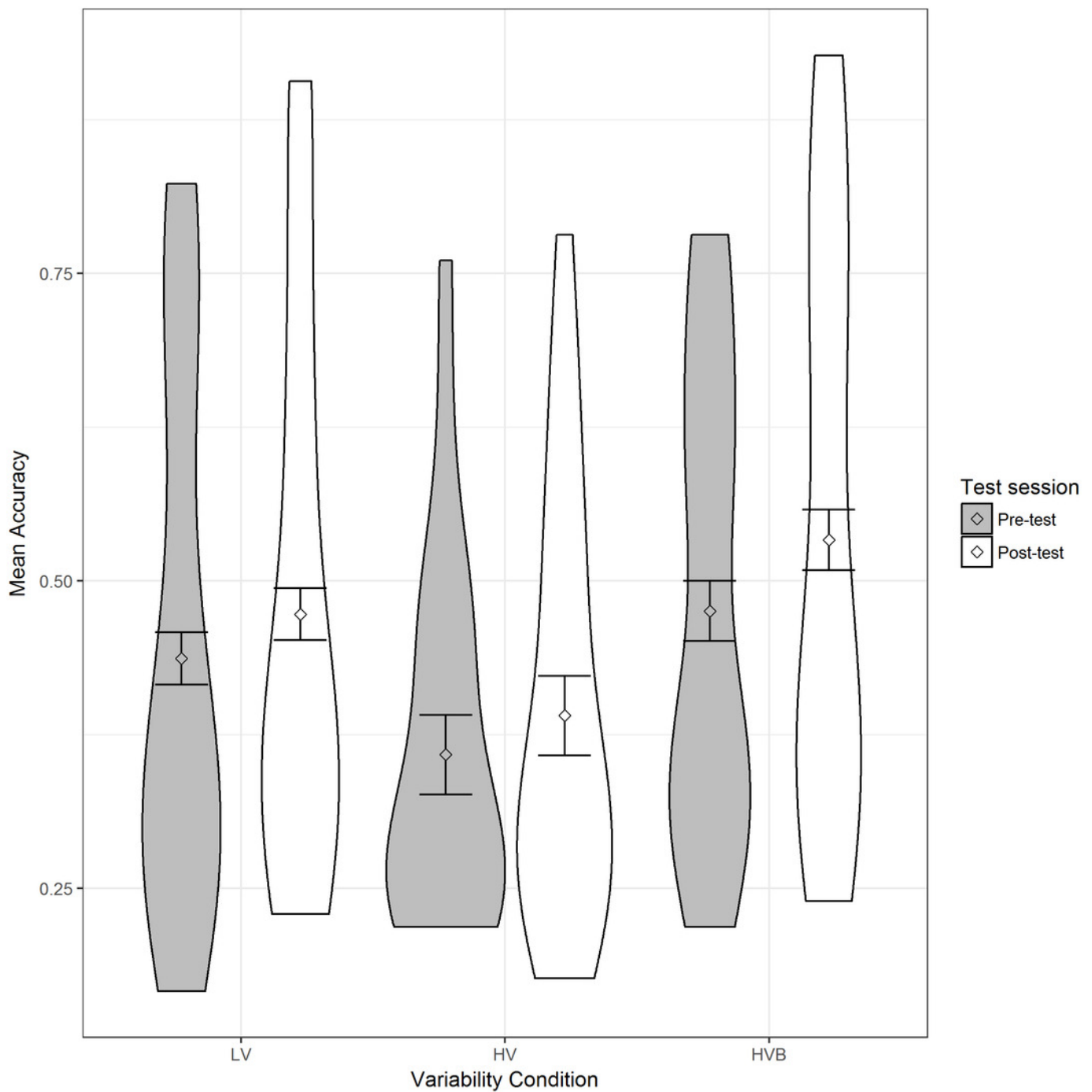
## Figure 2

Screen shot from the training task. The stimuli heard is 'dì', tone 4, [earth]. The foil picture on the right is 'dí' tone 2, [siren].



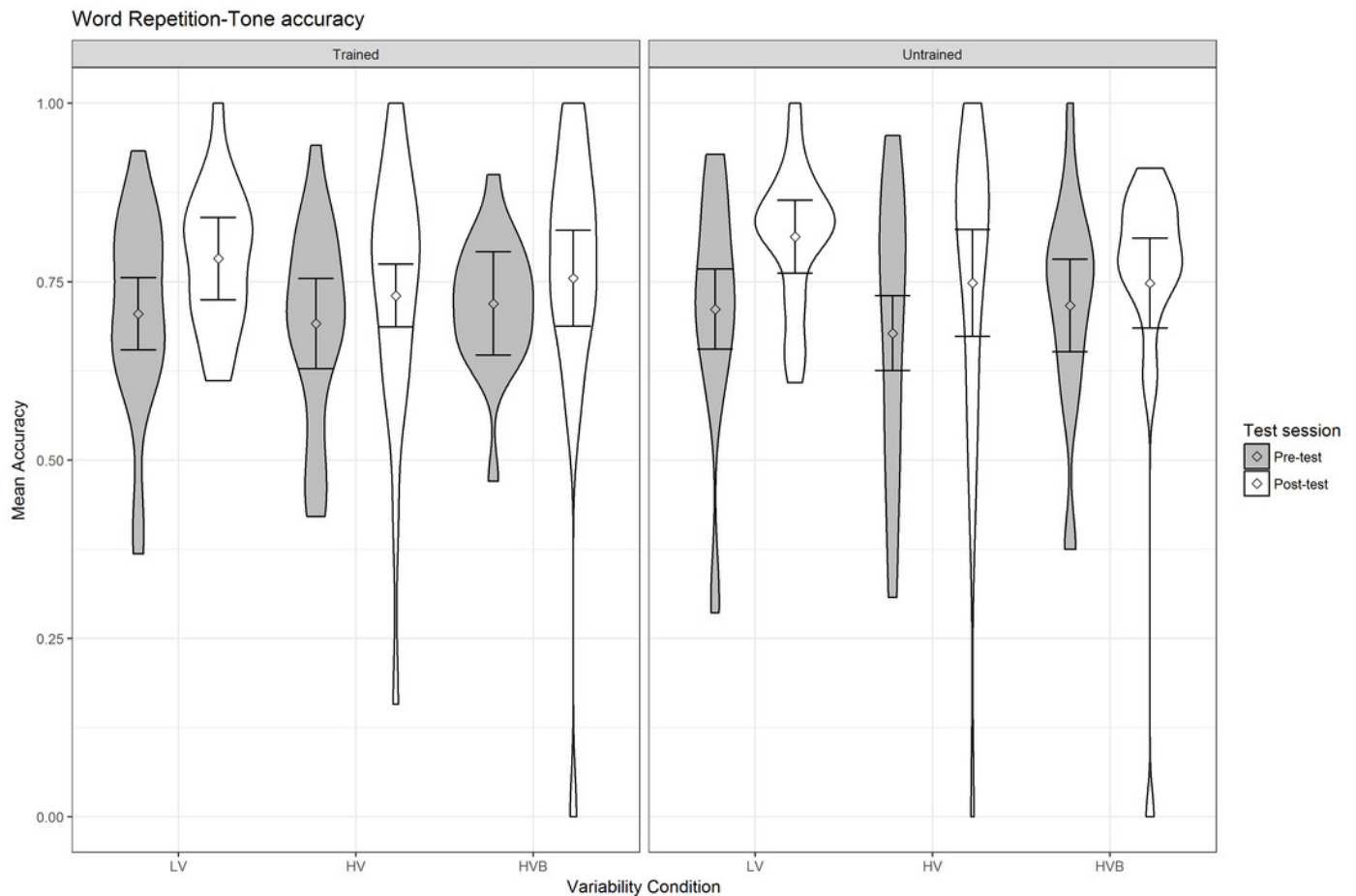
## Figure 3

Mean Accuracy from LV (Low Variability), HV (High Variability) & HVB (High Variability Blocking) groups in Pitch Contour Perception Test. Error bars represents the 95% confidence intervals.



## Figure 4

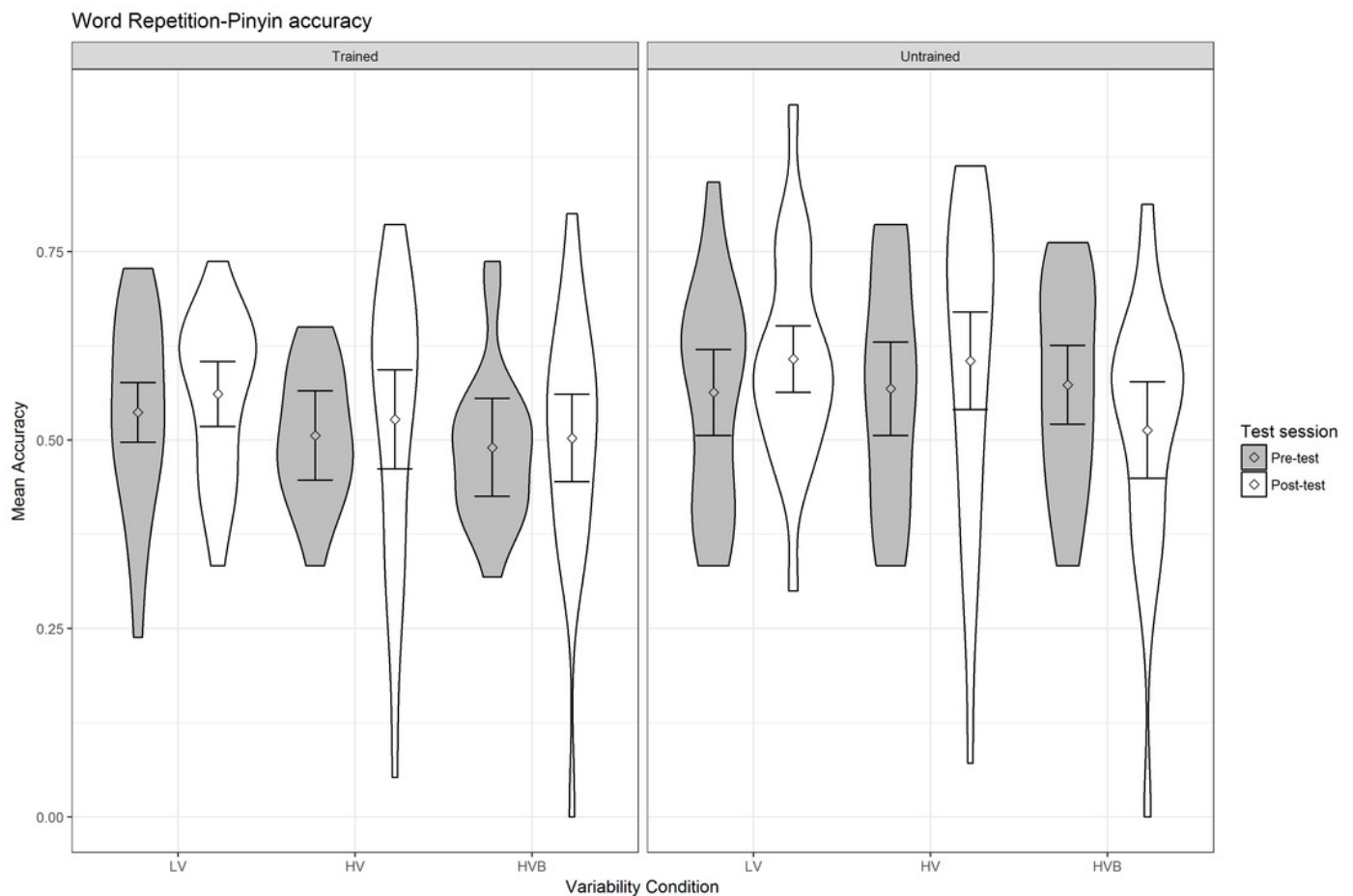
Accuracy of Word Repetition for LV (Low Variability), High Variability (HV) and High Variability Blocking (HVB) training groups in Pre- and Post-tests. Error bars show 95% confidence intervals.





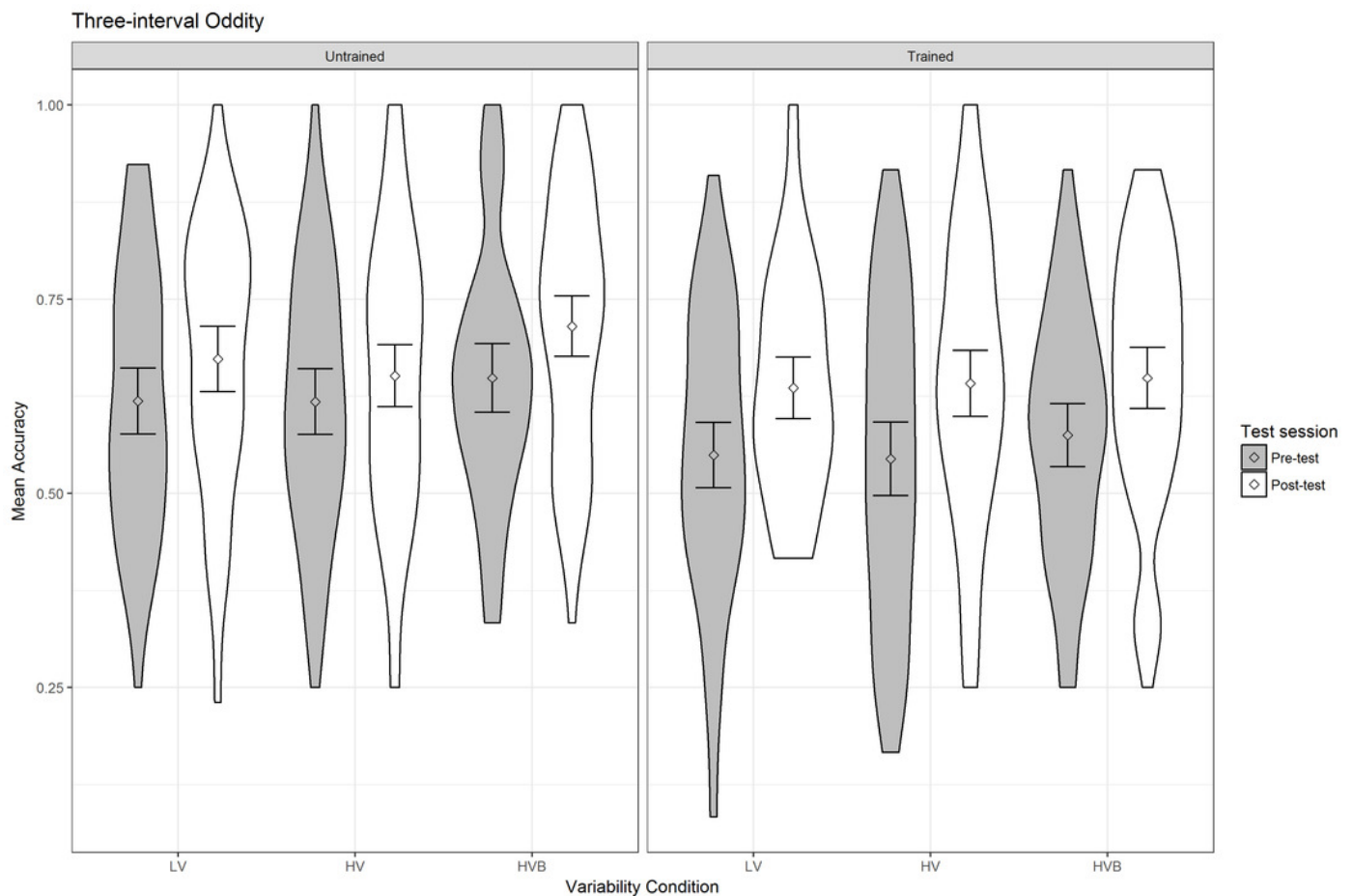
## Figure 5

Mean pinyin accuracy of Word Repetition for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups in Pre- and Post-tests. Error bars show 95% confidence intervals.



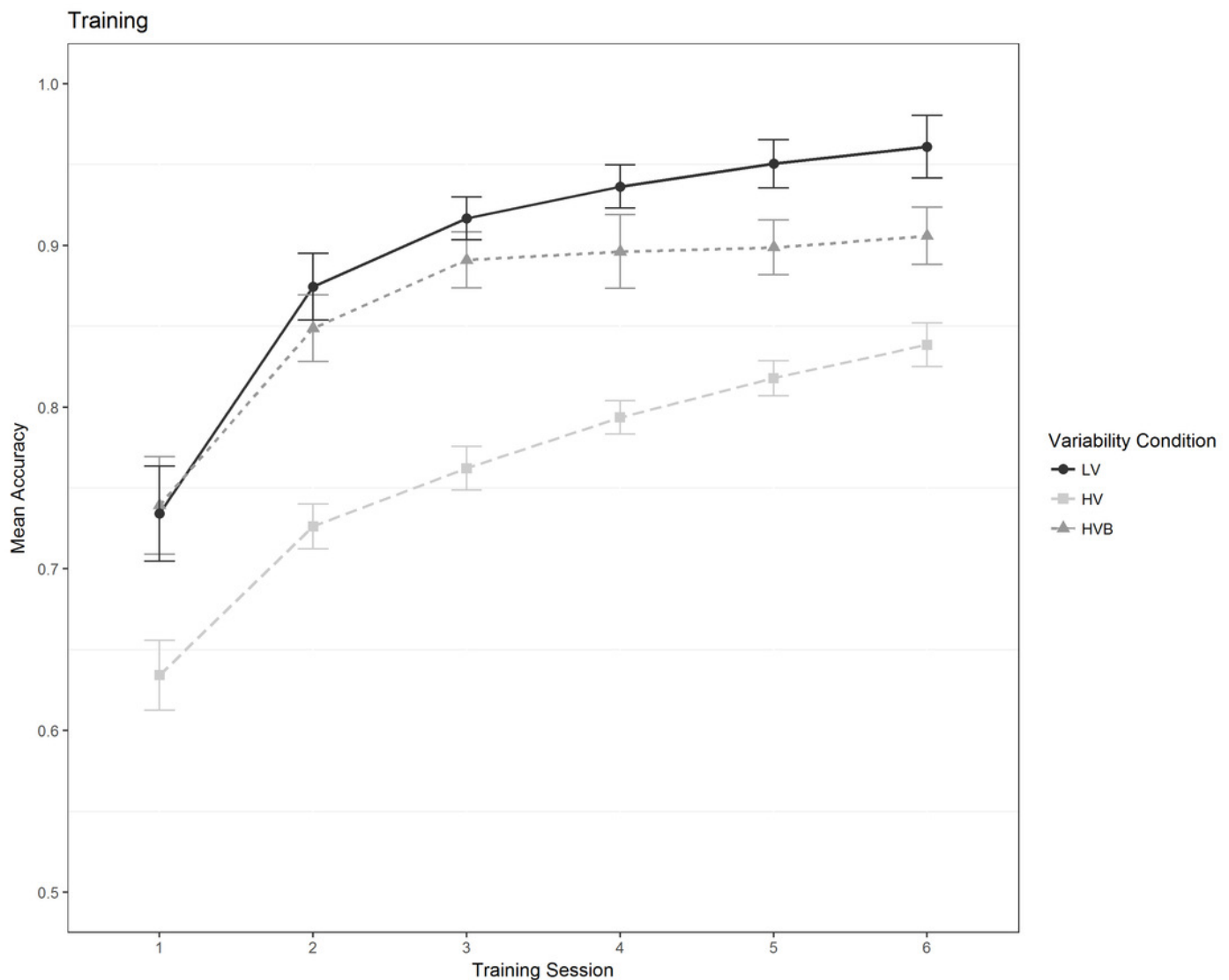
## Figure 6

Mean accuracy in Three Interval Oddity task for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups in Pre- and Post-tests. Error bars show 95% confidence intervals.



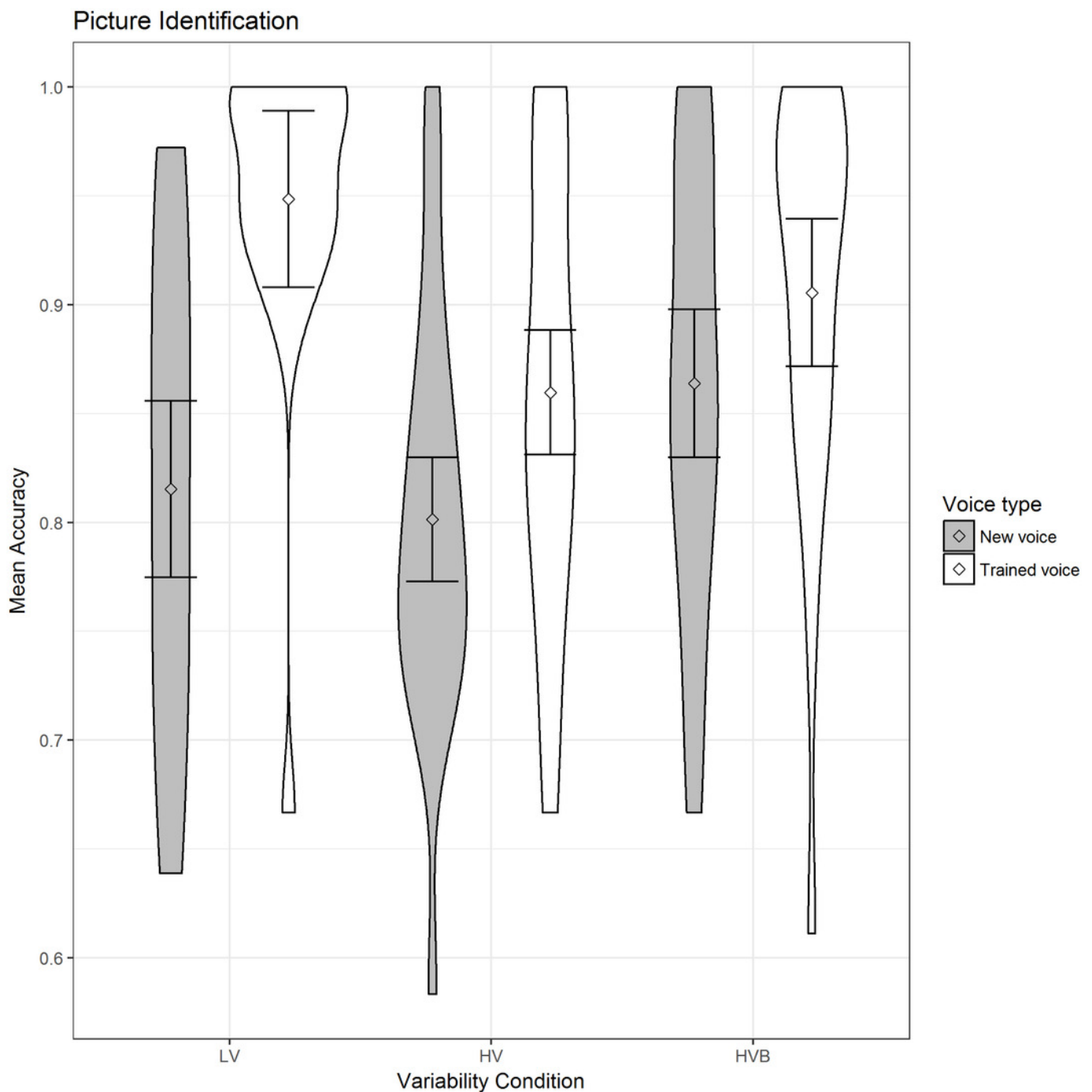
## Figure 7

Mean accuracy of Training for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups for each session. Y-axis starting from chance level. Error bars show 95% confidence intervals.



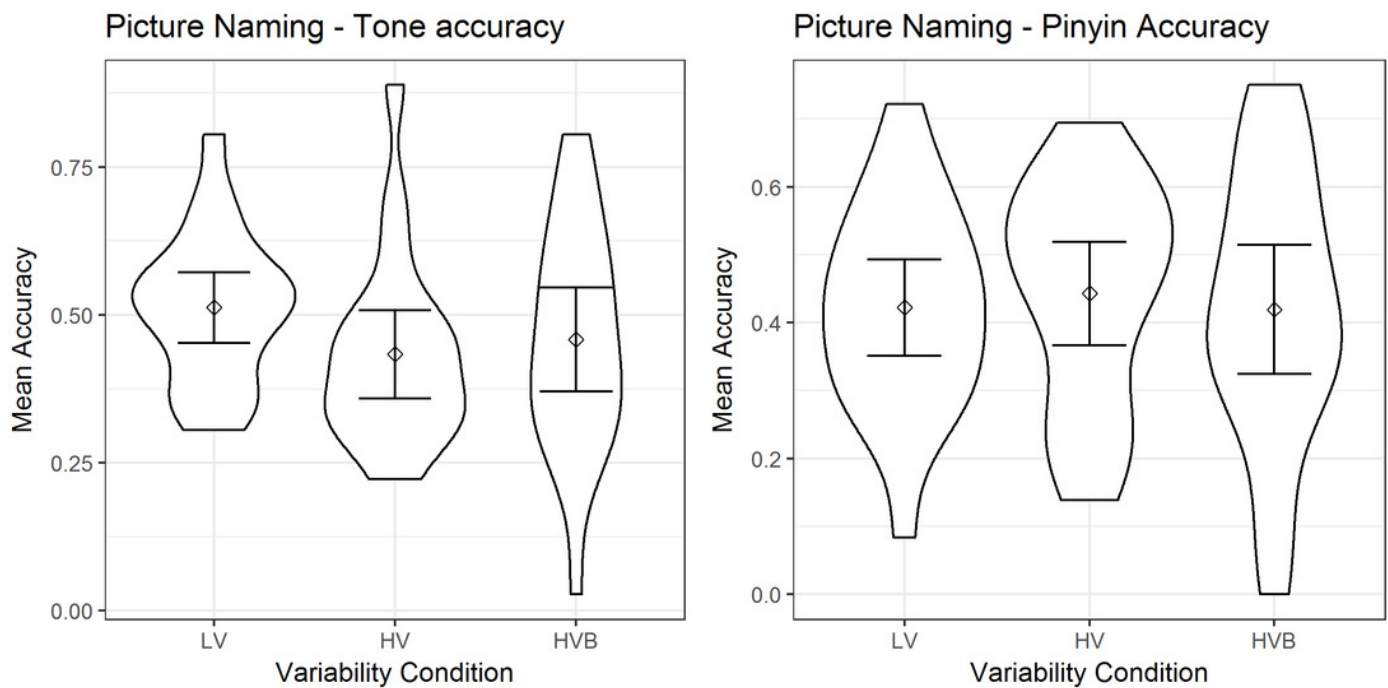
## Figure 8

Mean accuracy of Picture Identification for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups for new voices and trained voices. Error bars show 95% confidence intervals.



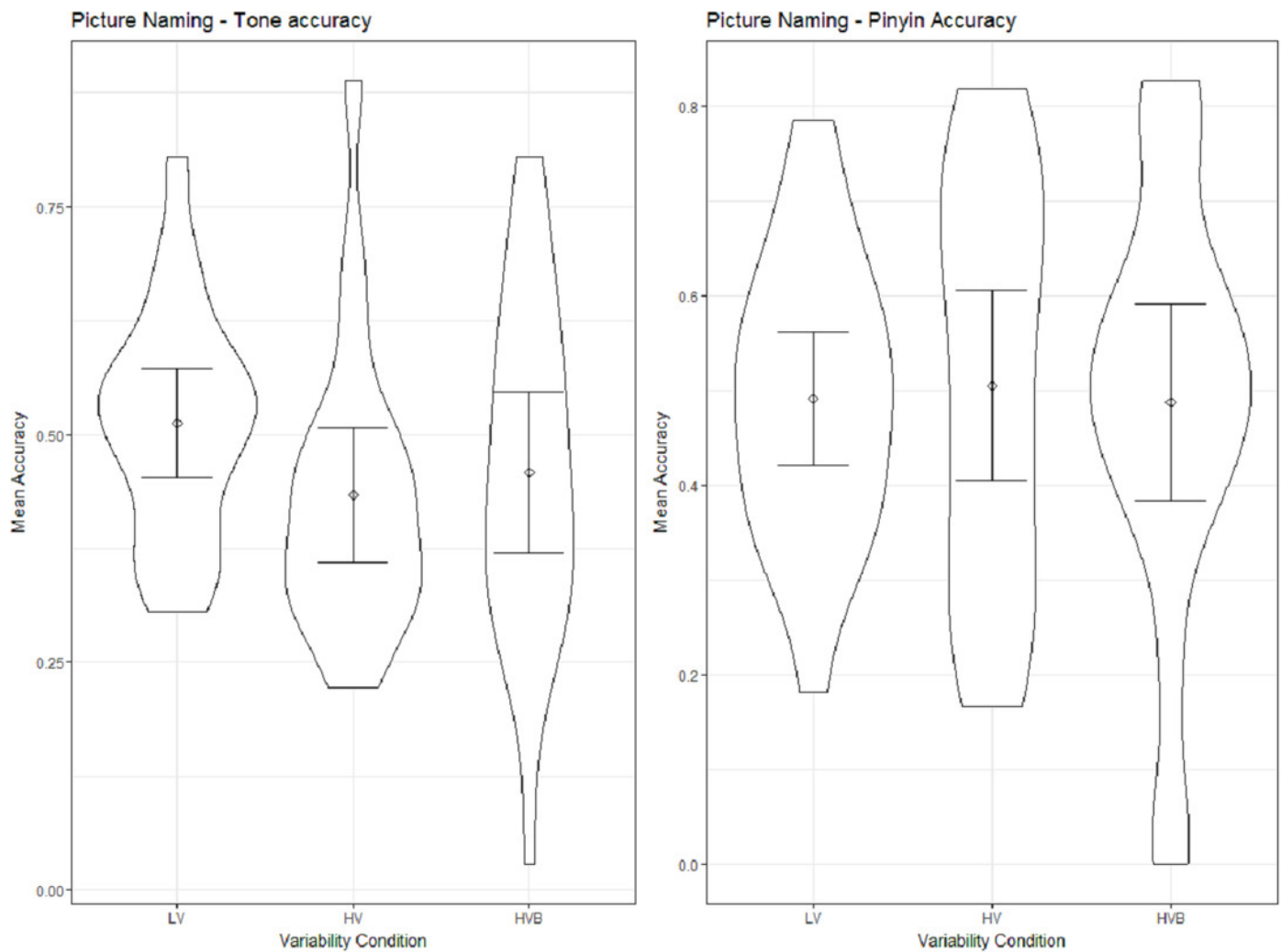
## Figure 9

Mean tone accuracy and pinyin accuracy of Picture Naming for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups. Error bars show 95% confidence intervals.



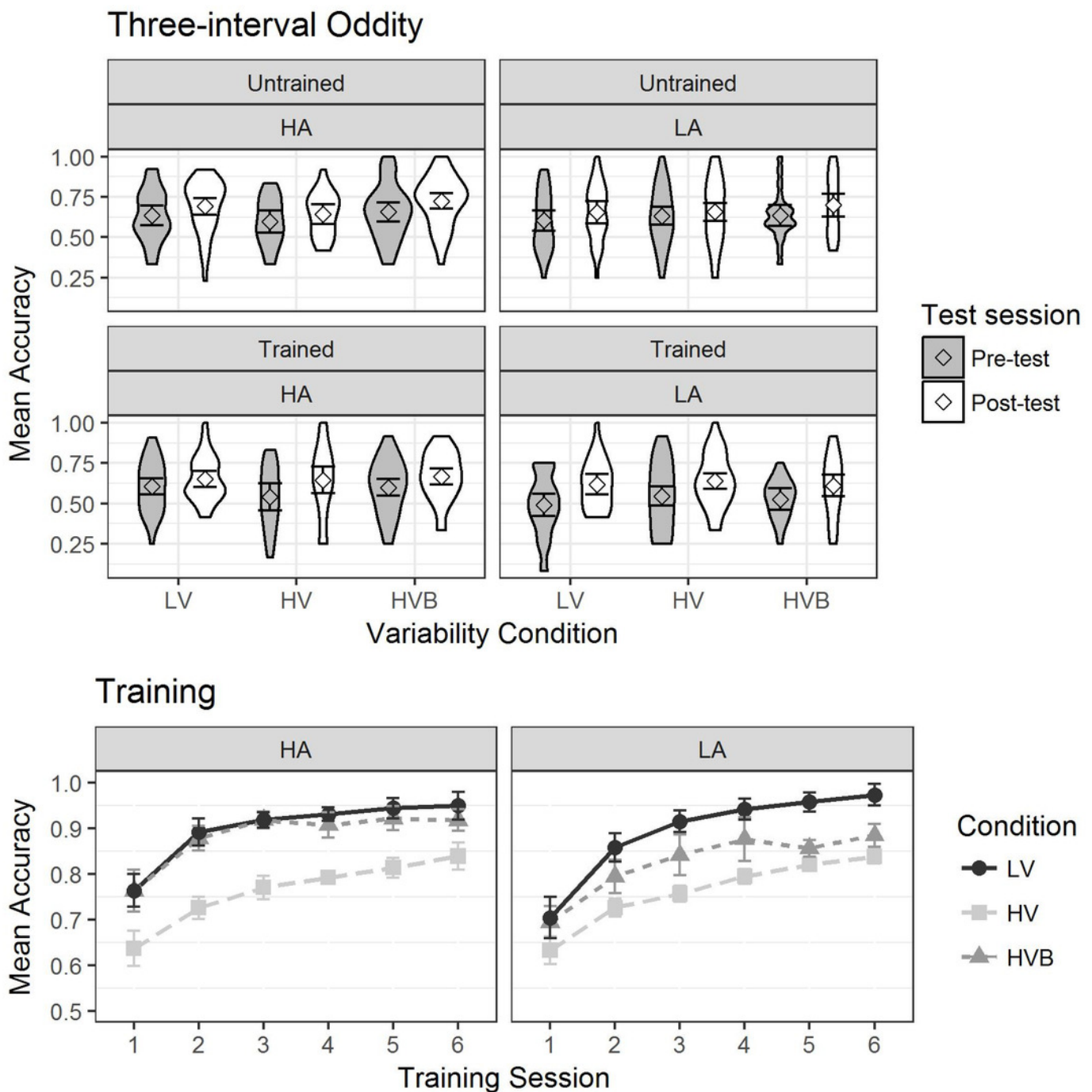
## Figure 10

Violin plot for Tone accuracy and Pinyin accuracy of Picture Naming for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups. Error bars show 95% confidence intervals.



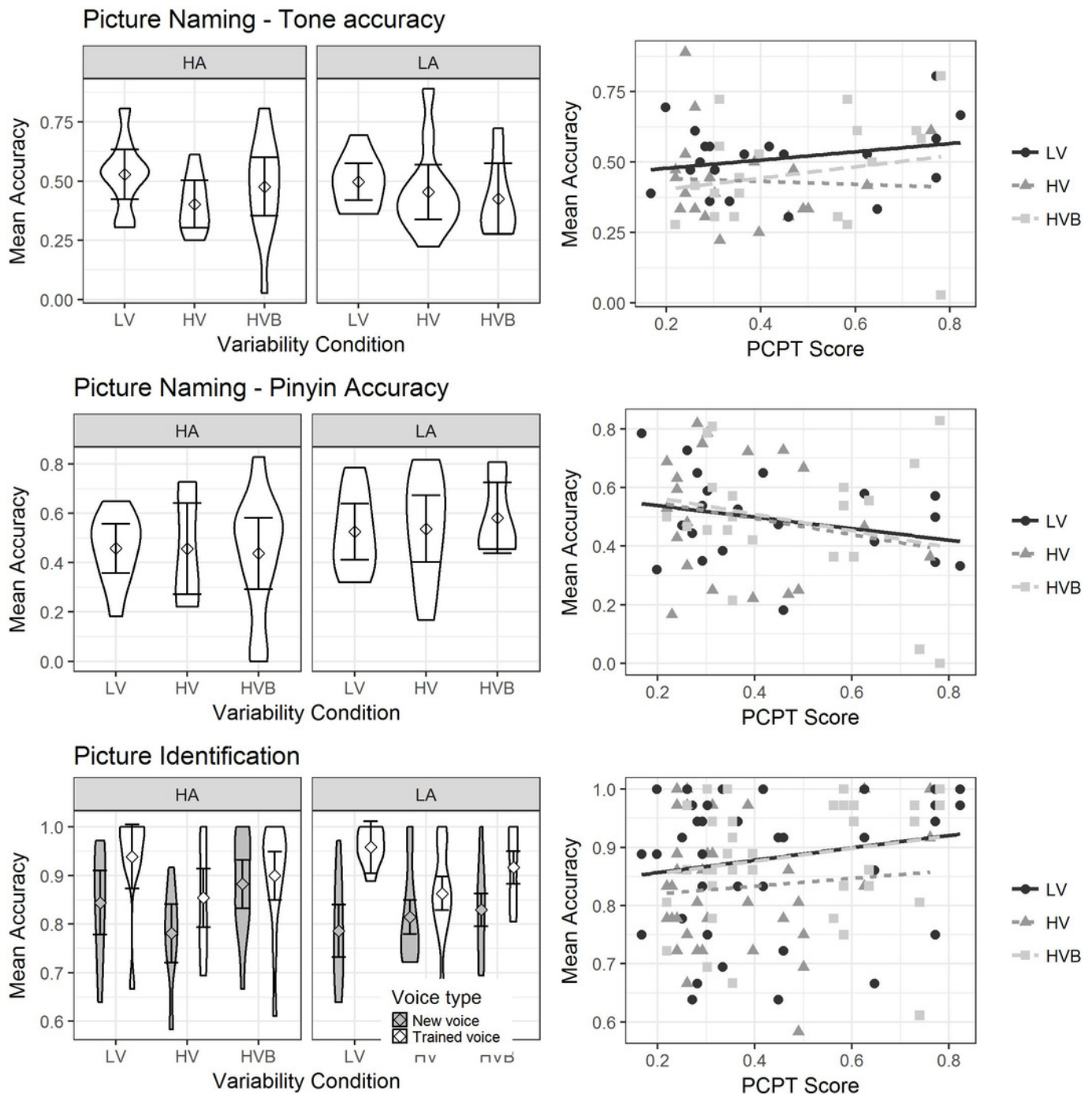
# Figure 11

[i]Accuracy in the Three Interval Oddity and Training data for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups, split by high versus low aptitude in the PCPT task. Error bars show 95% confidence intervals.[i



## Figure 12

[i]Accuracy in the Picture Naming and Picture Identification data for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups, split by high versus low aptitude in the PCPT task. Error bars show 95% confidence inter





## Figure 13

Accuracy in the Word Repetition data for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocking) training groups, split by high versus low aptitude in the PCPT task. Error bars show 95% confidence intervals.

