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

Hanyu Dong [Corresp., 1] , Meghan Clayards [2] , Helen Brown [3] , Elizabeth Wonnacott [Corresp. 1]

[1] Division of Psychology and Language Sciences, University College London, London, United Kingdom

[2] Department of Linguistics, School of Communications Sciences and Disorders, McGill University, Montreal, QC, Canada

[3] 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 to be more effective than low variability training when 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 second language 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-blocked 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 generalisation. All groups improved in both the production and perception of tones which transferred to untrained 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 for a benefit of high-variability 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. Additional Bayes Factor analyses indicated substantial evidence for the null for the hypotheses of a benefit of high-variability in generalisation, however the evidence regarding the interaction was ambiguous. We discuss these results in light of previous findings.

1

2

**The effects of high versus low talker variability and individual aptitude on phonetic**

**training of Mandarin lexical tones**

Hanyu Dong[1], Meghan Clayards[2], Helen Brown[3], & Elizabeth Wonnacott[1]

*[1]Division of Psychology and Language Sciences, University College London, London, UK*

*[2]Department of Linguistics, School of Communications Sciences and Disorders, McGill*

*University, Montreal, QC, Canada*

*[3] Department of Psychology, Nottingham Trent University, Nottingham, UK*

10

Correspondence concerning this article should be addressed to Elizabeth Wonnacott, Division of

Psychology and Language Sciences, Chandler House, 2 Wakefield Street, London, WC1N 1PF.

Email: e.wonnacott@ucl.ac.uk

14

15

16

17

18

19

20                              Acknowledgement

23

24    **Abstract**

25    High variability training has been found to be more effective than low variability training when

26    learning various non-native phonetic contrasts. However, little research has considered whether

27    this applies to the learning of tone contrasts. The only two relevant studies suggested that the

28    effect of high variability training depends on the perceptual aptitude of participants (Perrachione,

29    Lee, Ha, & Wong, 2011; Sadakata & McQueen, 2014). The present study extends these findings

30    by examining the interaction between individual aptitude and input variability using natural,

31    meaningful second language input (both previous studies used pseudowords). Sixty English

32    speakers took part in an eight session phonetic training paradigm. They were assigned to

33    high/low/high-blocked variability training groups and learned real Mandarin tones and words.

34    Individual aptitude was measured following previous work. Learning was measured using one

35    discrimination task, one identification task and two production tasks. All tasks assessed

36    generalisation. All groups improved in both the production and perception of tones which

37    transferred to untrained voices and items, demonstrating the effectiveness of training despite the

38    increased complexity compared with previous research. Although the low variability group

39    exhibited an advantage with the training stimuli, there was no evidence for a benefit of high-

40    variability in any of the tests of generalisation. Moreover, although aptitude significantly

41    predicted performance in discrimination, identification and training tasks, no interaction between

42    individual aptitude and variability was revealed. Additional Bayes Factor analyses indicated

43    substantial evidence for the null for the hypotheses of a benefit of high-variability in

44    generalisation, however the evidence regarding the interaction was ambiguous. We

45    discuss these results in light of previous findings.

46   **1      Introduction**

47      One challenging aspect of learning a second language (L2) is learning to accurately

48   perceive non-native phonetic categories. This task is particularly difficult when the L2 relies on

49   the same acoustic dimensions as the first language (L1), but for different purposes (Bygate,

50   Swain, & Skehan, 2013), suggesting that it is challenging to adjust existing acoustic properties in

51   the L1 to learn new L2 categories. This challenge is compounded by the fact that speech is

52   highly variable in the natural linguistic environment. Variability comes not only from the

53   phonetic context but also from differences between speakers. Thus, learners must learn to

54   distinguish the new L2 categories despite all the variability present in the learning input. There is

55   evidence that native listeners can process this variability in speech faster and more accurately

56   than non-native listeners (Bradlow & Pisoni, 1999), indicating that variability is indeed a

57   challenge for L2 learners. Despite this, it has been suggested that input variability may be

58   beneficial for L2 learning and generalisation (Barcroft & Sommers, 2005; Lively, Logan &

59   Pisoni, 1993). However recent evidence suggests that the ability to benefit from variability may

60   depend on individual learner aptitude (Perrachione et al., 2011; Sadakata & McQueen, 2014), at

61   least in the learning of lexical tones (i.e. the distinctive pitch patterns carried by the syllable of a

62   word which, in certain languages, distinguish meaningful lexical contrasts). The current paper

63   further explores how and when variability supports or impedes learning of new L2 phonetic

64   categories, focusing on English learners of Mandarin tone contrasts.

65   *1.1      High Variability L2 Phonetic Training for Non-Tonal Contrasts*

66      A substantial body of literature has explored whether phonetic training can be used to

67   improve identification and discrimination of non-native phonetic contrasts in L2 learners. An

68   early study by Strange and Dittman (1984) attempted to train Japanese speakers on the English

69    /r/- /l/ distinction, a phoneme contrast that does not exist in Japanese. Participants were trained

70    on stimuli from a synthetic *rock-lock* continuum. The key result was that although performance

71    increased both for trained and novel synthetic items, participants failed to show any

72    improvement for naturally produced minimal pair items. Later research suggested that a key

73    factor which prevented generalisation to natural speech tokens was a lack of variability in the

74    training materials: Variability was present in the form of the ambiguous intermediate stimuli

75    along the continuum, however, there was a single phonetic context and a single (synthesised)

76    speaker. Logan, Lively, and Pisoni (1991) also trained Japanese learners on the English /r/-/l/

77    contrast, but included multiple natural exemplars spoken by six speakers, with the target speech

78    sounds appearing in a range of phonetic contexts. In contrast to Strange and Dittman, they found

79    that participants successfully generalised both to new speakers and new words at test. This was

80    the first study to indicate the importance of variability within the training materials. A follow up

81    study by Lively et al. (1993) provided further evidence for this by contrasting a condition with

82    *high variability* input to one with *low variability* input in which the stimuli were spoken by a

83    single speaker (although still exemplified in multiple phonetic contexts). Participants in the low

84    variability group improved during the training sessions but failed to generalise this learning to a

85    new speaker.

86         Following Lively et al. (1993) *high variability phonetic training* (HVPT) has become

87    standard in L2 phonetic training. This methodology has been successfully extended to training a

88    variety of contrasts in various languages such as learning of the English /u/-/ʊ/ distinction by

89    Catalan/Spanish bilinguals (Aliaga-García & Mora, 2009), learning of the English /i/-/ɪ/ contrast

90    by native Greek speakers (Giannakopoulou, Uther & Ylinen, 2013; Lengeris & Hazan, 2010),

91   and learning of the English /w/-/v/ distinction by native German speakers (Iverson, Ekanayake,

92   Hamann, Sennema, & Evans, 2008).

93       There is also some evidence that this type of perceptual training benefits production in

94   addition to perception. Bradlow, Akahane-Yamada, Pisoni and Tohkura (1999) found that

95   production of the /r/-/l/ contrast improved in Japanese speakers following HVPT, with this

96   improvement being retained even after three months. Similar improvement on the production of

97   American English mid to low vowels by Japanese speakers following HVPT was also reported

98   by Lambacher, Martens, Kakehi, Marasinghe, and Molholt (2005). However, the evidence here

99   is mixed: A recent study by Alshangiti and Evans (2014) employed HVPT to train Arabic

100   learners on non-native English vowel contrasts and found no improvements in production,

101   although participants receiving additional explicit production training did show some limited

102   improvement.

103       Although the studies reviewed above all used HVPT, only the original work by Logan

104   and colleagues directly contrasted the use of high and low variability materials. It is notable these

105   seminal experiments used small samples (the tests of generalisation were administered to only

106   three of the participants in Logan et al., 1991). Since then, few studies have explicitly contrasted

107   high and low variability training. One such study was Sadakata and McQueen (2013), who

108   trained native Dutch speakers with geminate and singleton variants of the Japanese fricative /s/.

109   Participants were trained with either a limited set of words recorded by a single speaker (low

110   variability) or with a more variable set of words recorded by multiple speakers (high variability).

111   Both types of training led to increases in generalisation to untrained fricatives and speakers.

112   However, in an identification task, the improvement was greater for participants receiving high

113   variability training than those receiving low variability training. Similar results were reported by

114    Wong (2014) who trained native Cantonese speakers with the English /e/ - /æ/ contrast. Both low

115    variability (1 speaker) and high variability (6 speakers) training lead to increased performance

116    from pre- to post- test, but the improvement was greater for the high variability group. This was

117    found in tests of generalisation to new speakers and new items, and from perception to

118    production. In contrast, a recent phonetic training study did not find the same benefit.

119    Giannakopoulou, Brown, Clayards, and Wonnacott (2017) compared matched high variability

120    (four speakers) and low variability (one speaker) training for adult and child (8-year-old) native

121    Greek speakers who were trained on the English /i/-/ɪ/ contrast. This study did *not* show a benefit

122    for high variability compared to low variability training in either age group, even for

123    generalisation items. However, for adult participants, it is unclear the extent to which this was

124    due to ceiling effects. To our knowledge, the only other previous studies that specifically

125    manipulated variability during learning of non-native phonetic categories are those by

126    Perrachione et al. (2011), and Sadakata and McQueen (2014), which both looked at the learning

127    of lexical tone. We discuss these studies in more detail in the following section.

128         Although there is a relatively small evidence base regarding a benefit of high over low

129    phonetic training for non-native phoneme categories, there is further evidence for this benefit in

130    related areas of speech and language learning, specifically accent categorisation and adaptation

131    (Bradlow & Bent, 2008; Clopper & Pisoni, 2004), and L2 vocabulary learning (Barcroft &

132    Sommers, 2005, 2014; Sommers & Barcroft, 2007, 2011). Benefits of HVPT are generally seen

133    in tasks of generalisation, suggesting that exposure to variation across speakers and/or items

134    boosts the ability to generalise across these dimensions. This intuitively sensible result is in line

135    with the predictions of computational models in which irrelevant contextual/speaker identity

136    cues compete with phonetically relevant cues, so that dissociation of these irrelevant cues is the

137   key mechanism which underpins generalisation (Apfelbaum & McMurray, 2011; Ramscar &

138   Baayen, 2013; Ramscar, Yarlett, Dye, Denny & Thorpe, 2010).

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

140         Each of the phonetic training studies discussed above involved training a *segmental*

141   contrast (consonantal or vocalic). Lexical tone is another type of phonological contrast in some

142   natural languages, whereby the pitch contour is used to distinguish lexical or grammatical

143   meanings (Yip, 2002). For example, Mandarin Chinese has four lexical tones: level-tone (Tone

144   1), rising-tone (Tone 2), dipping-tone (Tone 3) and falling-tone (Tone 4). These pitch contours

145   combine with syllables to distinguish meanings. For instance, the syllable *ba* combines with the

146   four tones to mean: *eight* (*bā*, Tone 1), *pluck* (*bá*, Tone 2), *grasp* (*bǎ*, Tone 3) and *father* (*bà*,

147   Tone 4). Each of these words thus forms a minimal pair with each of the others. Note that while

148   non-tonal languages such as English use pitch information extensively for intonation (e.g.

149   forming a question, or for emphasis), and that pitch plays a role in marking stress at the lexical

150   level in (e.g. '*import*/*im'port*), this is quite different from a lexical tone system, causing

151   difficulties for L2 learners of Mandarin.

152         The first study examining lexical tone training was conducted by Wang, Spence,

153   Jongman, and Sereno (1999). A similar paradigm to that used by Logan et al. (1991) was

154   adopted using four speakers for training. Training materials were all real monosyllabic Mandarin

155   words that varied in the consonants, vowels and syllable structure. During training participants

156   heard a syllable whilst viewing two of the four standard diacritic representations (i.e., →, ↗, ∨, ↘,

157   which are iconic in nature). They were asked to pick out the picture of the arrow that

158   corresponded to the tone. At test, participants chose which tone they had heard out of a choice of

159   all four diacritics. There were also two generalisation tasks, one testing generalisation to

160    untrained items and one testing generalisation to a new speaker. Native American English

161    speakers showed significant improvement in the accuracy of tone identification after eight

162    sessions of high variability training over two weeks, and this generalised to both new words and

163    a new speaker. In a follow up study, Wang, Jongman and Sereno (2003) used the same training

164    paradigm to test whether learning transferred to production. They recruited participants taking

165    Mandarin courses and asked them to read through a list of 80 Mandarin words written in Pinyin

166    (an alphabetic transcription) before and after training. They found improvements in production,

167    although these were mainly seen in pitch contour rather than pitch height.

168        These studies suggested that as with segmental phoneme contrasts, high variability

169    training may also facilitate the learning of tone contrasts. However, Wang and colleagues (1999,

170    2003) did not directly contrast high and low variability training materials. Perrachione et al.

171    (2011) investigated this contrast directly. They trained native American English speakers with no

172    previous knowledge of Mandarin (or any other tonal language), using English monosyllabic

173    pseudowords combined with Mandarin tones 1 2, and 4 (→, ↗ & ↘). The training task used

174    either low variability (one speaker) or high variability (four speaker) input. During the training,

175    participants matched the sound they heard with one of three pictures of concrete objects

176    presented, where the three words associated with these pictures were minimal trios that differed

177    only in tone. Participants were tested on their ability to generalise their learning to new speakers.

178    Importantly, Perrachione et al. (2011) were also interested in the role of individual differences in

179    learning. Therefore, they also determined participants' baseline ability to perceive the tone

180    contrasts prior to training using a *Pitch Contour Perception Test.* In this task, participants heard a

181    vowel produced with either Mandarin tone 1, 2 or 4 whilst viewing pictures of standard diacritics

182    associated with these tones (→, ↗ & ↘), and were asked to select the arrow that corresponded to

183    the tone. Based on performance in this task, the researchers grouped participants into high and

184    low aptitude groups. The results showed that whilst the low variability group outperformed the

185    high variability group during training (presumably due to accommodation to a repeated speaker

186    throughout the task), there were no differences between the high and low variability groups

187    during test. Critically however, there was an interaction between an individuals' aptitude

188    categorisation and the type of variability training: Only participants with high aptitude benefitted

189    from high variability training, while those with low aptitude actually benefitted more from low

190    variability training. It is important to note that this interaction was seen in a task which relied on

191    participants' ability to generalise their learning[1] of tones to an untrained speaker. That is, in a

192    task where we would expect that exposure to multiple speakers would be beneficial since it

193    should allow learners to better dissociate the tones from the particular speakers used in training.

194    These results, therefore, suggest that only the high aptitude learners can take advantage of this

195    benefit. Another training study by Sadakata and McQueen (2014) also explored the relationship

196    between input variability and individual aptitude in lexical tone training, though using different

197    training and testing materials. They trained native Dutch speakers (with no prior knowledge of

198    Mandarin or any other tonal language) using naturally produced bisyllabic Mandarin

199    pseudowords. The two syllables in each word either had Tone 2 followed by Tone 1, or Tone 3

200    followed by Tone 1, and each tone pair was randomly assigned one of two numeric labels (e.g.

201    for one participant Tone 2-Tone 1 was labelled "1", Tone 3-Tone 1 was labelled "2"). During the

202    training task, participants identified the tone pair type of each stimulus by choosing the correct

[1] In their paper, Perrachione et al (2011) do *not* refer to this task as a generalisation task. Instead they report a generalisation measure which is a ratio of performance on this test with novel speakers to performance in training (test-performance/training-performance). Note that this ratio will increase not only if participants are better at test, but also if they are *worse* in training. Using this measure, Perrachione et al. found a benefit of high variability training. However on inspection of the means, it seems that this relationship is driven by the *poorer* performance in training in the high variability condition, rather than by *better* performance in the test with novel speakers. We therefore do not see the ratio measure as providing evidence for an overall benefit of HV training on generalisation.

203    numeric label (e.g. hear /pasa/ with Tone 2-Tone 1, correct response is 1). Thus, in contrast to the

204    study by Perrachione et al. (2011), participants did not need to learn the meaning of each word.

205    Input variability was manipulated, with three levels (low/medium/high). In contrast to the work

206    by Perrachione et al., where the high variability and low variability conditions differed only in

207    terms of the number of speakers, in this study variability was increased both by including more

208    speakers and more items. Specifically, the number of different vowels used in the bi-syllabic

209    sequences was manipulated: the low variability group encountered only one vowel (.e.g. pasa,

210    casa, lasa, etc.) whereas the medium and high variability groups encountered four different

211    vowels (pasa, pesa, pisa, pusa; casa, cesa, cisa, cusa; lasa, lesa, lisa, lusa etc.). Participants were

212    tested on the trained items (i.e. using trained speakers and trained items). Generalisation was also

213    examined in a number of ways by looking at (1) trained items spoken by an untrained talker; (2)

214    pseudowords containing untrained vowels (3) pseudowords in which the order of tones in the bi-

215    syllables were reversed (i.e. a novel position), and (4) items where the tone was embedded in a

216    sentence context. As in the study by Perrachione et al. (2011), Sadakata and McQueen (2014)

217    also tested individual aptitude but with a different method. They employed a categorisation task

218    using stimuli from a six step Tone 2 to Tone 3 continuum (created using natural productions of

219    the two tones with the Mandarin vowel /a/ as endpoints and linearly interpolating between these

220    endpoints). Participants were asked to identify if the sound they heard was more like Tone 2 or

221    Tone 3, and a categorisation slope was obtained for each participant providing a measure of their

222    ability to discriminate this contrast, which is generally found to be the most challenging tone

223    contrast for L2 learners of Mandarin. Participants were grouped according to their slopes, and

224    this grouping was entered as a factor in the analyses of tests of learning, along with the effect of

225    training condition (high-medium-low) and the interaction between factors. For the test with

226   trained speakers and items, there was no group level effect of variability condition, however

227   there *was* an interaction between variability and aptitude similar to that reported by Perrachione

228   et al.: Participants with high aptitude benefitted from high variability training, while those with

229   lower aptitude benefitted more from low variability training. For the generalisation tests,

230   participants showed above chance performance in all but the new position condition,

231   demonstrating an ability to generalise their learning of tone across different dimensions.

232   However, they did *not* demonstrate an overall benefit of higher variability in any of the transfer

233   tests, nor, did variability interaction with aptitude. Note that the overall lack of a high variability

234   benefit is again surprising, particularly for test items with untrained talkers and novel items,

235   since the manipulations in training should specifically work to increase generalisation along

236   these dimensions.

237        In sum, the two studies which have directly compared high and low variability input in

238   training Mandarin tone contrasts have *not* found the predicted benefit of high variability on

239   generalisation, either when varying just speakers or when varying speakers and items. However,

240   both of these studies found an interaction between participant aptitude and variability condition.

241   The results of these studies thus provide mutually corroborating evidence – using somewhat

242   different training and testing methods – that the ability to learn from high variability input is

243   dependent on learner aptitude, although it should be noted that this interaction was found in a

244   task with untrained speakers in one study (Perrachione et al., 2011), but in a task with trained

245   stimuli in the other (Sadakata & McQueen, 2014).

246        Why might the ability to benefit from varied training materials depend on participant

247   aptitude? Perrachione et al. (2011) suggest that one reason why low aptitude participants may

248   struggle with multi-speaker input is that the speakers were intermixed during training: This

249 requires trial-by-trial adaptation to each speaker, which was not required in the corresponding

250 single speaker low variability conditions. This may place a burden on learners (see Mattys &

251 Wiget, 2011; Nusbaum & Morin, 1992, for evidence that intermixed multi-speaker stimuli are

252 difficult even for L1 processing and that this interacts with constraints on working memory and

253 attention). To test this, Perrachione et al. (2011) conducted a second experiment in which items

254 from each speaker were presented in separate blocks (as is more common in HVPT). This

255 improved performance during the training task compared with unblocked training for low

256 aptitude learners only, confirming the hypothesis that switching between speakers on a trial-by-

257 trial basis during training interferes with learning for low aptitude learners. On the other hand,

258 Sadakata and McQueen (2014) employed a training paradigm in which speakers were blocked in

259 the high variability condition, yet they still found the interaction with aptitude. However, recall

260 that in their experiment they also manipulated item variability, yet only speakers were blocked

261 by session, not items. Thus, it remains possible that trial-by-trial inconsistency at the level of

262 items could explain some of the greater difficulty of low aptitude learners in their study.

263 *1.3 The Current Study*

264 The fact that neither of the tone training studies found an overall benefit of high over low

265 variability in tone generalisation is surprising in light of the phonetic literature and the

266 predictions of the computational model (Apfelbaum & McMurray, 2011; Ramscar & Baayen,

267 2013) mentioned above. Moreover, as the previous authors point out, if it is actually the case that

268 learning from multiple voices is more or less effective for different groups of learners, this has

269 important implications for the design of L2 training tools. For this to be the case, it is important

270 to establish the generalizability of the findings to different contexts and materials, particularly

271 those which are relevant in an L2 learning context. We suggest that what L2 learners are most

272    interested in developing is their ability to use tone when mapping a word's phonological form to

273    its meaning (and vice versa). In this light, the paradigm used by Sadakata and McQueen (2014)

274    lacks ecological validity in looking only at mapping to abstract tone categories. On the other

275    hand, Perrachione et al. (2011) do train form-meaning mappings, yet, unlike Sadakata and

276    McQueen (2014) they use English pseudo-word stimuli, which has the consequence that learners

277    do not simultaneously have to deal with non-native segments and tones, as in a real world L2

278    learning situation. Furthermore, although there is limited data on the differences between words

279    and non-words in production, it has been noted that non-words may have different properties

280    from real words even within the same language (Scarborough ,2012) and may be more clearly

281    articulated (Hay, Drager & Thomas, 2013; Maxwell, Baker, Bundgaard-Nielsen & Fletcher,

282    2015). Thus, using non-words might make stimuli slightly easier to learn than if real words were

283    used.

284         The current training study addresses these issues in a partial replication of the previous

285    work: We use stimuli produced by native Mandarin speakers which are real words in that

286    language. This design choice follows earlier studies such as Wang et al. (1999) using a paradigm

287    in which participants are trained to identify word meaning on the basis of tone. In contrast to the

288    previous studies, we also trained the contrasts between all four tones (six tone contrasts) rather

289    than just three (on the assumption that learners are interested in learning the complete set of

290    contrasts within a particular language). We note that these design choices potentially increase the

291    difficulty of our training materials compared to previous work. A key question was whether

292    these choices would impact the interaction between learner aptitude and the benefits of more

293    variable training materials.

294    We followed Perrachione et al. (2011) in varying variability along one dimension only –

295    speaker variability, keeping training items identical across conditions. We also followed

296    Perrachione et al. (2011) in comparing high variability input which was blocked by speaker, with

297    input that was not, making three training conditions: low variability  (one speaker), high

298    variability (four speakers intermixed within each training session) and blocked training (four

299    speakers each presented in separate blocks). Note that our choice to manipulate only talker-

300    variability means that the high variability blocked condition is matched to the low variability

301    condition in terms of trial-by-trial inconsistency, unlike in Sadakata and McQueen (2014) where,

302    even though they blocked by speaker, the high variability condition contained more trial-by-trial

303    variability in terms of items. We predicted that the difficulty of high variability input for lower

304    aptitude participants would be greater in the unblocked condition, thus potentially increasing the

305    likelihood of seeing the predicted interaction between variability and learner aptitude. On the

306    other hand, blocked input is more usual of HVPT (e.g. Iverson, Hazan & Bannister, 2005; Logan

307    et al. 1991) and may increase the possibility of seeing an overall benefit of speaker variability on

308    generalisation.

309    We used two perceptual tasks designed to tap individual aptitude. These were adapted

310    from those used in Perrachione et al. (2011) and Sadakata and McQueen (2014). However, while

311    the previous studies grouped participants into one of two categories (high aptitude *vs.* low

312    aptitude) based on the aptitude score, in the current study they were used as continuous

313    measures. This allowed us to avoid assigning an arbitrary "cut off" for high *v*ersus low aptitude

314    groups, and the loss of information which occurs when an underlying continuous variable is

315    turned into a binary measure. Note that the statistical approach used in the current paper (logistic

316    mixed effect models) allowed us to include continuous predictors and look at their interactions

317    with other factors.

318         A further extension in the current study is that we use several new outcome measures to

319    test learning and generalisation. First, most similar to the task used in Perrachione et al. (2011)

320    was a picture identification task which was a version of the training task (2AFC picture

321    identification) without feedback. Following Perrachione et al. (2011) we included untrained-

322    speaker items, where benefits of speaker variability in training should be most apparent.

323    However, bearing in mind that Sadakata and McQueen (2014) actually found the key interaction

324    with aptitude *only* in the test with trained stimuli, we also included trained-speaker test items.

325         We also included a second perceptual task which did *not* involve knowing specific form-

326    meaning mappings and thus had the benefit that it could be conducted both pre- and post- test.

327    This was a three interval oddity task which required participants to pick the odd-one-out after

328    hearing three words spoken aloud, each by a different speaker. Two of the tokens were

329    productions of the same word and the third differed only in the tone (e.g. *bā,* Tone 1; *bā,* Tone 1;

330    *bà*, Tone 4). Because all three tokens are physically different, it requires the listener to focus on

331    the phonological level ignoring irrelevant acoustic differences. Furthermore, the use of three

332    speakers forces the listener to ignore irrelevant speaker-specific differences, making it especially

333    challenging (Strange & Shafer, 2008). This task used untrained speakers in every trial, so that

334    every test-item required generalisation to new speakers[2]. In addition, here it was possible to use

335    both trained and untrained *items*. Note that even though the variability over items is matched

336    across conditions, it is possible that varying speaker specific cues might also thus promote

_____

[2] If we wished to use trained speakers, in order to be able to the use the same test with the low variability condition, we would have to use a single speaker across all three test trials. Our pilot work suggested that participants performed at ceiling on a single-speaker version of this task, even at pre-test.

337  generalisation across this dimension. If this is the case, a high variability benefit may be stronger

338  for untrained items than trained items.

339       Finally, we also tested production using a picture naming task at post-test, in which

340  participants were required to name the pictures used in training in Mandarin. We also conducted

341  a word repetition task, which had the benefit that it could also be conducted at pre-test, and that

342  we could use both trained and untrained words (as for the three-interval oddity task discussed

343  above). Although there is evidence HVPT can benefit the production of tones (Wang et al.,

344  2003), there has been no direct examination of whether high variability training materials are

345  more effective than low variability training materials for production. However, more generally in

346  the L2 vocabulary learning literature, training with multiple speakers has been found to lead to

347  better recall in a picture naming task (Barcroft & Sommers, 2005), suggesting that the HVPT

348  advantage should extend to production measures.

349       In sum, the current experiment assessed whether individuals benefit from high over low

350  variability perceptual training when learning novel L2 tone contrasts, and whether this interacts

351  with learner aptitude. We used measures of aptitude taken from previous studies, but a training

352  paradigm with real Mandarin stimuli embedded in a vocabulary learning task, which trained

353  discrimination of all six Mandarin tone contrasts. Learning and generalisation were measured in

354  multiple tests of both perception and production. In general, the current design increased

355  ecological validity and likely also increased the difficulty of the learning task relative to previous

356  work. It is possible that increasing difficulty could exacerbate differences between learners of

357  different aptitudes, potentially increasing the effect. On the other hand, it is also possible that the

358  increased difficulty might make high variability input much harder for all participants,

359  decreasing or removing the specific benefit of HVPT for high aptitude learners.

360 **2     Method**

361 *2.1     Participants*

362      Sixty adults recruited from UCL Psychology Subject Pool participated in the experiment,

363 twenty in each of the three conditions (low variability, high variability, high variability blocked).

364 Participant information is summarised in *Table 1*. There was no difference between these groups

365 in age, $F$ (2,57) = 1.95, $p$ = .15. Participants had no known hearing, speech, or language

366 impairments. Written consent was obtained from participants prior to the first session. Each

367 participant was paid £45 at the end of the study.

368      All participants except three were native English speakers. Of the remaining three, one

369 participant (low variability condition) was a native bilingual of English and Hindi, one participant

370 (high variability condition) was a native French speaker, and one participant (high variability

371 condition) was a native Finnish speaker. Critically, participants had no prior experience of

372 Mandarin Chinese or any other tonal language. On average, participants had learned 2.4 ($SD$ = 0.8)

373 languages and the average age for starting to learn the first L2 was 12.6 years ($SD$ = 1.3).

374      Ethical approval was given by the UCL Research Ethics Committee with the approval

375 number 6176/002.

376 *2.2     Stimuli*

377 *2.2.1   Stimuli used in Training and in the Picture Identification, Three Interval Oddity, Word*

378      *Repetition and Picture Naming Tests*

379      These stimuli consisted of 36 minimal pairs of Mandarin words (6 minimal pairs for each

380 of the six tone contrasts generated by the four Mandarin tones). The words in each pair contained

381 the same phonemes, differing only in tone (e.g. *māo*, Tone 1 [*cat*] vs. *mào*, Tone 4 [*hat*]). All

382 words were picturable and started with a wide range of phonemes (see Appendix A). In order to

383  examine generalisation across items, half of the word pairs (3 per tone contrast) were designated

384  "trained" words and other half were designated "untrained" words. Trained words were

385  encountered in both training and test tasks; untrained words were only encountered in the three

386  interval oddity and word recognition tests.

387      The full set of 72 Mandarin words was recorded by two groups of native Mandarin

388  speakers using a Sony PCM-M10 handheld digital audio recorder. The first group consisted of

389  three female and two male speakers. These stimuli were used in the Training, Word Repetition

390  and Picture Identification tasks. The second group consisted of three new female speakers and

391  two new male speakers. These stimuli were used in the three interval oddity task (making all new

392  speakers in that task). See Table 2 for a summary of the manipulation of item and speaker

393  novelty across the different test tasks, and Table 3 for the tasks in which speakers are

394  counterbalanced.

395      In the low variability condition only one speaker (Trained voice 1) was used in training,

396  and this same speaker was also used as the test voice in the Word Repetition test and for trained

397  items in the picture identification test. In the high variability conditions, four speakers (Trained

398  voice 1 plus three others) were used in training. Only one of these speakers (Trained voice 1)

399  was used in the word repetition test and for trained items in the picture identification test. In all

400  conditions, a further speaker (Untrained voice 1) was assigned to the untrained test items in the

401  picture identification test. The assignment of speakers was rotated across participants, resulting

402  in five counterbalanced versions of each condition (see Table 3). This ensured that any

403  difference found between the low and high variability conditions, and between trained and

404  untrained voices, were not due to idiosyncratic difference between speakers. There was no

405  counterbalancing of speaker in other tasks.

406        All words were edited into separate sound files, and peak amplitude was normalized

407    using Audacity (Audacity team, 2015, http://audacity.sourceforge.net/). Any background noise

408    was also removed. All recordings were perceptually natural and highly distinguishable as judged

409    by native Chinese speakers. Clipart pictures of the 72 words were selected from free online

410    clipart databases.

411    *2.2.2   Stimuli used in the Aptitude Tests:*

412        Pitch Contour Perception Test: Six Mandarin vowels (/a/, /o/, /e/, /i/, /u/, /y/) were

413    repeated in the four Mandarin tones by two male and two female native Mandarin speakers from

414    talker set 2, making 96 stimuli in total. Stimuli were identical across conditions and participants.

415        Categorisation of Synthesised Tonal Continua: Natural endpoints were chosen from a

416    native Mandarin male speaker producing the word '*wan*' with both Tone 2 and Tone 3. A neutral

417    vowel was also recorded by a native male English speaker producing the 'father vowel' /a/. This

418    vowel was edited slightly to remove portions containing creaky voice at the end. The three

419    syllables (wan [Tone 2], wan [Tone 3], /a/) were then manipulated in Praat (Boersma &

420    Weenink, 2015). All three syllables were normalized to be approximately 260 ms long using the

421    PSOLA method. The neutral vowel was manipulated to have a flat fundamental frequency (148

422    Hz) and a flat intensity contour (75dB). The pitch contours of the two natural endpoints were

423    extracted and a 6-step pitch continuum (Step 1: Tone 2, Step 6: Tone 3) was generated by

424    linearly interpolating between the endpoints. These six pitch contours were then each

425    superimposed on a copy of the neutral vowel using the PSOLA method. Stimuli were identical

426    across participants and conditions.

427    **2.3    Procedure**

428    The experiment involved three stages (see *Figure 1*): Pre-test (session 1), training

429    (sessions 2-7), and post-test (session 8). Participants were required to complete all eight sessions

430    within two weeks, with the constraint of one session per day at most. The majority of sessions

431    took place in a quiet, soundproof testing room in Chandler House, UCL. The remaining sessions

432    took place in a quiet room in a student house.

433    Participants were given a brief introduction about the aim of the study and told that they

434    were going to learn some Mandarin tones and words. They were explicitly told that Mandarin

435    has four tones (flat, rising, dipping and falling) and that the tonal differences were used to

436    distinguish meanings. The experiment ran on a Dell Alienware 14R laptop with a 14-inch screen.

437    The experiment software was built using a custom-built software package developed at the

438    University of Rochester.

439    The specific instructions for each task were displayed on-screen before the task started.

440    After each task, participants had the opportunity to take a 1-minute break. The tasks completed

441    in each session are listed in *Figure 1* and described in more detail below. Note that the Pitch

442    Contour Perception Test and Categorisation of Synthesized Tonal Continua were carried out at

443    the beginning of the first session as they provided the measure of individual aptitude prior to

444    exposure to any Mandarin stimuli. There was no time limit for making responses in any of the

445    tasks. Participants wore a pair of HD 201 Sennheiser headphones throughout the experiment with

446    audio stimuli presented at a comfortable listening level.

447    *2.3.1    Individual Aptitude Measures*

448    *2.3.1.1    The Pitch Contour Perception Test*

449    This test was based on the work of Wong and Perrachione (2007). Participants heard a

450    tone (e.g. /a/ [Tone 1]), while viewing pictures of four arrows indicating the different pitch

451    contours. Participants clicked on the arrow that they thought matched the tone heard. No

452    feedback was provided. There were 96 stimuli in total (4 speakers * 4 tones * 6 vowels). This

453    task provided another measure of individual differences in tone perception prior to training.

454    Although Perrachione et al. only conducted this task at pre-test, for consistency with the

455    Categorisation of Synthesised Tonal Continua (described below) we also repeated the test at

456    post-test and conducted analyses to identify whether performance on this task was itself

457    improved as a result of training (see Section 3.2.2).

458                    *2.3.1.2   Categorisation of Synthesised Tonal Continua*

459            This test was based on Sadakata and McQueen (2014). Participants first practiced

460    listening to Tone 2 and Tone 3 while viewing the corresponding picture of an arrow depicting the

461    pitch change. Each tone was repeated 10 times. In each test trial, participants then decided

462    whether the sound they heard was closer to Tone 2 or Tone 3 by clicking on the corresponding

463    arrow. No feedback was provided. The speech continua consisted of 6 steps (Step 1: Tone 2,

464    Step 6: Tone 3) with each step repeated 10 times per block. Participants completed two blocks,

465    with an optional 1 minute break in the middle, resulting in 120 trials in total. This task provided a

466    measure of individual differences in tone perception prior to training. In line with Sadakata and

467    McQueen's procedure, participants completed the task both before and after training and we

468    conducted analyses to explore whether there was improvement from pre to post-test (Section

469    3.2.1).

470    *2.3.2   Training Task*

471            Participants completed the training task in Session 2-7. On each trial, participants heard a

472    Mandarin word and selected one of two candidate pictures displayed on the computer screen.

473    The two pictures always belonged to the same minimal pair. Feedback was provided about

474    whether the answer was correct (a green happy face appeared) or incorrect (a red sad face

475    appeared). If the correct choice was made, a picture of a coin also appeared in a box on the left-

476    hand side of the screen, with the aim of motivating participants to try to earn more coins in each

477    subsequent session of training. After that, everything but the correct picture was removed from

478    the screen and the participant heard the correct word again. In the lower right corner of the

479    screen a trial indicator of X/288 was displayed where X indicated the number of trials completed.

480    This tool helped participants to keep track of their performance (see *Figure 2*).

481        There were 18 picture/word pairs used. Each word was used as the target four times.

482    Thus, each picture pair appeared eight times, resulting in 288 trials per session. Participants were

483    assigned to one of the following conditions: low variability, high variability and high variability

484    blocked (with the assignment of speakers counterbalanced – see Table 3). Each training session

485    lasted for approximately 30 minutes.

486        In the low variability condition, only *one* speaker was used. In the high variability

487    conditions, *four* speakers were used. For each participant, each of their six training sessions was

488    identical. In the high variability condition without blocking, all of the speakers were heard in

489    each of the training sessions, with the order randomized so that speaker varied from trial to trial.

490    In contrast, in the high variability blocked condition, from Day 1 to Day 4 of training (i.e.,

491    Session 2-5), only one speaker was involved on each day's training session, (with the trained

492    speaker that was used in the test tasks (e.g. F1 for Version 1) always occurring on Day 3 (i.e.,

493    Session 4)); on Days 5 and 6 of training (i.e., Sessions 6 and 7), participants heard all four

494    speakers, each in a separate block, with each word being repeated twice in each voice on these

495    days. In all three conditions, the order of items was randomized within each session.

496    *2.3.3    Perceptual Tests*

497            *2.3.3.1    Three Interval Oddity Test (pre- post test)*

498        This task required participants to identify the odd one out (i.e. the stimulus with a

499    different tone) from a choice of three Mandarin words, each spoken by a different speaker. Four

PeerJ

500    untrained speakers were used (3 female, 1 male). Each trial used one of the 36 minimal pairs

501    from the main stimuli set (18 trained pairs, 18 untrained pairs). Preliminary work suggested that

502    trials differed in difficulty depending on whether the "different" stimulus was spoken by the

503    single male speaker, or one of the three female speakers. We therefore ensured that there were

504    equal numbers of the following trial types: (i) "Neutral" - all three words were spoken by female

505    speakers (ii) "Easy" - the "different" word was spoken by a male speaker and the other two were

506    spoken by female speakers; (iii) "Hard" - the "different" word was spoken by a female speaker

507    and the other two were spoken by one male speaker and one female speaker. Each of the words

508    in the minimal pair was used once as the target ("different") word, making 72 trials in total.

509        During the task, three frogs were displayed on the screen. Participants heard three words

510    (played with ISIs of 200ms) and indicated which word was the odd one out by clicking on the

511    appropriate frog, which could be in any of the three positions. They could not make their

512    response until all three words had been heard, at which point a red box containing the instruction

513    "Click on the frog that said the different word" appeared at the bottom of the screen. No

514    feedback was provided. Participants completed this task twice – once in the pre-test, and once in

515    the post-test.

516        *2.3.3.2   Picture Identification Test (post- only test)*

517        This task was the same as the training task with the following changes. Firstly, each word

518    was only repeated twice, once by a trained speaker (trained voice 1) and once by an untrained

519    speaker (Untrained voice 1), making 72 trials in total. Secondly, no feedback was given. This

520    task was completed only in the post-test.

521 *2.3.4   Production Test*

522 *2.3.4.1   Word Repetition Test (pre-post test)*

523    All seventy-two Mandarin words from the main stimulus set (18 trained pairs, 18

524 untrained pairs) set were presented one at a time in a randomised order. They were always

525 spoken by the same speaker and this speaker was also used in their training stimuli (training

526 voice 1; see Table 3). After each word, two seconds of white noise was played. This was

527 included to make sure that participants had to encode the stimulus they were repeating and could

528 not access the information in echoic storage (Flege, Takagi & Mann, 1995). Participants were

529 instructed to listen carefully to the word and then to repeat the word aloud after the white noise.

530 Verbal responses were digitally recorded and were later transcribed and rated by native speakers

531 of Mandarin (see Section 3.5.1). This task was completed once in the pre-test and once in the

532 post-test.

533    *2.3.4.2   Picture Naming Test (post-only test)*

534    All 36 pictures from the training words were presented in a randomised order.

535 Participants were instructed to try to name the picture using the appropriate Mandarin word.

536 Verbal responses were recorded and were later transcribed and rated by native Mandarin

537 speakers (see Section 3.5.1). This task was completed only in the post-test.

538 *2.3.5   Other tasks*

539    *2.3.5.1   English Introduction Task*

540    This task was included in the batch of tasks administered at pre-test in case the meaning

541 of some pictures were ambiguous (not all items were concrete nouns – e.g. "*to paint*").

542 Participants saw each of the 36 pictures from the training set presented once each in a random

543 order and heard the corresponding English word. No response was recorded. Participants

544 completed this task only once, at the end of the pre-test session.

545    *2.3.5.2   Questionnaires*

546        Participants completed a language background questionnaire after the experiment.

547    Participants were asked to list all the places they had lived for more than 3 months and any

548    languages that they had learned. For each language the participant was asked: (a) to state how

549    long they learned the language for and their starting age; (b) to rate their own current proficiency

550    of the language.

551    **3      Results**

552    ***3.1      Statistical Approach***

553        Three different sets of frequentist analyses are reported. First, we conducted the analysis

554    on two individual measures Categorisation of Synthesized Tonal Continua (Section 3.2.1) and

555    Pitch Contour Perception Test (Section 3.2.2). The primary aim of these analyses was to ensure

556    that the three groups did not differ at pre-test, however we also looked for possible differences at

557    post-test. Second, separate analyses are reported on data from the tests administered pre- and

558    post- training (i.e. Word Repetition task (Section 3.5.2) and Three Interval Oddity task (Section

559    3.4.1)), the data collected during Training (Section 3.3) and the data from the two tasks

560    administered only at post-test (i.e. the Picture Identification task (Section 3.4.2) and Picture

561    Naming task (Section 3.5.3). These analyses explored the effects of our experimentally

562    manipulated conditions on the various measures of Mandarin tone learning. Third, analyses were

563    conducted exploring the role of aptitude in each of these tasks (Section 3.6). Specifically, we

564    wanted to see whether aptitude interacted with *variability-condition* in predicting the benefits of

565    training, in line with the predictions of previous research (Perrachione et al., 2011; Sadakata &

566    McQueen, 2014).

567        Except where stated, analyses used logistic mixed effect models (Baayen, Davidson, &

568    Bates, 2008; Jaeger, 2008; Quené & Van den Bergh, 2008) using the package lme4 (Bates,

569    Maechler, & Bolker, 2013) for the R computing environment (R Development Core Team,

570    2010). Logistic mixed effect models allow binary data to be analysed with logistic models rather

571    than as proportions, as recommended by Jaeger (2008). In each of the analyses, the factor

572    *variability-condition* has three levels (low variability [LV], high variability [HV], and high

573    variability blocked [HVB]) which we coded into two contrasts with LV as the baseline (LV

574    versus HV, LV versus HVB). An exception to this is the training data, where a model containing

575    all three conditions would not converge and we took a different approach, as described in Section

576    3.3. We also included the interactions between these contrasts and the other factors. We used

577    centred coding which ensured that other effects were evaluated as averaged over all three levels

578    of *variability-condition* (rather than the reference level of LV[3]). Similarly, for the Three Interval

579    Oddity task, we included a *trial-type* factor. The purpose of this was to control for the fact that

580    participants were likely to find some trial types easier than others due to the gender of the

581    speakers producing the stimuli (see Section 2.3.3.1). We therefore coded a factor *trial-type* with

582    three levels (neutral, easy, hard– see method) and included contrasts with neutral ("neutral versus

583    easy" and "neutral versus hard") using centered coding. In order to perform the analysis

584    comparing pre- and post-test performance, *test-session* was coded as a factor with two levels

585    (pre-test/post-test) with "pre-test" set as the reference level. This allowed us to look at the

586    (accidental) possible differences between the experimental conditions at the pre-test stage, as

587    well as whether post-test performance differed from this baseline. All other predictors, including

588    both discrete factor codings with two levels (*item-novelty* in the Word Repetition and Three

589    Interval Oddity tasks, and *voice-novelty* in the Picture Identification task) and numeric predictors

590    (*training-session*) in the Training data analyses and the individual difference measures in the

---

[3] 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.

PeerJ reviewing PDF | (2018:07:29837:2:1:NEW 18 May 2019)

591    models reported in Section 3.6), were centred (i) to reduce the effects of collinearity between

592    main effects and interactions, and (ii) so that the main effects were evaluated as the average

593    effects over all levels of the other predictors (rather than at a specified reference level for each

594    factor). We automatically put experimentally manipulated variables and all of their interactions

595    into the model, without using model selection (except for "*trial-type*" in the Three Interval

596    Oddity task which works as a control factor and for this factor we only used its main effect and

597    the interaction with *test-session*). However, we did not inspect the models for all main effects

598    and interactions. Instead, we report the statistics which were necessary to look for accidental

599    differences at pre-test, and those related to our hypotheses. We aimed to examine whether the

600    training improved participants' performance on both untrained items and untrained voices and

601    whether such improvement was modulated by their individual aptitudes. Participant is included

602    as a random effect and a full random slope structure was used (i.e., by-subject slopes for all

603    experimentally manipulated within-subject effects (*test-session*, *voice-novelty*, *item-novelty*) and

604    interactions, as recommended by Barr, Levy, Scheepers, and Tily (2013). In some cases the

605    models did not converge and in those cases correlations between random slopes were removed.

606    Models converged with Bound Optimization by Quadratic Approximation (BOBYQA

607    optimization; Powell, 2009). R scripts showing full model details can be found here:

608    https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f

609          In addition to the frequentist analyses, in order to aid interpretation of key null results we

610    also included Bayes factor analyses. Our approach for these is described within the relevant

611    section (Section 3.7).

612    ***3.2    Individual Aptitude Tasks***

613   *3.2.1   The Pitch Contour Perception Test*

614        The predicted variable was whether a correct response was given (1/0) on each trial. The

615   predictors were the contrasts between *variability-conditions* (LV versus HV; LV versus HVB)

616   and *test-session* (pre-test, post-test). There was no significant difference between the LV and HV

617   groups ($\beta$ = -0.35, *SE* = 0.26, *z* = -1.38, *p* = 0.17) or between the LV and HVB groups ($\beta$ = 0.17,

618   *SE* = 0.26, *z* = 0.66, *p* = 0.51) at pre-test on this measure. Participants showed significant

619   improvement after training ($\beta$ = 0.21, *SE* = 0.05, *z* = 4.13, *p* < 0.001), which can be seen in

620   *Figure 3*.

621        Thus, the three participant groups did not differ in their pre-test performance and the

622   groups showed equivalent improvement from pre- to post- test. Given that this measure is

623   affected by training, we used participants scores at pre-test as our measure of individual

624   differences in the analyses reported in Section 3.6..

625   *3.2.2   Categorisation of Synthesised Tonal Continua*

626        We estimated individual's performance on the Categorisation of Synthesised Tonal

627   Continua task following Sadakata and McQueen (2014). We used the Logistic Curve Fit function

628   in SPSS to calculate a slope coefficient for each participant (Joanisse, Manis, Keating &

629   Seidenberg, 2000). The slope (standardized $\beta$) indicates individual differences in tone perception.

630   The smaller the slope, the better the performance. Sadakata and McQueen, removed data from

631   participants with a slope measuring greater than 1.2. Using this threshold 43/60 participants

632   failed the threshold in the current study. This is consistent with the observation that most of the

633   participants were not able to consistently categorise the endpoints of the continua, indicating that

634   this was not a good test of aptitude. We do not report further analyses involving this aptitude

635     variable however they can be found in the supplemental materials

636     (https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f).

637     *3.3*     *Training*

638         A model containing data from all three conditions did not converge; however two

639     separate models, one including the LV and HV conditions, and the other the LV and HVB

640     conditions (with condition as a factor with two levels), did converge. In each case the predicted

641     variable was whether a correct response was given (1/0) on each trial. The predictors were the

642     numeric factor *training-session* (1:6) and the factor *variability-condition* which had two levels

643     (Model 1: LV versus HV; Model 2, LV versus HVB). The mean accuracy is displayed in *Figure*

644     *4.*

645     In both models, there was an effect of *training-session* (Model 1: $\beta = 0.49$, $SE = 0.04$, $z = 11.52$,

646     $p < .001$; Model 2: $\beta = 0.53$, $SE = 0.04$, $z = 12.17$, $p < .001$): Participants' performance increased

647     significantly over time, with additional training sessions. Overall, the LV group performed better

648     than both the HV group ($\beta = -0.79$, $SE = 0.16$, $z = -5.03$, $p < .001$) and the HVB group ($\beta = -$

649     $0.83$, $SE = 0.32$, $z = -2.61$, $p < .01$). However, the LV versus HV contrast was also modulated by

650     an interaction with *test-session* ($\beta = -0.19$, $SE = 0.04$, $z = -4.59$, $p < .001$), as was the LV versus

651     HVB contrast ($\beta = -0.35$, $SE = 0.08$ $z = -4.33$, $p < .001$). From *Figure 4* it can be seen that the

652     LV and the HVB group did not differ in the first session (i.e. where they get identical input) but

653     the difference gradually increased over the next few sessions. For the LV and the HV group, they

654     differed starting from the first session and this difference continued to increase throughout

655     training.

656     *3.4*     *Perceptual tests*

657   *3.4.1   Three Interval Oddity Task*

658       The predicted variable was whether a correct response was given (1/0) on each trial. The

659   predictors were *test-session* (pre-test, post-test), *variability-condition* (LV versus HV, LV versus

660   HVB), *trial-type* (neutral versus easy, neutral versus hard) and *item-novelty* (trained item,

661   untrained item). The mean accuracy is displayed in *Figure 5*.

662       At pre-test, there was no significant difference between the LV and HV groups ($\beta$ = -

663   0.002, $SE$ = 0.14, $z$ = -0.01, $p$ = .99) nor between the LV and HVB groups ($\beta$ = 0.12, $SE$ = 0.14, $z$

664   = 0.86, $p$ = .39), suggesting that the groups started at a similar level. However, performance with

665   the "untrained" was significantly greater than performance on the "trained" items at pre-test ($\beta$ =

666   -0.31, $SE$ = 0.06, $z$ = -4.95, $p$ < 0.01), suggesting incidental differences between item sets. As

667   expected, at pre-test participants performed significantly better on "*easy*" trials (where the target

668   speaker had a different gender) than "neutral" trials (where all three speakers had the same

669   gender, $\beta$ = 0.40, $SE$ = 0.08, $z$ = 5.09, $p$ < 0.01) and "neutral" trials were marginally easier than

670   "*hard*" trials (where one of the foil speakers had the odd gender out, $\beta$ = -0.14, $SE$ = 0.08, $z$ = -

671   1.81, $p$ = 0.07).

672       Overall, participants' performance increased significantly after training (*Mpre* = 0.59,

673   *SDpre* = 0.21, *Mpost* = 0.66, *SDpost* = 0.19, $\beta$ = 0.31, $SE$ = 0.05, $z$ = 6.54, $p$ < .001). The

674   interaction between *test-session* and *item-novelty* was not significant ($\beta$ = 0.14, $SE$ = 0.09, $z$ =

675   1.49, $p$ = .14), suggesting no evidence that training had a greater effect for trained words than for

676   untrained words. Critically, there was no interaction with *test-session* for either the contrast

677   between the LV versus the HV conditions ($\beta$ = -0.01, $SE$ = 0.12, $z$ = -0.12, $p$ = .90) or the

678   contrast between the LV versus the HVB conditions ($\beta$ = 0.01, $SE$ = 0.12, $z$ = 0.11, $p$ = .91) and

679   they were not qualified by any higher level interactions with *item-novelty* (LV versus HV: $\beta$ = -

680 0.1, *SE* = 0.22, *z* = -0.64, *p* = 0.52; LV versus HVB: $\beta$ = 0.13, *SE* = 0.22, *z* = 0.57, *p* = 0.57).

681 This suggests no evidence that the extent to which participants improved on this task between

682 pre and post-test differed according to *variability-conditions,* or that this differed for *trained*

683 versus *untrained* items.

684    Although not part of our key predictions, we also looked to see if there was evidence that

685 participants improved more with the easier or harder trials. In fact, the interaction between *test-*

686 *session* and the contrast between "*easy*" and "neutral" was significant ($\beta$ = -0.27, *SE* = 0.11, *z* = -

687 2.39, *p* = .02) while the contrast between "neutral" and "*hard*" was not ($\beta$ = 0.12, *SE* = 0.11, *z* =

688 1.06, *p* = .29). This was due to the fact that there was improvement for "neutral" (*Mpre* = 0.57,

689 *SDpre* = 0.14, *Mpost* = 0.65, *SDpost* = 0.15) and "hard" trials (*Mpre* = 0.54, *SDpre* = 0.16,

690 *Mpost* = 0.65, *SDpost* = 0.15) but not for "easy" trials (*Mpre* = 0.66, *SDpre* = 0.16, *Mpost* = 0.68,

691 *SDpost* = 0.15).

692 *3.4.2   Picture Identification*

693    The predicted variable was whether a correct response was given (1/0) on each trial. The

694 predictors were the factor *voice-novelty* (Trained voice, Untrained voice) and the factor

695 *variability-condition* which had two contrasts (LV versus HV, LV versus HVB). The mean

696 accuracy is displayed in *Figure 6.*

697    There was a main effect of *voice-novelty* ($\beta$ = 1.07, *SE* = 0.16, *z* = 6.53, *p* < .001)

698 reflecting higher performance in trials with trained voices. Although participants in the LV group

699 performed better than those in the HV group ($\beta$ = -0.71, *SE* = 0.32, *z* = -2.23, *p* =.03), there was

700 no significant difference between the LV and the HVB group ($\beta$ = -0.14, *SE* = 0.32, *z* = -0.44, *p*

701 =.66) and there was a significant interaction between *voice-novelty* and both the LV-HV contrast

702 ($\beta$ = -1.19, *SE* = 0.35, *z* = -3.43, *p* < .01) and the LV-HVB contrast ($\beta$ = -1.11, *SE* = 0.36, *z* = -

703 3.08, $p < .01$). Breaking this down by *variability-condition*: for each condition there was

704 significantly better performance with trained than untrained voices (LV: $\beta = 1.83$, $SE = 0.29$, $z$

705 $= 6.42$, $p < 0.001$; HV: $\beta = 0.64$, $SE = 0.23$, $z = 2.86$, $p < 0.01$; HVB: $\beta = 0.73$, $SE = 0.26$, $z =$

706 2.82, $p < 0.01$), indicating greater ease with the familiar voice. Breaking down by *voice-novelty*:

707 For the trained voice, performance was higher in the LV condition than in either the HV or HVB

708 conditions, although this was only significant for the LV versus HV contrast (LV versus HV: $\beta =$

709 -1.30, $SE = 0.44$, $z = -2.97$, $p < 0.01$; LV versus HVB: $\beta = -0.70$, $SE = 0.45$, $z = -1.55$, $p = 0.12$).

710 Importantly, for untrained voices, neither of the contrasts between conditions was significant

711 (LV versus HV: $\beta = -0.12$, $SE = 0.26$, $z = -0.45$, $p = 0.65$; LV versus HVB $\beta = 0.41$, $SE = 0.27$, $z$

712 $= 1.51$, $p = 0.13$), indicating no evidence for greater generalisation following high variability

713 training.

714 ### 3.5 *Production tests*

715 *3.5.1 Coding and inter-rater reliability analyses*

716 The same methods were used for both production tests. The files were combined into a

717 single set, along with the 360 stimuli which were used in the experiment (and which were

718 produced by native Mandarin speakers). The latter items were included in order to examine

719 whether the raters were reliable. All stimuli were rated by two raters: Rater 1 was the first author

720 and Rater 2 was recruited from the UCL MA Linguistics program and was naïve to the purposes

721 of the experiment. Raters were presented with recordings in blocks in a random sequence (blind

722 to test-type, condition, whether the stimulus was from pre-test or post-test and whether it was

723 produced by a participant or was one of the experimental stimuli). For each item, raters were

724 asked to (i) identify the tone, (ii) give a rating quantifying how native-like they thought the

725 pronunciation was compared (1-7 with 1 as not recognizable and 7 as native speaker level), and

726 (iii) transcribe the pinyin (segmental pronunciation) produced by the participants.

727       If there was no sound or the tone was unrecognizable, the rater coded 0 when identifying

728    the tone. Data from these trials were removed from the dataset before analyses were conducted.

729    In addition, all of the data from one participant was removed from the analyses due to bad

730    recording quality resulting from a technical error. In total, this resulted in 3.38% (359/10620) of

731    production trials being removed from analysis (*Word Repetition:* Pre-test 1.98% (84/4248); Post-

732    test 3.72% (158/4248); *Picture Naming* 5.51% (117/2124)). Three measurements were taken

733    from the production tasks: mean accuracy of tone identification (Tone accuracy), mean tone

734    rating (Tone rating) and mean accuracy of production in pinyin (derived by coding each

735    production as correct (1= the entire string is correct) or incorrect (0 = at least one error in the

736    pinyin)). As a first test of rater reliability, performance with the native speaker stimuli was

737    examined– these were near ceiling: Rater 1: Tone accuracy = 98%, Tone rating = 6.7, Pinyin

738    accuracy = 80%; Rater 2: Tone accuracy = 87%, Tone rating = 6.5, Pinyin accuracy = 80%).

739       Furthermore, for the remaining data (i.e. the experimental data) inter-rater reliability was

740    examined for all three measures for the two production tasks. For the binary measures (Tone

741    accuracy and Pinyin accuracy), kappa statistics were calculated using the "fmsb" package in R

742    (Cohen, 2014). For the Word Repetition data, for Tone accuracy *kappa* = 0.39 ("fair

743    agreement"), and for Pinyin accuracy *kappa* = 0.33 ("fair agreement"; Landis & Koch, 1977).

744    For the Picture Naming test, for Tone accuracy *kappa* = 0.67 ("substantial agreement") and for

745    Pinyin accuracy *kappa* = 0.53 ("moderate agreement"); For the Tone rating, the package "irr" in

746    R was used to assess the intra-class correlation (McGraw & Wong, 1996) based on an average-

747    measures, two-way mixed-effects model. For Word Repetition, *ICC* = 0.22 and for Picture

748    Identification *ICC* = 0.37; according to Cicchetti (1994), values less than .40 are regarded as

749    "poor". Given this, we do not include analyses with Tone Rating as the dependent variable

750    (though these data are included in the data set

751    [https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f](https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f)). All of the analyses

752    presented in Sections 3.5.2 and 3.5.3 were based on Rater 2 (the naive rater).

753    *3.5.2   Word Repetition*

754                          *3.5.2.1   Tone accuracy*

755          The predicted variable was whether a correct response was given (1/0) on each trial (as

756    identified by the coder). The predictors were *test-session* (pre-test, post-test), *variability-*

757    *condition* (LV versus HV, LV versus HVB) and *item-novelty* (trained, untrained). The mean

758    accuracy, split by test-session and training condition, is shown in *Figure 7*

759          At pre-test, there was no significant difference between the LV and the HV group ($\beta$ =

760    0.01, *SE* = 0.18, *z* = 0.06, *p* = .95) nor between the LV and the HVB group ($\beta$ = 0.11, *SE* = 0.18,

761    *z* = 0.64, *p* = .53), suggesting the groups started at a similar level. There was also no difference

762    between trained and untrained words at pre-test ($\beta$ = -0.02, *SE* = 0.07, *z* = 0.-0.26, *p* = 0.80).

763          Across the three groups, participants' performance increased significantly after training

764    (*Mpre* = 0.71, *SDpre* = 0.09, *Mpost* = 0.79, *SDpost* = 0.09, $\beta$ = 0.40, *SE* = 0.08, *z* = 5.29, *p* <

765    .001). There was no significant difference in the improvement for trained and untrained items

766    (*word-type* by *test-session* interaction: $\beta$ = 0.13, *SE* = 0.10, *z* = 1.22 *p* = .22). Critically, the

767    interactions between the variability contrasts and *test-session* were not significant (LV versus

768    HV: $\beta$ = -0.10, *SE* = 0.18, *z* = -0.55, *p* = .58; LV versus HVB: $\beta$ = -0.11, *SE* = 0.18, *z* = -0.62, *p* =

769    .54), and they were not qualified by any higher level interactions with *item-novelty* (LV versus

770    HV: $\beta$ = 0.15, *SE* = 0.25, *z* = 0.61, *p* = .54; LV versus HVB: $\beta$ = -0.31, *SE* = 0.26, *z* = -1.21, *p* =

771    .23). This suggests there is no evidence that participants' improvement in their production of

772    tones was affected by their *variability-condition,* or that this differed for *trained* versus *untrained*

773    items.

774                 *3.5.2.2   Pinyin accuracy*

775         The predicted variable was whether the participants produced the correct string of

776   phonemes (1/0) in each trial (as determined by Rater 2). The predictors were *test-session* (pre-

777   test, post-test)*, variability-condition* (LV versus HV, LV versus HVB) and *item-novelty* (trained,

778   untrained). Mean pinyin accuracy is displayed in *Figure 8*.

779         At pre-test, there was no significant difference between the LV and the HV group ($\beta$ = -

780   0.01, *SE* = 0.11, *z* = -0.11, *p* = .91) nor between the LV and the HVB group ($\beta$ = -0.03, *SE* =

781   0.11, *z* = -0.24, *p* = .81), suggesting that the groups started at a similar level. However,

782   participants did better on untrained words than trained words at pre-test ($\beta$ = 0.21, *SE* = 0.07, *z* =

783   3.11, *p* < .01), suggesting potential accidental differences in these items. Participants showed

784   significant improvement after training (*Mpre* = 0.54, *SDpre* = 0.09, *Mpost* = 0.58, *SDpost* = 0.19,

785   $\beta$ = 0.15, *SE* = 0.05, *z* = 3.38, *p* < .01). However, there was no evidence that different variability

786   conditions resulted in different amounts of improvement  (*test-session* by LV versus HV: $\beta$ =

787   0.05, *SE* = 0.11, *z* = 0.46, *p* = .65; *test-session* by LV versus HVB: $\beta$ = -0.12, *SE* = 0.11, *z* = -

788   1.08, *p* = .28) or any interaction between *variability condition*, *test-session* and *item-novelty* (LV

789   versus HV: $\beta$ = 0.11, *SE* = 0.22, *z* = 0.51, *p* = .61; LV versus HVB: $\beta$ = -0.14, *SE* = 0.22, *z* = -

790   0.64, *p* = .52). This suggests there is no evidence that participants' improvement in pinyin

791   accuracy was affected by their *variability-condition,* or that this differed for *trained* versus

792   *untrained* items.

793   *3.5.3   Picture Naming*

794                *3.5.3.1   Tone accuracy*

795         The predicted variable was whether a correct response was given (1/0) on each trial (as

796   identified by the coder). There was only one predictor, *variability-condition* (LV versus HV, LV

797   versus HVB) for both models. The descriptive statistics are displayed in *Figure 9*.

798     Participants in the LV group showed no significant difference compared with the HV

799     group ($\beta$ = -0.34 $SE$ = 0.19, $z$ = -1.81, $p$ = 0.07) and the HVB group ($\beta$ = -0.10, $SE$ = 0.19, $z$ = -

800     0.52, $p$ = .61. This suggests there is no evidence that participants' ability to produce the tones

801     accurately differed according to their *variability-condition*.

802                             *3.5.3.2   Pinyin Accuracy*

803     The predicted variable was whether the participants produced the correct string of phonemes

804     (1/0) in each trial and there was a single predictor *variability-condition* (LV versus HV, LV

805     versus HVB). For both models there was no significant difference between variability conditions

806     (LV versus HV: $\beta$ = 0.09, $SE$ = 0.23, $z$ = 0.41, $p$ = 0.68; LV versus HVB: $\beta$ = 0.12, $SE$ = 0.23, $z$ =

807     0.51, $p$ = 0.61). This suggests there is no evidence that participants' pinyin accuracy differed

808     according to their *variability-condition*.

809     **3.6     Analyses with Individual Aptitude**

810     In order to look at the effect of learner aptitude and the interaction between this factor

811     and variability condition, we first calculated the mean accuracy at pre-test on the Pitch Contour

812     Perception Test for each participant. This score (scaled by a factor of 10, so that each one unit

813     increase in aptitude corresponded to a 10% higher performance in the Pitch Contour Perception

814     test) was centered and used as a continuous predictor (*aptitude*) and added to each of the models

815     reported above. In addition, we added the interaction between this factor and key experimental

816     factors (see Table 4). Based on Perrachione et al. (2011) and Sadakata and McQueen (2014), for

817     our measures of tone-learning, high variability should benefit high aptitude participants only,

818     while low variability would benefit low aptitude participants only. In our design, we used a

819     continuous measure of individual ability rather than a binary division of high and low variability.

820     We therefore predicted a stronger positive correlation between *aptitude* and amount of learning

821     in the high variability condition than in the low variability condition. In the tests administered

822     only post training (i.e. Picture Identification and Picture Naming) this would show up as an

823     interaction between aptitude and condition. In the models for the pre- and post-test data (i.e.

824     Three Interval Oddity and Word Repetition) this would show up as a three-way interaction

825     between *condition, test-session* and *aptitude*. We also looked at the interactions between these

826     factors and *voice-novelty* (Picture Identification) and *item-novelty* (Three Interval Oddity and

827     Word Repetition). Note that there are no clear directional hypotheses here: Perrachione et al.

828     (2011) found the interaction in a test with untrained voices and trained items, and Sadakata and

829     McQueen (2014) found the interaction in a test with trained voices and trained items.  For

830     training, in principal both the two-way interaction of *aptitude* by *condition* and the three-way

831     interaction of *aptitude* by *condition* by *training-session* are of interest. However, it was not

832     possible to fit a converging model containing the three-way factor[4].

833             Each model reported in Table 4 contained all the fixed effects included in the original

834     models in addition to the fixed effects listed in the table (note that to avoid convergence issues

835     due to over complex models, we did *not* attempt to include the complete set of interactions for

836     every combination of experimental variables with aptitude – only those for which we had

837     predictions). We attempted to have full random effects structure for these fixed effects however

838     in some cases we had to remove correlations between slopes due to problems with convergence

839     and for one of the models with the training data we had to remove the random slope for training

840     session). Note that we don't include models for the pinyin measures, since our measure of

841     aptitude is relevant to tone learning only. For each of the new models we first confirmed that

842     adding in the new effects and interactions with the individual measures did not change any of the

---

[4] This was the case even if we split the data into two models, as we did in Section 3.3.

843    previously reported patterns of significance for the experimental effects (see script

844    https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f) for full models.

845        The results are shown in Table 4. *Aptitude* is a positive predictor of performance in each

846    of the tests and in training, with p-values significant or marginal in each case. However there was

847    no interaction between *aptitude* and any other factor. Thus, there was no evidence that this

848    measure of aptitude correlated with participants ability to benefit from training (no interaction

849    with *test-session*), nor - critically for our hypothesis - did this differ by training condition (no

850    interaction with *condition* or with *test-session* by *condition*).

851        Although the analyses use a continuous measure of Pitch Contour Perception Test, for the

852    purposes of visualisation, *Figure 10* (Three Interval Oddity task and Training task)*, Figure 11*

853    (Picture Naming and Picture Identification) and *Figure 12* (Word Repetition) use the mean

854    accuracy for participants split into aptitude groups using a median split based on their Pitch

855    Contour Perception Test score.

856        In sum, participants with higher aptitude measures were better at the tasks, but there is no

857    evidence either that this affected their improvement due to training, or, critically, their ability to

858    benefit from the different variability exposure sets.

859    **3.7    Bayes Factor Analyses**

860        In the analyses reported above, we did not find evidence – in any of our tests – for either

861    of two key hypotheses: (1) the hypothesis that training with multiple speakers leads to greater

862    generalization to new speakers than training with a single speaker *or* (2) the hypothesis that there

863    is an interaction between the variability of the training materials and participant aptitude, such

864    that higher aptitude participants benefit more from training with multiple speakers while lower

865    aptitude participants benefit more from training with a single speaker. However, there is a

866    difficulty in interpreting these null results since a non-significant result (p > .05) does *not* tell us

867    whether we have evidence for the null, as opposed to no evidence for any conclusion at all, or

868    even evidence against the null. Thus, we should not reduce our confidence in either of our

869    hypotheses on the basis of the null results reported above (despite the fact that reducing

870    confidence in a theory following non-significant results is common practice) – see Dienes (2014)

871    for discussion. An alternative statistic is a Bayes Factor, which are used to assess the strength of

872    evidence for one theory (H1) over another (the null hypothesis). We therefore supplemented the

873    analyses above by computing Bayes factors for contrasts relating to these two key hypotheses.

874    These are reported in sections 3.7.1 and 3.7.2 below.

875    *3.7.1    H1: Greater generalization - to novel voices and in production - in the multiple speaker*

876    *conditions (HV and HVB) than in the low variability condition (LV)*

877    We aimed to compute Bayes Factors comparing this hypothesis to the null for each of our

878    data sets. To have maximum evidence, we pool the HV and HVB conditions and contrast this

879    with the LV condition. For the post-tests we are interested in the evidence for a main effect of

880    this contrast. For the pre-post tests, we are interested in the interaction between this contrast and

881    session. To further maximize evidence, for the Three Interval Oddity test and Word Repetition

882    tests we look at trained and untrained items combined (since *both* types of item involve

883    generalisation to an untrained voice and thus should benefit from high variability training),

884    however in the Picture Identification test we excluded trained *voice* test items, since the benefit

885    of high variability training was not predicted for these items. For the production measures, we

886    are interested in whether there is a high variability benefit for our tone learning measure and our

887    pinyin measure (the latter given that Barcroft and Sommers, 2014, found a benefit of multi-

888    speaker training in their vocabulary recall task).

889     We computed Bayes factors following Dienes (2014) and Dienes, Coulton and Heather

890     (2018). To compute a Bayes factor (*B*) it is necessary to have both a model of the data and a

891     model of H1. The model of the data is an estimate of the mean difference for the contrast in

892     question, and of the standard error. Here, we get these estimates by running logistic mixed

893     models and taking the betas and standard errors for the relevant coefficients (note that this allows

894     us to meet normality assumptions by continuing to work within log-odds space). The models we

895     ran here were similar to the previous analyses but with variability-condition coded as a centered

896     contrast between LV and the HV+HVB conditions, and other factors combined/excluded as

897     described in the previous paragraphs. The full set of models is in

898     https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f.

899     We model H1 using a half-normal distribution with a mode of 0 and a standard deviation

900     *x* which is set to be a rough estimate of the predicted difference for this contrast. This allows for

901     possible effects between 0 and twice the predicted effect, with values closer to 0 being more

902     likely (Dienes, 2014).

903     In the absence of any prior data using sufficiently similar materials, and since we did not

904     wish to use unprincipled default values, we estimated *x* for each contrast using the scale and/or

905     values from elsewhere in the data (see Dienes 2014, 2015 for a related approach). Specifically,

906     for each of the cases where we predicted a main effect (Picture Identification and Picture

907     Naming), we set *x* as the difference between the grand mean (the Intercept - since we use a

908     centered coding) and an estimate of minimal possible performance on the task. The logic is as

909     follows[5]: The *maximum* difference between conditions is seen if *low variability* participants

910     show baseline performance and *high variability* participants show performance greater than

---

[5] Further details of the logic of these computations is spelt out in the script available at
https://osf.io/wdh8a/?view_only=ad8455b30b2e4271aaa4cc55fc94a40f

911  baseline. In this case, if performance on this test is $p$ (so the grand mean is $\bar{p}$) and the baseline is

912  $b$, the difference in $p$ between the two conditions will be equal to: $2(\bar{p}\text{-}b)$. This gives us an

913  estimate of the *maximum* value of $x$; since we are using a half normal distribution with a mean of

914  zero, we assume the maximum value is equal to approximately *2SD*, so we can set our estimate $x$

915  of the standard deviation to be equal to *half* of this value (i.e. $x = \bar{p}\text{-}b$.). Baseline performance

916  depends on the task: for the 2AFC Picture Identification task it is chance (50% = 0 in log odds

917  space); for the Picture Naming, tone measure, we assume a ¼ chance of identifying the correct

918  one (25% = -1.099 in log odds space); for Picture Naming, Pinyin measure, there is no chance

919  and we therefore took minimal performance as making one correct response in the test[6] (i.e. 1/72

920  = -4.263 in log odds space). For the cases where we are estimating an interaction between *test-*

921  *session* and *variability-condition* we set $x$ as equal to the mean increase in performance from *pre-*

922  and *post-* test across conditions (main effect of *test-session*). The logic is as follows: the

923  *maximum* difference is seen if *low variability* participants show no effect of *test-session* (no

924  improvement) and *high variability* participants show a positive effect of *test-session*. In this case,

925  if the mean effect of *test-session* is $\bar{t}$, the difference in $t$ between the two conditions will be equal

926  to $2\bar{t}$. Again, we can set our estimate of $x$ to be half this value (i.e. $x = \bar{t}$).

927          We interpret BFs using the following conventions: $B < 1/3$ indicates substantial evidence

928  for the null, $B > 3$ indicates substantial evidence for H1, values between 1/3 and 3 indicate that

929  the data collected do not sensitively distinguish H0 from H1 (Jeffreys 1961; Dienes 2008). Since

930  there is subjectivity in how the values for H1 are determined, we indicate the robustness of

931  Bayesian conclusions by reporting a robustness region for each $B$, which gives the range of

932  values of the scale factor $x$ that qualitatively support the same conclusion (i.e. evidence as

---

[6] Note that we cannot compute log-odds of 0.

933 supporting H0, or as supporting H1, or there not being much evidence at all). Note that for

934 evidence for H0, the maximum $x$ is always infinity. The results are reported in Table 5. It can be

935 seen we have substantial or strong evidence for the null for every test except for the Word

936 Repetition test for the Pinyin accuracy measure, where the evidence is ambiguous, and that the

937 robustness regions indicate that we would continue to have evidence for the null even with

938 smaller estimates of the scale factor $x$.

939 *3.7.2   H1: There is an interaction between an individual's tone-aptitude and variability-*
940 *condition, such that participants with greater tone-aptitude show greater performance*
941 *following the multiple speaker conditions (HV and HVB) and those with lesser tone*
942 *aptitude show greater performance in the single speaker condition (LV)*
943       We aimed to compute Bayes Factors comparing this hypothesis to the null for each of our

944 data sets. We take the same approach as above except that we also compute Bayes factors for

945 Training data, and for the Picture Identification test we look at both trained voice and untrained

946 voice data – pooling the two in order to maximize available evidence. This is because this

947 interaction has been reported with trained items (Sadakata & McQueen, 2013) as well as

948 untrained items (Perrachione et al., 2011). We again combine the HV and HVB conditions

949 except for training where we look at the LV versus HV and LV versus HVB contrasts separately,

950 since we have seen in our previous analyses that HV and HVB are quite different (HVB

951 participants show higher performance).We again combine the evidence from trained and

952 untrained *items* in the pre-post tests. For the post-session only tests, we are interested in the

953 evidence for an interaction between the *variability-condition* contrast and *aptitude*. For the tests

954 which appeared both pre- and post- training, we are interested in the interaction between the

955 *variability-condition* contrast, *aptitude* and *test-session*. For training we look at the evidence for

956 an interaction between each *variability-condition* contrast and *aptitude* (a more complex model

957   containing the interaction with training-session did not converge). As in our frequentist analyses

958   of aptitude, for the production measures – Word Repetition and Picture Naming – we do *not* look

959   at the pinyin measures since our aptitude measure is relevant only to tone learning.

960         We computed Bayes factors following the same procedure as in Section 3.7.1 and again

961   derived our estimates of the scale factor $x$ - the difference predicted under H1 - using the scale

962   and/or values from elsewhere in the data. Specifically, for each of the cases where we predicted a

963   two-way interaction between *variability-condition* and *aptitude* we set $x$ as equal to the mean

964   effect of *aptitude* across conditions (main effect of *aptitude*)[7]. The logic is as follows: The

965   *maximum* difference is seen if *low variability* participants show no effect of *aptitude* and the *high*

966   *variability* participants show a positive effect of *aptitude* (note that a negative effect of aptitude

967   is not expected in any condition). In this case, if the mean effect of *aptitude* is $\bar{a}$, the difference in

968   $a$ between the two conditions will be equal to $2\bar{a}$. Again, we can set our estimate of $x$ – the *SD* of

969   the half normal – to be half this maximum value i.e. $x = \bar{a}$. For the cases where we are interested

970   in the three-way interaction between aptitude, test-condition and test-session, we based our

971   estimate on half the difference between the maximal effect of aptitude *(maxA* – taken from the

972   scale*)* and their actual aptitude score at pre-test (*baselineA* – taken from the data). The logic is as

973   follows: The maximal effect of the interaction would be seen if participants in the *low variability*

974   condition showed the same baseline effect of aptitude at *pre-test* and at *post-test* (*ba*), whereas

975   participants in the high variability condition showed maximal improvement at post-test (*maxa*).

976   In this case, the interaction between aptitude and session for the high variability group would be

---

[7] An alternative which would be more equivalent to the other BF analyses would be to inform the effect using the value of the two-way interaction of aptitude: test-session. We do not do this since we did not find an effect of this two-way interaction in either data set.

977     equal to: $maxa - ba$. Again, we can set our estimate of $x$ – the $SD$ of the half normal - to be half

978     this maximum value, i.e. $x = \frac{ma - ba}{2}$

979          The maximum effect of aptitude was computed from the scale and the length of the

980     aptitude predictor. Specifically, we assumed that the maximal effect of aptitude would be

981     obtained if participants with maximal aptitude were at ceiling (71/72 correct – log odds 4.263)

982     and those with minimal aptitude were at chance (25% in Word Repetition, Tone Accuracy, log

983     odds= 1.099; 33.33% in Three Interval Oddity, log odds = 0.693). We divided this range by the

984     length of the aptitude predictor to obtain a measure of a one-step change in aptitude.

985     The results are summarised in Table 6. It can be seen that although there is more evidence for the

986     null than H1 in each case (i.e. BF < 1) we do *not* have substantial evidence for the null over H1

987     in any case. Thus, we cannot draw any inferences about the interaction from this data. Note that,

988     in most cases, the robustness regions indicate that even if the scale factor $x$ was twice as large,

989     i.e. corresponding to the *maximum* value we might expect, the $B$ would be ambiguous.

990     **4**       **Discussion**

991          The current study investigated the effect of different types of phonetic training on English

992     speakers' learning of novel Mandarin words and tones. To our knowledge, this is the first study

993     to train naive participants on all four Mandarin tones, using real language stimuli embedded in a

994     word learning task. Learning was examined using a range of perception and production tasks.

995     Following previous literature, we compared three training conditions: low variability (single

996     speaker), high variability (four speakers, presented intermixed) and high variability blocked (four

997     speakers, presented in blocks). We also administered tests designed to tap individual aptitude in

998     the perception of pitch contrasts, adapted from the previous literature. The results indicated that

999     participants' performance increased during training and that training also led to improved

1000   performance on pre- to post- tests of discrimination and production, with evidence of

1001   generalisation to untrained voices and items. Participants also showed some ability to recall

1002   trained words – including their tones – in a picture naming task administered at post-test.

1003   However, the only place where we saw any effect of the variability manipulation was in the

1004   training task (and with trained items in the picture identification task, which was highly similar

1005   to training), where the *low* variability group outperformed both of the high variability groups.

1006   Critically, we found no evidence in any of our tests that high variability input benefitted learning

1007   or generalisation, nor did we find any evidence of an interaction between individual aptitude and

1008   the ability to benefit from high variability training. In the following discussion, we first consider

1009   the findings from each task in turn before turning to a more general discussion of our findings in

1010   relation to the predicted benefit of high variability input.

1011   ***4.1    Tests of individual aptitude***

1012        In the current work, we conducted two tests with the purpose of capturing individual

1013   aptitude: The Pitch Contrast Perception Test (following Perrachione et al 2011) and the

1014   Categorisation of Synthesised Tonal Continua, following Sadakata and McQueen (2014).

1015   Although our goal was to measure participants' baseline aptitude, the tests were conducted both

1016   at pre- and post- test, following Sadakata and McQueen, who did not find differences from pre-

1017   to post- tests with their categorisation measure, and who used combined data from pre- and post-

1018   test to compute participants slopes. Unfortunately, the performance of our own participants

1019   suggested that the Categorisation of Synthesised Tonal Continua test was not a good test of

1020   aptitude, with the majority of participants failing to meet the slope threshold used in Sadakata

1021   and McQueen, and most being unable to consistently categorise the end points of the continua. It

1022   is unclear why our results differ from the previous study (we aimed to follow their procedures),

1023    but this meant that we were unable to use this as an aptitude measure in our later analyses. The

1024    scores on the Pitch Contrast Perception Test alone therefore served as our measure of individual

1025    aptitude. Interestingly, preliminary analyses (Section 3.2.1) demonstrated that performance in

1026    this test improved from pre- to post- training. This suggests that this measure is not a "pure"

1027    measure of individual differences since it also appears to be affected by experience. Given this,

1028    we only used participants' scores on this test from *pre-test* as the measure of aptitude in

1029    subsequent analyses.

1030    *4.2    Performance in Training*

1031        The training task employed in this study was a 2AFC task, where participants had to

1032    identify the correct meaning of a Mandarin word based on its tone. The results from training

1033    indicate that participants performed better in the single speaker LV training than in either the

1034    multiple speaker HV or HVB groups. This difference was present from the first session for the

1035    LV-HV contrast, and from the second session for the LV-HVB contrast (i.e. the first session

1036    where the two conditions differ), and increased over time for both contrasts. Greater difficulty

1037    with multiple speaker input is in line with the findings of Perrachione et al. (2011), although the

1038    differences did not emerge so rapidly in that study, possibly due to there being fewer trials per

1039    session. Intuitively, repeated exposure to the single speaker in the LV condition allows for

1040    greater adaptation to speaker specific cues, whereas in the HV conditions participants have to

1041    adapt to multiple speakers. This is particularly difficult in the unblocked HV condition, where

1042    trial-by-trial adaptation is needed, which is effortful for participants (Magnuson & Nusbaum,

1043    2007). Importantly, however, for all three groups, their performance gradually increased over

1044    each session. In combination with the fact that their performance on the other tasks increased

1045    after training, this indicates that the training task and materials were effective. We also explored

1046    the role of learner aptitude in this task (as measured by performance on the Pitch Contour

1047    Perception Test at pre-test) and whether this influenced participant's performance differently in

1048    the different variability conditions. Overall, aptitude was found to be a significant predictor of

1049    performance during training. However, there was no evidence for an interaction with training

1050    condition, although our Bayes Factor Analyses suggests that the data here are inconclusive. We

1051    return to this finding in Section 4.5 below.

1052    **4.3    *Perception Tests***

1053        We included two perceptual tasks which tapped learning and generalisation due to

1054    training: A *Picture Identification* administered at post-test and a *Three Interval Oddity* task

1055    administered at both pre- and post-test. The *Picture Identification* task was a version of the

1056    training task without feedback, and is the most similar to the tests used by Perrachione et al.

1057    (2011), and Sadakata and McQueen (2014). We used this test to look at learning of the trained

1058    stimuli, comparing trained and untrained voices. The three interval oddity task had not been used

1059    in the previous studies, but allowed us to use a pre- /post- test design, and also to look at

1060    participants' performance with *untrained items*. These tests provided evidence that participants

1061    improved in their perception of tones following training: They were above chance in using the

1062    tone to identifying the correct picture in the picture identification task at post-test, and they

1063    improved in their ability to discriminate tones in the three interval oddity task (59% performance

1064    prior to training, 66% post training). There was also evidence of generalisation across both

1065    voices and items: Participants were above chance in identifying the correct pictures even with an

1066    untrained voice (although they did show significantly weaker performance than with the trained

1067    voiced) and they improved in their ability to discriminate the between minimal pair items in the

1068    three interval oddity task, even for untrained items.

1069    Our key questions concerned the role of variability in training. First, we were interested in

1070    whether there was evidence that exposure to multiple voices during training led to greater ability

1071    to generalise across voices at test – i.e. greater performance with novel in the high variability

1072    conditions than in the low variability condition. We did not see this. In fact, the only effect of

1073    variability in this data was a *low* variability benefit, which we saw in the Picture Identification

1074    task for the trained-voice items (seen in the contrast between LV and HV condition). This

1075    mirrors what we saw in training and reflects the greater exposure to this particular speaker in the

1076    low variability training. However, in the tests tapping generalisation to a novel speaker – i.e. in

1077    untrained voice trials in the Picture Identification task, and with all of the test-items in the Three

1078    Interval Oddity task, there was no difference between variability-training conditions. Bayes

1079    factor analyses indicate that in both cases, there was substantial evidence for the null.

1080         The second hypothesis was that there would be an interaction between learner aptitude

1081    (as measured by the Pitch Contour Perception Test at pre-test) and variability training condition,

1082    such that high aptitude participants would benefit more from high variability training. Note that

1083    previous work had found this interaction both in tests involving generalisation (Perrachione et

1084    al., 2011) and with trained items (Sadakata & McQueen, 2014) so we considered both in our

1085    analyses here. There was no evidence of such an interaction in either the Picture Identification or

1086    Three Interval Oddity tasks. However, Bayes Factor analyses suggest that the data are

1087    inconclusive. We return to these points in Section 4.5 below.

1088         Another finding from the Three Interval oddity test that is worth noting, although it did

1089    not concern our hypotheses, is that some trial types were harder than others. Recall that this test

1090    involved participants hearing three different stimuli each produced by a different speaker, which

1091    makes noting the similarity across two of the stimuli much harder - something we discovered in

1092    pilot work, where even before training participants were near ceiling with an equivalent task in

1093    which the same speaker produced all three stimuli within a single trial. However, analyses of

1094    *trial-type* demonstrated that participants were additionally affected by the gender of the three

1095    speakers producing each of the stimuli. Specifically, at pre-test, participants showed best

1096    performance for trials where one of the speakers was male and the other two were female, and

1097    the target "odd man" was the male speaker ("easy" trials). In contrast, they showed worst

1098    performance if there was one male and two female speakers, but the "odd man" was one of the

1099    female speakers ("hard" trials). Middle level performance was shown for trials where all three

1100    speakers were female ("neutral" trials). This is presumably due to participants relying on

1101    perceptual cues associated with speaker gender to do the task. Interestingly, our analyses showed

1102    that performance only increased for the trials where the odd one was not the lone male (the

1103    "neutral" and "hard" ones), but not for those where the male was the odd man. Given that

1104    participants are not near ceiling at pre-test (67%), it is perhaps surprising that their trained

1105    knowledge of the tone contrasts does not boost their performance. One possibility is although

1106    they are now better able to use tone cues, they are also *less* likely to use gender based cues,

1107    which they may now realize are less reliable, masking improvement based on tone for these

1108    particular test items.

1109    **4.4    Production Tasks**

1110        In this study, we used two production tasks, a word repetition task administered pre and

1111    post training, in which participants repeated back Mandarin words, and a Picture Naming task

1112    testing vocabulary recall, which was administered at post-test only. High variability perceptual

1113    training for tones has been previously found to transfer to production (Bradlow and Pisoni, 1999;

1114    Zeromskaite, 2014), however the benefits of high variability and low variability training have not

1115    been contrasted.

1116    In the Word Repetition task, there was a significant, though relatively modest

1117    improvement in participants' ability to reproduce the tone of the stimuli, such that it could be

1118    identified by a native speaker (from pre- to post- test: 70% to 76%) and in the Picture Naming

1119    task, participants showed an ability to recall and reproduce the correct tone, although

1120    unsurprisingly with less accuracy than in the repetition task (50%). For Word Repetition, we

1121    were also able to look at transfer to untrained words: As in the perception tasks, there was once

1122    again equivalent improvement for both trained and untrained items. Together, these results

1123    provide evidence that purely perceptual training on tone contrast can transfer to production, as

1124    well as to novel items.

1125    In addition to looking at the production of *tones*, we also looked at participants' ability to

1126    produce the correct segmental phonology (pinyin-score). Participants showed a small but

1127    significant improvement on this measure in Word Repetition (54% correct at pre-test, 58% at

1128    post-test), and some ability to recall the segments in the Picture Naming test (50% correct). This

1129    indicates some learning of segmental phonology due to training, despite the fact that the focus of

1130    the training task was on training tonal information through the presentation of tonal minimal-

1131    pairs.

1132    Turning to the role of variability, the predicted benefit of high variability training was *not*

1133    evident in any of the measures in either of the production tasks, with Bayes factor analyses

1134    indicating substantial evidence for the null except for the Word Repetition pinyin-measure,

1135    where the evidence was ambiguous. With regard to aptitude, although performance on the Pitch

1136    Contour Perception Test at pre-test was predictive of participants' ability to produce tones in

1137    both tasks (indicating a relationship between participants perceptual and production ability), we

1138    did *not* find the predicted interaction between aptitude and variability condition in either task.

1139    Here however, Bayes Factor analyses suggests that the results are inconclusive. We return to

1140    these points about variability below.

1141    *4.5      The Role of High Variability Materials in Training and Generalisation*

1142        In the current study, across all of the different tests, we did not find either an overall

1143    benefit of exposure to high variability training materials for generalisation, or any interaction

1144    between such a benefit and individual aptitude.

1145        We consider first the lack of *overall* variability benefit for generalisation. Importantly, in

1146    addition to finding a pattern of null results (i.e. p < .05) in the frequentist analyses, additional

1147    Bayes Factor analyses also found substantial evidence for the null (BF < .33) in all but one of the

1148    test measures (Word Repetition, Pinyin, where BF = .421). Thus, there is good evidence that, at

1149    least for these training and test materials, exposure to stimuli from multiple speakers does *not*

1150    lead to greater generalisation in either perception or production. This finding is consistent with

1151    the lack of a main effect of variability condition in the transfer tasks in either Sadakata &

1152    McQueen (2014) or Perrachione et al. (2011) (see also footnote 1). However it is at odds with

1153    other phonetic training studies focused on segmental contrasts (Clopper & Pisoni, 2004; Logan et

1154    al. 1991, Lively et al., 1993; Sadakata & McQueen 2013) and with the literature demonstrating a

1155    high variability benefit in vocabulary learning (Barcroft & Sommers, 2005, 2014; Sommers &

1156    Barcroft, 2007, 2011). This suggests that this overall variability benefit may be restricted to

1157    segmental rather than tonal phonetic learning, at least for speakers of a non-tonal L1.

1158        It is difficult to reconcile the lack of benefit for vocabulary learning in the picture naming

1159    task, given the findings of Barcroft, Sommers and colleagues (2005, 2007, 2011, 2014), since

1160    this test is quite similar to that used in their experiments. However, one possibility is that this is

1161    due to differences in our training set up (i.e. focused on training tonal contrasts) compared with

1162    the earlier vocabulary studies. Nonetheless it remains unclear why *tone learning* should be

1163    different from other types of phonetic learning in terms of benefiting from talker-variability.

1164    Theoretically speaking, in a framework where all cues compete, variation in idiosyncratic

1165    speaker-specific cues would be expected to provide key evidence as to which cues are irrelevant

1166    to the phonetic contrast in question (Apfelbaum & McMurray, 2011; Ramscar & Baayen, 2013;

1167    Ramscar, Yarlett, Dye, Denny & Thorpe, 2010). This raises the question of how participants in

1168    our low variability condition are able to generalize at all – i.e. how can they identify the

1169    phonetically relevant cues compared with the idiosyncratic cues associated with the single

1170    speaker to which they were exposed? One possibility is that other variation in our materials aided

1171    generalisation, in particular in our real word stimuli, each tone-contrast is encountered with

1172    multiple consonants and vowels. If item variability also aids generalisation to new speakers, this

1173    might explain why we found equivalent generalisation across conditions instead of seeing greater

1174    generalisation in the HV conditions (i.e. even the LV condition is really a high variability

1175    condition, because of the item variability). On the other hand, Sadakata and McQueen (2014)

1176    also saw generalisation even for their low variability condition, and in their study this condition

1177    lacked variation in terms of both speakers and phonetic contexts. This suggests that the relevant

1178    cues for the tone contrasts may be sufficiently acoustically salient for learners to identify them,

1179    even when exposure occurs in limited contexts.

1180         Another possibility – and the one suggested by the findings of Sadakata and McQueen

1181    (2014) and Perrachione et al. (2011) – is that benefits of high variability for generalisation are

1182    masked by individual differences. In their studies, only high aptitude participants showed a high

1183    variability benefit, while low aptitude participants did not. It is possible that for lower aptitude

1184    participants, the benefits of exposure to varying, idiosyncratic cues are offset by the greater

1185    difficulty that these participants have in attuning to the different speakers during training, as

1186    discussed above (Section 4.1). This explanation is supported by evidence from a study by

1187    Goldinger, Pisoni and Logan (1991) who explored the effect of increasing the processing cost of

1188    multi-speaker input in the context of word recall (in the L1). Specifically, they exposed

1189    participants to single versus multi-speaker word lists, manipulating presentations rates. They

1190    found that single-speaker lists produced better word recall than multiple-speaker lists at short

1191    inter-word intervals (less than 2000 ms) whereas this effect was reversed for longer inter-word

1192    intervals. This suggests that increasing encoding difficulty can remove the benefits of multi-

1193    speaker exposure. Relatedly, Sinkeviciute, Brown, Brekelmans, & Wonnacott (in press; preprint)

1194    found that young learners have greater difficulty processing multi-speaker training materials in

1195    L2 vocabulary learning, and subsequently fail to show a speaker-variability benefit at test. One

1196    interpretation of these findings is that age-related capacity limitations may constrain the ability to

1197    benefit from speaker variability, supporting the notion that differences in capacity limitations can

1198    affect an individual's ability to benefit from multi-talker training.

1199            Returning to the current study, we did not find an interaction between variability-training

1200    and learner aptitude. However, it is important to acknowledge the results of our Bayes factor

1201    analyses, which did not find substantial evidence in support of the null over H1 (or H1 over H0)

1202    for any of the test tasks. This means that we cannot draw conclusions about this hypothesis from

1203    the current data. In theory, we could continue collecting data until we had substantial evidence

1204    for either H0 or H1. To explore the feasibility of this, we conducted supplementary analyses to

1205    estimate the sample size that might be needed to see substantial evidence for the null (based on

1206    the assumption that the error term would reduce in proportion to √SE). Taking the Picture

1207    Identification test (the test most similar to previous studies) our results suggests that it would

1208    require N > 300 – i.e. over five times our current sample size. This suggests that this

1209    experimental paradigm is not sufficiently sensitive to address this hypothesis.

1210          Given the ambiguity of our findings with regard to the interaction, it is not appropriate to

1211    extensively interpret why we do not find the interaction while the previous studies did. However,

1212    we note that there are a variety of differences across the studies which could underpin the

1213    different findings, if it holds true. For example, the test of individual differences which we use is

1214    harder than that used by Perrachione et al. (2011) since it uses all six Mandarin vowels (whereas

1215    the original study used five, without /u/) and all of the Mandarin tones (where Perrachione et al.

1216    used three, without Tone 3). This change also means that that we cannot easily contrast the range

1217    of participant scores in the two studies and it may be that the spread of ability of our participant

1218    is different from theirs. In addition, our training task is potentially harder than both of the

1219    previous studies, i.e. involving all four tones in the context of natural Mandarin stimuli in the

1220    context of a word learning tasks. Finally, we also note that our statistical analyses are different

1221    from both of the previous studies in that they took their continuous aptitude measures and turned

1222    these into binary factors using a "cut off", whereas our statistical approach allows us to use them

1223    as continuous variables. However, this should in principle make our approach more powerful

1224    than in previous studies.

1225    *4.6     Future Directions*

1226          If the interaction between aptitude and training condition reported in Sadakata and

1227    McQueen (2014) and Perrachione et al., (2011) is to have implications for educational materials,

1228    it is important to establish whether it extends to other more naturalistic materials. Given the

1229    relatively small samples in these original studies, and the increasing recognition that psychology

1230    experiments have been routinely underpowered (Maxwell, Lau & Howard, 2015; and see

1231    Vasishth, Mertzen, Jäger, & Gelman, (2018) for a recent demonstration in the area of reading)

1232    and that can lead to increases in both type 1 and type 2 error, we suggest that it would be useful

1233    to implement a direct, high powered replication of these previous studies. We note that having a

1234    sufficient sample to provide substantial evidence for H1/H0 using Bayesian methods, or to obtain

1235    90% power for frequentist methods, would likely require a much larger sample than is standard

1236    in these types of studies. Given the time-consuming nature of these multiple session training

1237    studies, moving to online testing may be necessary to make this feasible (see Xie et al. 2018 for

1238    an example of an acoustic training study done over the web), or alternately multi-lab

1239    collaboration may be necessary. Note that this would also allow us to see whether the fact that

1240    Perrachione et al., (2011) found their interaction with *untrained* voices, whereas Sadakata &

1241    McQueen (2014) found it only for *trained* voices, is a true difference (due to the different

1242    paradigms) or due to power. Critically, successful replication would allow us to then extend the

1243    paradigms in such a way as to explore the factors above. For example, would increasing the

1244    number of tones to use all four Mandarin tones and/or using natural Mandarin stimuli affect the

1245    interaction between variability in the input and learner aptitude?

1246        Although direct replication will play a useful role in establishing these effects, we believe

1247    that ultimately it will also be important to develop a more nuanced approach to measuring the

1248    factors leading to different levels of aptitude both in tone learning, and in other types of phonetic

1249    learning. We note that here in addition to not seeing the predicted interaction with variability, we

1250    also didn't see interactions between aptitude and training session in any of our tasks, suggesting

1251    that our aptitude measure predicted baseline performance on the task and *not* the ability to

1252    improve due to training. In addition, the tasks used to measure "aptitude" are quite similar in

1253    nature to the training and test tasks, decreasing their explanatory value. Our ongoing work

1254    explores the combined predictive value of a range of measures including measures of attention,

1255    working memory and musical ability. Identifying factors which are predictive of aptitude for

1256    tone learning has clear implications for teaching and the personalisation of teaching methods.

1257    **5        Conclusion**

1258        We trained naive participants on all four Mandarin tones, using real language stimuli

1259    embedded in a word learning task. We found improvements in both production and perception of

1260    tones which transferred to novel voices and items. We found that learning was greatest for

1261    training with a single voice but that training with a single voice versus four voices (whether

1262    intermixed or blocked) lead to equal amounts of generalisation. Although learner aptitude

1263    predicted performance in most tasks, there was no evidence that different levels of aptitude lead

1264    to better or worse learning from different types of training input.

1265                                    **References**

1266    Aliaga-García, C., & Mora, J. C. (2009). Assessing the effects of phonetic training on L2 sound

1267          perception and production. In M. A. Watkins, A. S. Rauber, & B.O. Baptista (Eds.). *Recent*

1268          *Research in Second Language Phonetics/Phonology: Perception and Production* (pp. 2-

1269          31). Newcastle upon Tyne, UK: Cambridge Scholars Publishing.

1270    Audacity Team. (2015). Audacity (Version 2.1.1). *Computer Program*. Retrieved May, 2015,

1271          from http://audacityteam.org/

1272    Alshangiti, W., & Evans, B. G. (2014, May). Investigating the domain-specificity of phonetic

1273          training for second language learning: Comparing the effects of production and perception

1274          training on the acquisition of English vowels by Arabic learners of English. In *the*

1275          *Proceedings of the International Seminar for Speech Production, Cologne, Germany*.

1276    Apfelbaum, K. S., & McMurray, B. (2011). Using variability to guide dimensional weighting:

1277          Associative mechanisms in early word learning. *Cognitive Science*, *35*(6), 1105-1138.

1278          https://doi.org/10.1111/j.1551-6709.2011.01181.x

1279    Baayen, R. H., Davidson, D. J., & Bates, D. M. (2008). Mixed-effects modeling with crossed

1280          random effects for subjects and items. *Journal of Memory and Language*, *59*(4), 390-412.

1281          https://doi.org/10.1016/j.jml.2007.12.005

1282    Barcroft, J., & Sommers, M. S. (2005). Effects of acoustic variability on second language

1283          vocabulary learning. *Studies in Second Language Acquisition*, *27*, 387-414.

1284          https://doi.org/10.1017/S0272263105050175

1285    Barcroft, J., & Sommers, M. S. (2014). Effects of variability in fundamental frequency on L2

1286          vocabulary learning: A comparison between learners who do and do not speak a tone

1287        language. *Studies in Second Language Acquisition*, *36*(3), 423-449.

1288        https://doi.org/10.1016/j.neuropsychologia.2006.11.015

1289    Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for

1290        confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and*

1291        *Language*, *68*(3), 255-278. https://doi.org/10.1016/j.jml.2012.11.001

1292    Bates, D., Maechler, M., Bolker, B., & Walker, S. (2013). lme4: Linear mixed-effects models

1293        using Eigen and S4. R package version 1.0–5. 2013.

1294    Boersma, P., & Weenink, D. (2015). Praat: doing phonetics by computer [Computer program].

1295        Version 5.4.14, retrieved 24 July 2015 from http://www.praat.org/

1296    Bradlow, A. R., & Bent, T. (2008). Perceptual adaptation to non-native

1297        speech. *Cognition*, *106*(2), 707-729. https://doi.org/10.1016/j.cognition.2007.04.005

1298    Bradlow, A. R., & Pisoni, D. B. (1999). Recognition of spoken words by native and non-native

1299        listeners: Talker-, listener-, and item-related factors. *The Journal of the Acoustical Society*

1300        *of America*, *106*, 2074-2085. http://dx.doi.org/10.1121/1.427952

1301    Bradlow, A. R., Akahane-Yamada, R., Pisoni, D. B., & Tohkura, Y. I. (1999). Training Japanese

1302        listeners to identify English/r/and/l: Long-term retention of learning in perception and

1303        production. *Perception & Psychophysics*, *61*(5), 977-985.

1304        https://doi.org/10.3758/BF03206911.

1305    Bygate, M., Swain, M., & Skehan, P. (2013). *Researching pedagogic tasks: Second language*

1306        *learning, teaching, and testing*. London UK: Routledge.

1307    Cicchetti, D. V. (1994). Guidelines, criteria, and rules of thumb for evaluating normed and

1308        standardized assessment instruments in psychology. *Psychological Assessment*, *6*(4), 284.

1309        http://dx.doi.org/10.1037/1040-3590.6.4.284

1310 Clopper, C. G., & Pisoni, D. B. (2004). Some acoustic cues for the perceptual categorization of

1311     American English regional dialects. *Journal of Phonetics*, *32*(1), 111-140.

1312     https://doi.org/10.1016/S0095-4470(03)00009-3

1313 Cohen, A. D. (2014). *Strategies in learning and using a second language*. London UK:

1314     Routledge.

1315 Dienes, Z. (2008). *Understanding psychology as a science: An introduction to scientific and*

1316     *statistical inference*. Macmillan International Higher Education.

1317 Dienes, Z. (2014). Using Bayes to get the most out of non-significant results. *Frontiers in*

1318     *psychology*, *5*, 781. https://doi.org/10.3389/fpsyg.2014.00781

1319 Dienes, Z. (2015). How Bayesian statistics are needed to determine whether mental states are

1320     unconscious. In M. Overgaard (Ed.), *Behavioural methods in consciousness research* (pp.

1321     199–220). Oxford: Oxford University Press.

1322 Dienes, Z., Coulton, S., & Heather, N. (2018). Using Bayes factors to evaluate evidence for no

1323     effect: examples from the SIPS project. *Addiction*, *113*(2), 240-246.

1324     **https://doi.org/10.1111/add.14002**

1325 Flege, J. E., Takagi, N., & Mann, V. (1995). Japanese adults can learn to produce

1326     English/I/and/l/accurately. *Language and Speech*, *38*, 25-55.

1327     https://doi.org/10.1177/002383099503800102

1328 Giannakopoulou, A., Brown, H., Clayards, M., & Wonnacott, E. (2017). High or low?

1329     Comparing high and low variability phonetic training in adult and child second language

1330     learners. *PeerJ*, *5*, e3209. DOI:10.7717/peerj.3209

1331 Giannakopoulou, A., Uther, M., & Ylinen, S. (2013). Enhanced plasticity in spoken language

1332     acquisition for child learners: Evidence from phonetic training studies in child and adult

1333        learners of English. *Child Language Teaching and Therapy*, *29*, 201-218.

1334        https://doi.org/10.1177/0265659012467473

1335   Hay, J., Drager, K., & Thomas, B. (2013). Using nonsense words to investigate vowel

1336        merger. *English Language & Linguistics*, *17*(2), 241-269.

1337   Iverson, P., Ekanayake, D., Hamann, S., Sennema, A., & Evans, B. G. (2008). Category and

1338        perceptual interference in second-language phoneme learning: An examination of

1339        English/w/-/v/learning by Sinhala, German, and Dutch speakers. *Journal of Experimental*

1340        *Psychology: Human Perception and Performance*, *34*, 1305. https://doi.org/10.1037/0096-

1341        1523.34.5.1305

1342   Iverson, P., Hazan, V., & Bannister, K. (2005). Phonetic training with acoustic cue

1343        manipulations: A comparison of methods for teaching English /r/-/l/ to Japanese

1344        adults. *The Journal of the Acoustical Society of America*, *118*, 3267-3278.

1345        https://doi.org/10.1121/1.2062307

1346   Jaeger, T. F. (2008). Categorical data analysis: Away from ANOVAs (transformation or not) and

1347        towards logit mixed models. *Journal of Memory and Language*, *59*(4), 434-446.

1348        https://doi.org/10.1016/j.jml.2007.11.007

1349   Jeffreys, H. (1998). *The theory of probability*. OUP Oxford.

1350   Joanisse, M. F., Manis, F. R., Keating, P., & Seidenberg, M. S. (2000). Language deficits in

1351        dyslexic children: Speech perception, phonology, and morphology. *Journal of*

1352        *Experimental Child Psychology*, *77*, 30-60. https://doi.org/10.1006/jecp.1999.2553

1353   Lambacher, S. G., Martens, W. L., Kakehi, K., Marasinghe, C. A., & Molholt, G. (2005). The

1354        effects of identification training on the identification and production of American English

1355        vowels by native speakers of Japanese. *Applied Psycholinguistics*, *26*(2), 227-247.

1356        https://doi.org/10.1017/S0142716405050150

1357    Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical

1358        data. *Biometrics*, *33*(1),159-174. DOI: 10.2307/2529310

1359    Lengeris, A., & Hazan, V. (2010). The effect of native vowel processing ability and frequency

1360        discrimination acuity on the phonetic training of English vowels for native speakers of

1361        Greek. *The Journal of the Acoustical Society of America*, *128*, 3757-3768.

1362        https://doi.org/10.1121/1.3506351

1363    Lively, S. E., Logan, J. S., & Pisoni, D. B. (1993). Training Japanese listeners to identify

1364        English/r/and/l/. II: The role of phonetic environment and talker variability in learning new

1365        perceptual categories. *The Journal of the Acoustical Society of America*, *94*, 1242-1255.

1366        https://doi.org/10.1121/1.408177

1367    Logan, J. S., Lively, S. E., & Pisoni, D. B. (1991). Training Japanese listeners to identify

1368        English/r/and/l/: A first report. *The Journal of the Acoustical Society of America*, *89*, 874-

1369        886. https://doi.org/10.1121/1.1894649

1370    Magnuson, J. S., & Nusbaum, H. C. (2007). Acoustic differences, listener expectations, and the

1371        perceptual accommodation of talker variability. *Journal of Experimental Psychology:*

1372        *Human Perception and Performance*, *33*, 391. https://doi.org/10.1037/0096-1523.33.2.391

1373    Mattys, S. L., & Wiget, L. (2011). Effects of cognitive load on speech recognition. *Journal of*

1374        *Memory and Language*, *65*(2), 145-160. https://doi.org/10.1016/j.jml.2011.04.004

1375    Maxwell, O., Baker, B., Bundgaard-Nielsen, R., & Fletcher, J. (2015). A comparison of the

1376        acoustics of nonsense and real word stimuli: Coronal stops in Bengali. *Proceedings of the*

1377        *meeting of the International Congress of Phonetic Sciences,* Glasgow, UK.

1378   Maxwell, S. E., Lau, M. Y., & Howard, G. S. (2015). Is psychology suffering from a replication

1379        crisis? What does "failure to replicate" really mean? *American Psychologist*, *70*(6), 487.

1380        http://dx.doi.org/10.1037/a0039400

1381   McGraw, K. O., & Wong, S. P. (1996). Forming inferences about some intraclass correlation

1382        coefficients. *Psychological Methods*, *1*(1), 30. http://dx.doi.org/10.1037/1082-989X.1.1.30

1383   Nusbaum, H. C., & Morin, T. M. (1992). Paying attention to differences among talkers. *Speech*

1384        *perception, production and linguistic structure*, (pp. 113-134). Amsterdam, Netherlands:

1385        IOS press.

1386   Perrachione, T. K., Lee, J., Ha, L. Y., & Wong, P. C. (2011). Learning a novel phonological

1387        contrast depends on interactions between individual differences and training paradigm

1388        design. *The Journal of the Acoustical Society of America*, *130*, 461-472.

1389        https://doi.org/10.1371/journal.pone.0089642

1390   Powell, M. J. (2009). The BOBYQA algorithm for bound constrained optimization without

1391        derivatives. *Cambridge NA Report NA2009/06, University of Cambridge, Cambridge*, 26-

1392        46.

1393   Quené, H., & Van den Bergh, H. (2008). Examples of mixed-effects modeling with crossed

1394        random effects and with binomial data. *Journal of Memory and Language*, *59*(4), 413-425.

1395        https://doi.org/10.1016/j.jml.2008.02.002

1396   Ramscar, M., & Baayen, H. (2013). Production, comprehension, and synthesis: a communicative

1397        perspective on language. *Frontiers in psychology*, *4*, 233.

1398        https://doi.org/10.3389/fpsyg.2013.00233

1399    Ramscar, M., Yarlett, D., Dye, M., Denny, K., & Thorpe, K. (2010). The effects of

1400       feature-label-order and their implications for symbolic learning. *Cognitive science*, *34*(6),

1401       909-957. https://doi.org/10.1111/j.1551-6709.2009.01092.x

1402    R Development Core Team (2010). R: A Language and Environment for Statistical Computing,

1403       Version R 3.3.2. *Available at www. r-project. org. Accessed September* 2017.

1404    Sadakata, M., & McQueen, J. M. (2013). High stimulus variability in nonnative speech learning

1405       supports formation of abstract categories: Evidence from Japanese geminates. *The Journal*

1406       *of the Acoustical Society of America*, *134*, 1324-1335. https://doi.org/10.1121/1.4812767

1407    Sadakata, M., & McQueen, J. M. (2014). Individual aptitude in Mandarin lexical tone perception

1408       predicts effectiveness of high-variability training. *Frontiers in Psychology*, *5,* 1318.

1409       https://doi.org/10.3389/fpsyg.2014.01318

1410    Scarborough, R. (2012). Lexical similarity and speech production: Neighborhoods for nonwords.

1411       *Lingua*, *122*(2), 164-176. https://doi.org/10.1016/j.lingua.2011.06.006

1412    Sommers, M. S., & Barcroft, J. (2007). An integrated account of the effects of acoustic

1413       variability in first language and second language: Evidence from amplitude, fundamental

1414       frequency, and speaking rate variability. *Applied Psycholinguistics*, *28*(2), 231-249.

1415       https://doi.org/10.1017/S0142716407070129

1416    Sommers, M. S., & Barcroft, J. (2011). Indexical information, encoding difficulty, and second

1417    language vocabulary learning. *Applied Psycholinguistics*, *32*(2), 417-434.

1418    https://doi.org/10.1017/S0142716410000469

1419    Strange, W., & Dittmann, S. (1984). Effects of discrimination training on the perception of /r/-/l/

1420       by Japanese adults learning English. *Perception & Psychophysics*, *36*, 131-145.

1421       https://doi.org/10.3758/BF03202673

1422   Strange, W., & Shafer, V. L. (2008). Speech perception in second language learners: The re-

1423        education of selective perception. *Phonology and second language acquisition* (pp.153-

1424        192). Amsterdam, Netherland: John Benjamins Publishing Company.

1425   Vasishth, S., Mertzen, D., Jäger, L. A., & Gelman, A. (2018). The statistical significance filter

1426        leads to overoptimistic expectations of replicability. *Journal of Memory and Language*,

1427        *103*, 151-175. https://doi.org/10.1016/j.jml.2018.07.004

1428   Wang, Y., Jongman, A., & Sereno, J. A. (2003). Acoustic and perceptual evaluation of Mandarin

1429        tone productions before and after perceptual training. *The Journal of the Acoustical Society*

1430        *of America*, *113*, 1033-1043. https://doi.org/10.1121/1.1531176

1431   Wang, Y., Spence, M. M., Jongman, A., & Sereno, J. A. (1999). Training American listeners to

1432        perceive Mandarin tones. *The Journal of the Acoustical Society of America*, *106*, 3649-

1433        3658. https://doi.org/10.1121/1.428217

1434   Wong, J. (2014). The Effects of High and Low Variability Phonetic Training on the Perception

1435        and Production of English Vowels /e/-/æ/ by Cantonese ESL Learners with High and Low

1436        L2 Proficiency Levels. *Proceedings of the 15th Annual Conference of the International*

1437        *Speech Communication Association*, 524-528. Retrieved from

1438        https://repository.hkbu.edu.hk/hkbu_staff_publication/6234.

1439   Wong, P., & Perrachione, T. K. (2007). Learning pitch patterns in lexical identification by native

1440        English-speaking adults. *Applied Psycholinguistics*, *28*, 565-585.

1441        https://doi.org/10.1017/S0142716407070312

1442   Xie, X., Weatherholtz, K., Bainton, L., Rowe, E., Burchill, Z., Liu, L., & Jaeger, T. F. (2018).

1443        Rapid adaptation to foreign-accented speech and its transfer to an unfamiliar speaker. *The*

1444         *Journal of the Acoustical Society of America*, *143*, 2013-2031.

1445         https://doi.org/10.1121/1.5027410

1446   Yip, M. (2002). *Tone. Cambridge textbooks in linguistics*. Cambridge: Cambridge University

1447         Press.

1448   Zeromskaite, I. (2014). The potential role of music in second language learning: A review

1449         article. *Journal of European Psychology Students*, *5,* 78-88. http://doi.org/10.5334/jeps.ci

# Figure 1

Tasks completed in each of the eight sessions

This figure describes all tasks arranged through session 1-8

**SESSION 1**

1) Pitch Contour Perception Test

2) Categorisation of Synthesised Tonal Continua

3) Word Repetition

4) Three Interval Oddity

5) English Introduction

→

**SESSIONS 2- 7**

Training only

→

**SESSION 8**

1) Word Repetition

2) Three Interval Oddity

3) Picture Identification

4) Pitch Contour Perception Test

5) Categorisation of Synthesised Tonal Continua

6) Picture Naming

7) Questionnaire

# 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 for the LV (Low Variability), HV (High Variability) & HVB (High Variability Blocked) groups in Pitch Contour Perception Test. Error bars represents the 95% confidence intervals.
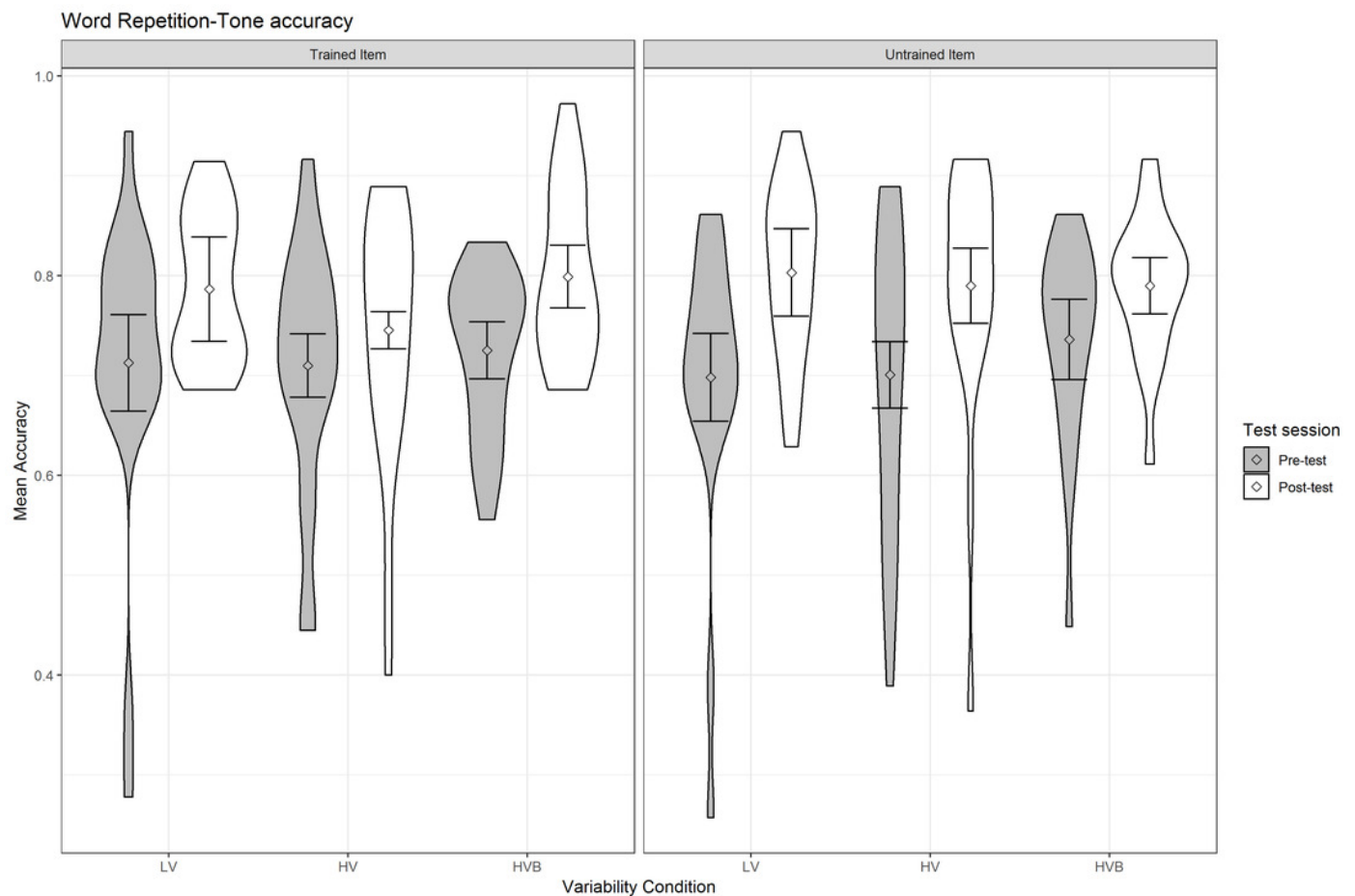
# Figure 4

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

# Figure 5

Mean accuracy in Three Interval Oddity task for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocked) training groups in Pre- and Post-tests for trained and untrained items. Error bars show 95% confidence intervals.
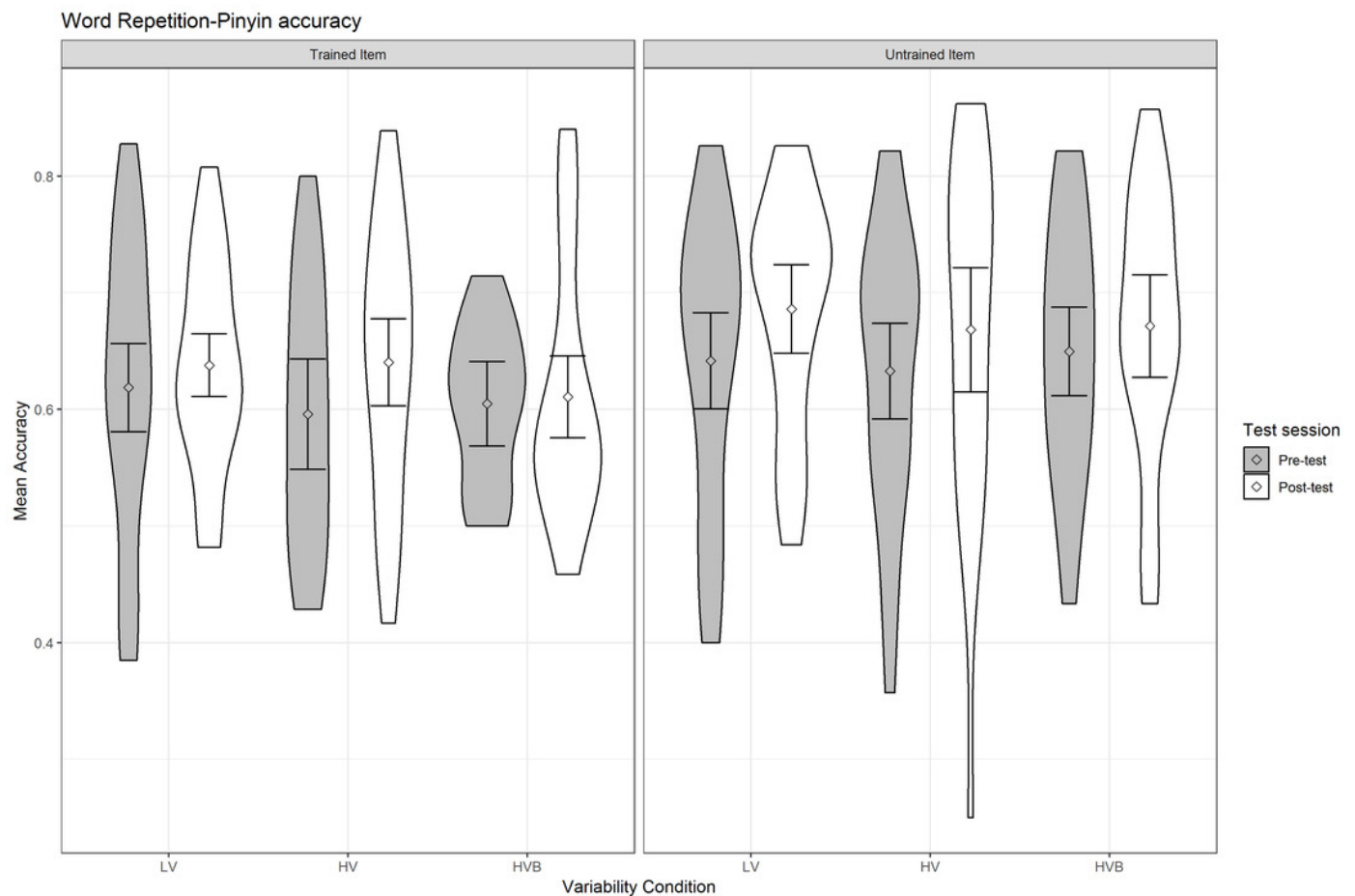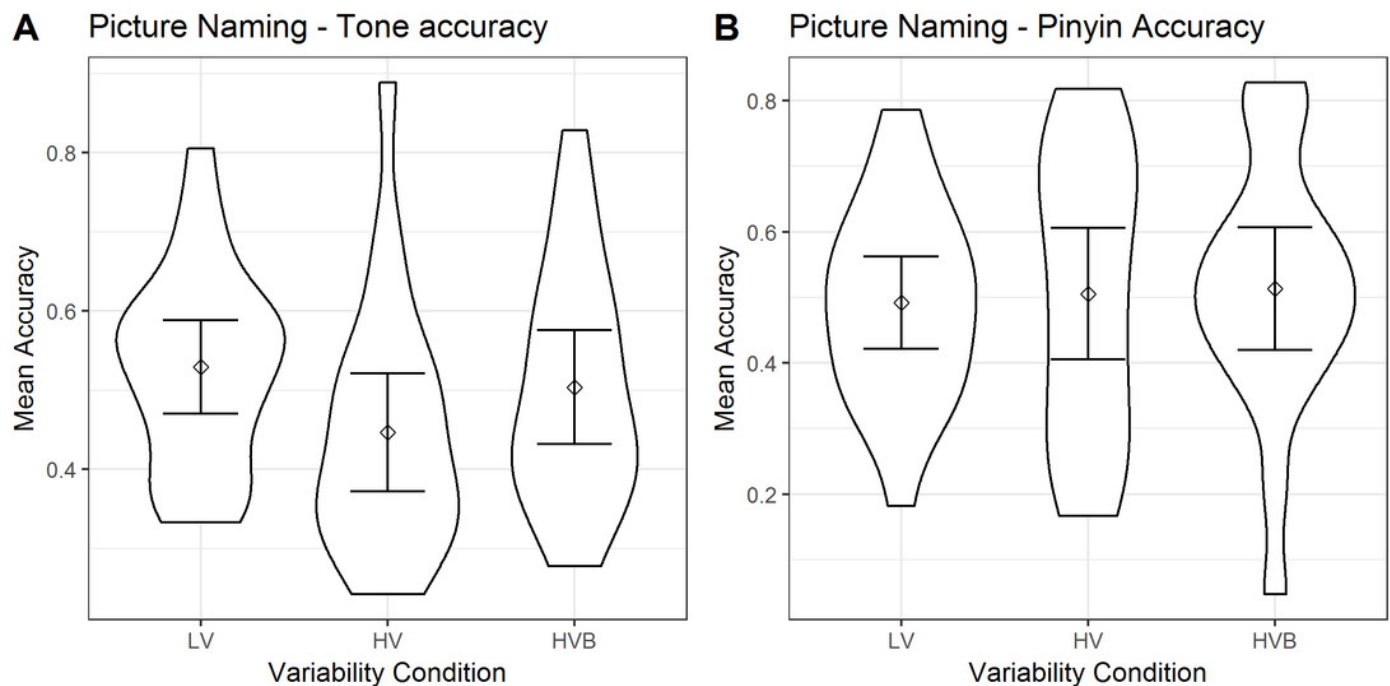
# Figure 6

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

# Figure 7

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

PeerJ

# Figure 8

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

# Figure 9

Tone accuracy and Pinyin accuracy of Picture Naming for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocked) training groups. Error bars show 95% confidence intervals.
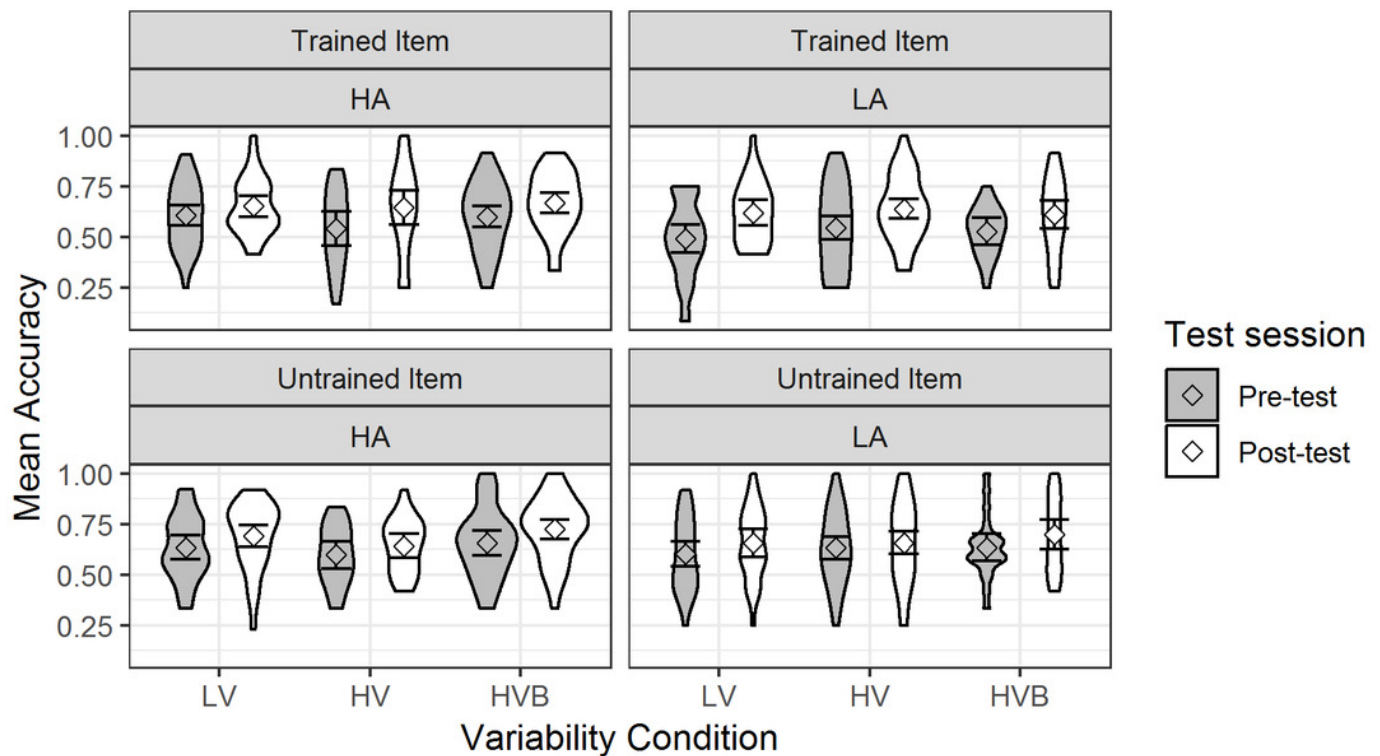
# Figure 10

Accuracy in the Three Interval Oddity and Training data for LV (Low Variability), HV (High Variability) and HVB (High Variability Blocked) training groups. Error bars show 95% confidence interval.

(A)Mean accuracy of Three Interval Oddity, split by high (HA) versus low (LA) aptitude in the

Pitch Contour Perception Test task (B) Mean accuracy of Training, split by high (HA) versus

low (LA) aptitude in the Pitch Contour Perception Test task

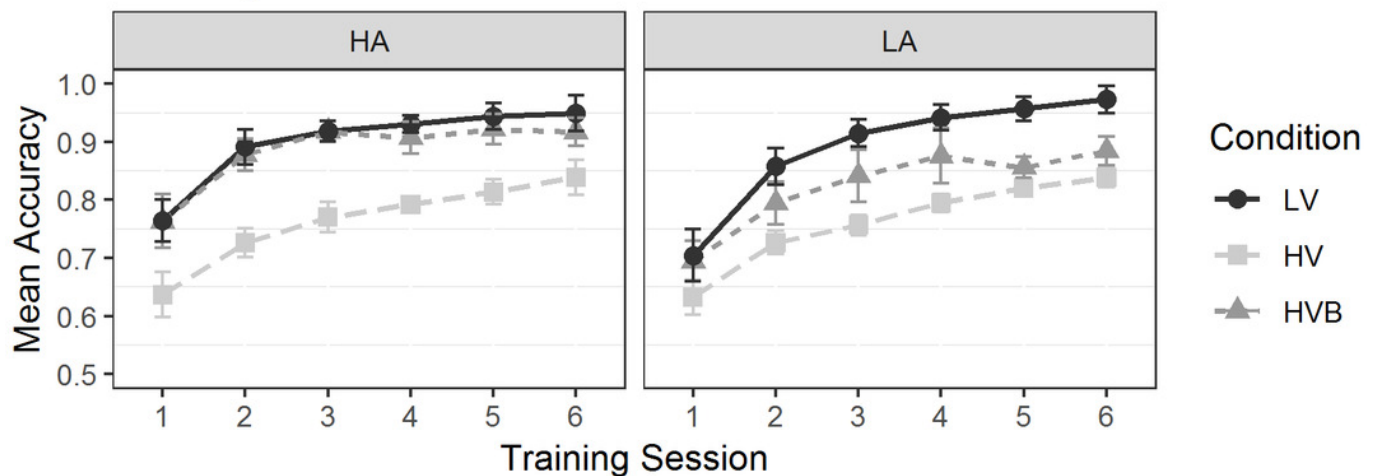**A**  Three-interval Oddity



**B**  Training

# Figure 11

Accuracy in the Picture Naming and Picture Identification data for LV , HV and HVB training groups, split by high (HA) versus low (LA) aptitude in the Pitch Contour Perception Test.

Error bars show 95% confidence interval. (A) Mean accuracy of Picture Naming tone accuracy measure (B) Scatter plot contrasting Mean accuracy of Picture Naming tone accuracy measure and corresponding aptitude measure from Picture Contour Perception Test (C) Mean accuracy of Picture Naming Pinyin accuracy measure (D) Scatter plot contrasting Mean accuracy of Picture Naming Pinyin accuracy measure and corresponding aptitude measure from Picture Contour Perception Test (E) Mean accuracy of Picture Identification (F) Scatter plot contrasting Mean accuracy of Picture Identification and corresponding aptitude measure from Picture Contour Perception Test

**A** Picture Naming - Tone accuracy



**C** Picture Naming - Pinyin Accuracy



**E** Picture Identification
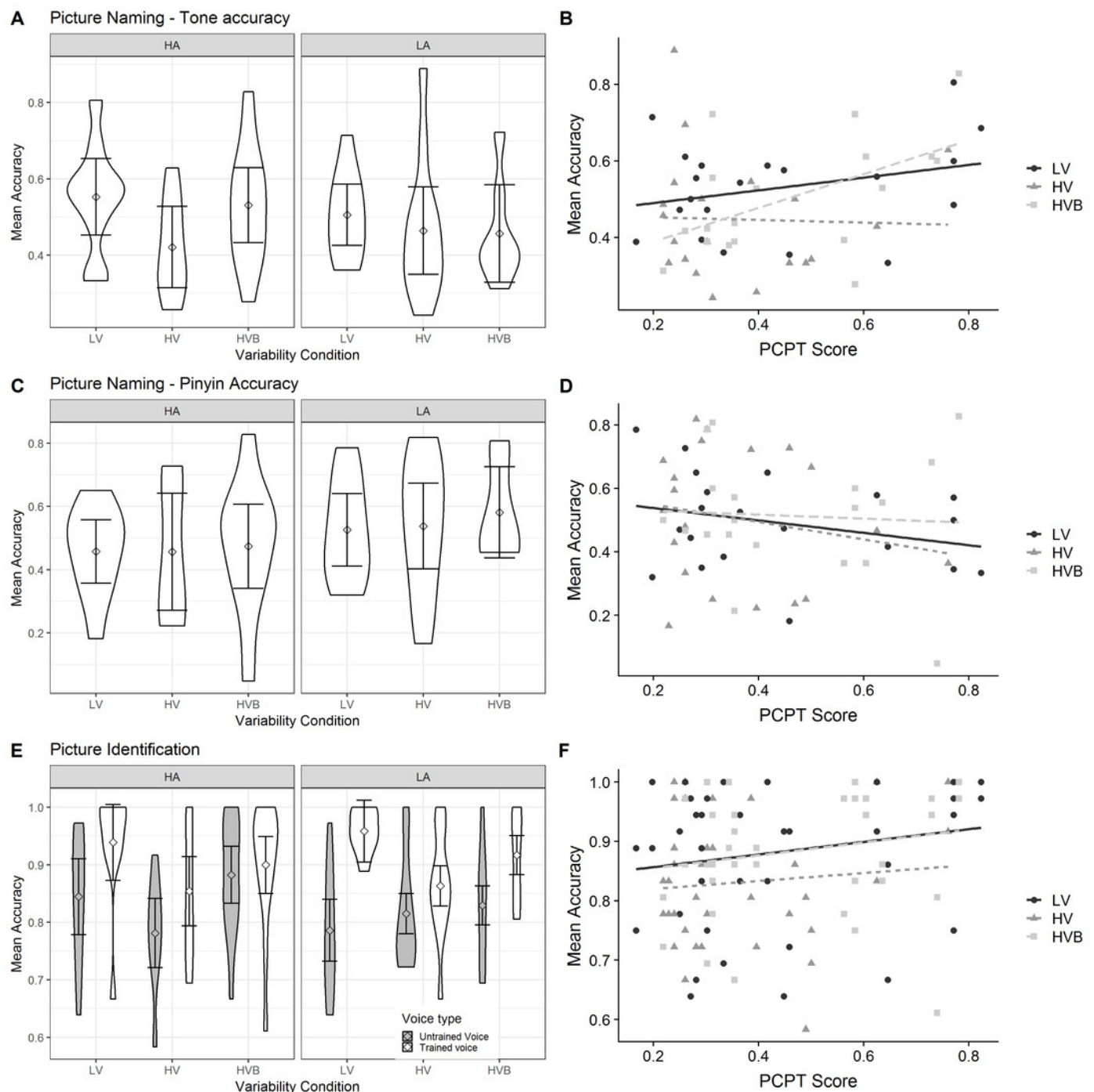


**B**



**D**



**F**

# Figure 12

Accuracy in the Word Repetition data for LV, HV and HVB training groups, split by high (HA) versus low (LA) aptitude in the Pitch Contour Perception Test. Error bars show 95% confidence intervals.

(A) Mean accuracy of Word Repetition tone accuracy measure (B) Mean accuracy of Word Repetition Pinyin accuracy measure
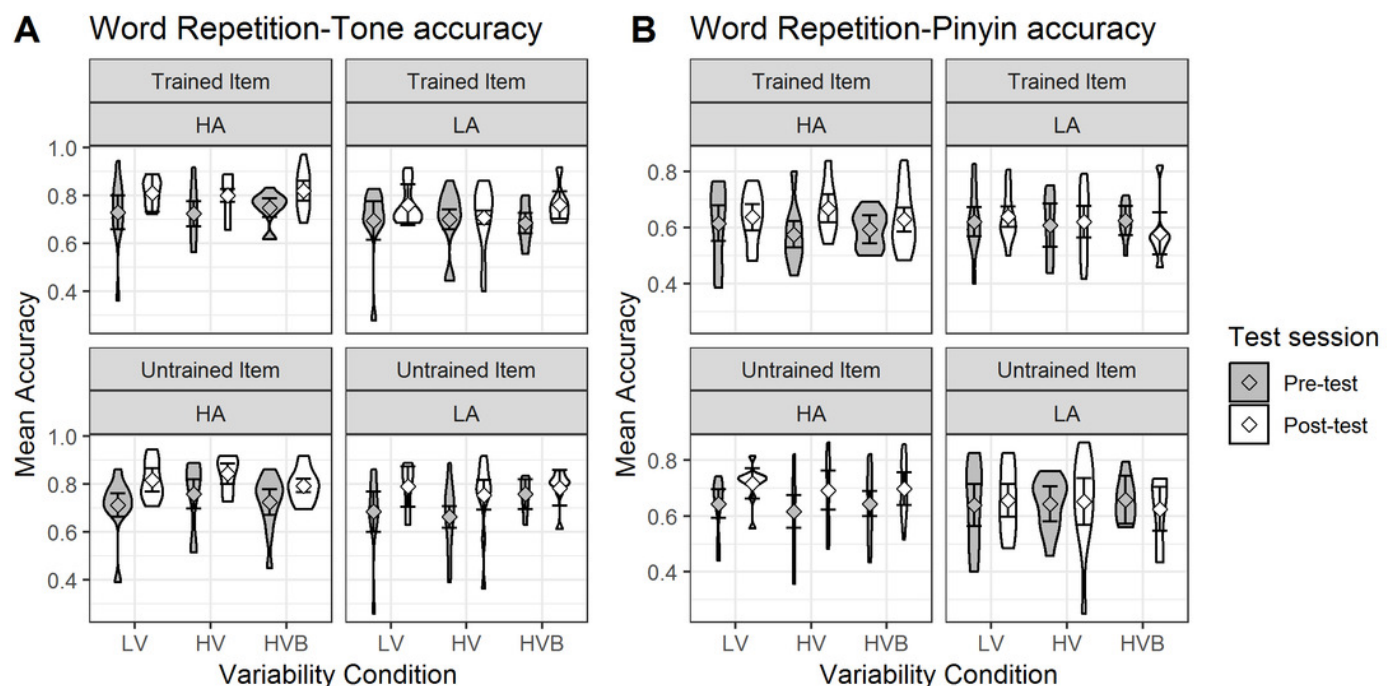
**Table 1**(on next page)

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

1

| Condition | Mean Age | Age Range | Languages Learned | Average Staring Age |
|---|---|---|---|---|
| Low Variability | 26.15 (2.2) | 19-53 | 2.7 (0.5) | 13.8 (1.1) |
| High Variability | 25.65 (0.7) | 19-47 | 2.5 (0.6) | 12.2 (0.5) |
| High Variability blocked | 22.05 (1.4) | 19-30 | 2.0 (1.3) | 11.8 (0.4) |

2

**Table 2**(on next page)

Use of trained and untrained items and voices in different tasks.

1

| Task | Items | Voice |
|---|---|---|
| Picture Identification | Trained | One trained voice (counterbalanced, see Table 3) One untrained voice (counterbalanced, see Table 3) |
| Three Interval Oddity (Pre and Post) | Trained and untrained | 4 new voices |
| Picture Naming | Trained | *NA* |
| Word Repetition (Pre and Post) | Trained and untrained | 1 trained voice (counterbalanced, see Table 3) |
| *Individual Aptitude test 1* Pitch Contour Perception Test (Pre and Post) | Vowels | 4 untrained voices |
| *Individual Aptitude test 2* Categorisation of Synthesised Tonal Continua (Pre and Post) | Synthesised voice | Synthesised voice |

2

**Table 3**(on next page)

Counterbalancing of voices across training conditions in the Picture Identification task (the only test in which trained and untrained voices are directly contrasted) and the Word Repetition tests.

1

| Task | Voice | | | | |
|---|---|---|---|---|---|
| | Version 1 | Version 2 | Version 3 | Version 4 | Version 5 |
| Training, LV | F1 | F2 | F3 | M1 | M2 |
| Training, HV/HVB | F1 | F2 | F3 | M1 | M2 |
| | F3 | F1 | M2 | F1 | F2 |
| | M1 | M1 | F1 | F2 | F3 |
| | M2 | M2 | F2 | F3 | M1 |
| Picture Identification | | | | | |
| *Trained voice* | F1 | F2 | F3 | M1 | M2 |
| *Untrained voice* | F2 | F3 | M1 | M2 | F1 |
| Word Repetition | F1 | F2 | F3 | M1 | M2 |

2

**Table 4**(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

| Data Set | Coefficient Name | Statistics |
|---|---|---|
| *Word Repetition:* | **Aptitude** | **β = 0.07, SE = 0.03, z = 2.35, p = .019** |
| *Tone Accuracy* | Aptitude by *Test-Session* | β = 0.03, SE = 0.04, z = 0.72, p = .473 |
| *(Pre/Post)* | Aptitude by LV-HV Contrast by *Test-Session* | β = 0.05, SE = 0.11, z = 0.47, p = .639 |
| | Aptitude by LV-HVB Contrast by *Test-Session* | β = 0.13, SE = 0.10, z = 1.35, p = .176 |
| | Aptitude by LV-HV Contrast by *Test-Session* by *Item-Novelty* | β = -0.14, SE = 0.15, z = -0.97, p = .334 |
| | Aptitude by LV-HVB Contrast by *Test-Session* by *Item-Novelty* | β = 0.07, SE = 0.13, z = 0.50, p = .61 |
| *Three Interval* | **Aptitude** | **β = 0.07, SE = 0.03, z = 2.19, p = .029** |
| *Oddity* | Aptitude by *Test-Session* | β = 0.01, SE = 0.23, z = 0.31, p = .757 |
| *(Pre/Post)* | Aptitude by LV-HV Contrast by *Test-Session* | β = 0.05, SE = 0.07, z = 0.77, p = .443 |
| | Aptitude by LV-HVB Contrast by *Test-Session* | β = 0.05, SE = 0.06, z = 0.83, p = .410 |
| | Aptitude by LV-HV Contrast by *Test-Session* by *Item-Novelty* | β = -0.12, SE = 0.13, z = -0.94, p = .346 |
| | Aptitude by LV-HVB Contrast by *Test-Session* by *Item-Novelty* | β = 0.06, SE = 0.11, z = 0.52, p = .604 |
| *Training* | **Aptitude** | **β = 0.13, SE = 0.048, z = 2.70, p = .007** |
| | Aptitude by LV-HV Contrast | β = -0.04, SE = 0.11, z = -0.332, p = .740 |
| | Aptitude by LV-HV Contrast | β = 0.03, SE = 0.10, z = 0.26, p = 0.795 |
| *Picture* | **Aptitude** | **β = 1.48, SE = 0.08, z = 1.96, p = .050** |
| *Identification* | Aptitude by Voice Novelty | β = -0.03, SE = 0.07, z = -0.33, p = .745 |

| | | |
|---|---|---|
| *(Post Only)* | Aptitude by LV-HV Contrast | β = -0.02, SE = 0.19, z = -0.12, p = .901 |
| | Aptitude by LV-HVB Contrast | β = 0.01, SE = 0.17, z = 0.09, p = .932 |
| | Aptitude by LV-HV Contrast by *Voice-Novelty* | β = 0.35, SE = 0.21, z = 1.63, p = .103 |
| | Aptitude by LV-HVB Contrast by *Voice-Novelty* | β = -0.11, SE = 0.19, z = -0.58, p = .566 |
| *Picture Naming:* | **Aptitude** | **β = 0.08, SE = 0.04, z = 1.89, p = .0059** |
| *Tone Accuracy* | Aptitude by LV-HV Contrast | β = -0.09, SE = 0.11, z = -0.84, p = .402 |
| | Aptitude by LV-HVB Contrast | β = 0.12, SE = 0.10, z = 1.22, p = .224 |

2

**Table 5**(on next page)

Bayes Factor results testing the hypothesis that there is greater generalisation following either of the high variability training conditions than the low variability condition

1

| Contrast | Mean difference | Stand. Error | H1 estimate $x$ | Bayes Factor ($B$) | Robustness Region |
|---|---|---|---|---|---|
| Picture ID (Novel voice only) *HV+ HVB > LV* | 0.13 | 0.228 | 1.71 | 0.219 | 1.11 : ∞ |
| Picture Naming, (Tone accuracy) *HV+ HVB > LV* | -0.225 | 0.168 | 1.076 | 0.067 | 0.202 : ∞ |
| Picture Naming (Pinyin Accuracy) *HV+ HVB > LV* | 0.104 | 0.196 | 4.05 | 0.08 | 0.101 : ∞ |
| Word Repetition (Tone accuracy) *test-session* by *HV+ HVB > LV* | -0.108 | 0.157 | 0.395 | 0.239 | 0.303 : ∞ |
| Word Repetition (Pinyin accuracy) *test-session* by *HV+ HVB > LV* | 0.095 | -0.034 | 0.152 | 0.421 | 0 : 0.202 |
| Three Interval Oddity *test-session* by *HV+ HVB > LV* | -0.001 | 0.1 | 0.31 | 0.303 | 0.303 : ∞ |

**Table 6**(on next page)

Bayes Factor results testing the hypothesis that there is an interaction between aptitude and variability-condition greater generalisation following either of the high variability training conditions than the low variability condition

1

| Contrast | Mean difference | Stand. Error | H1 estimate *x* | Bayes Factor (*B*) | Robustness Region |
|---|---|---|---|---|---|
| ID, (Tone accuracy)<br>*aptitude* by HV+ HVB > LV | 0.006 | 0.127 | 0.171 | 0.617 | 0: 0.354 |
| Picture Naming, (Tone accuracy)<br>*aptitude* by HV+ HVB > LV | 0.042 | 0.083 | 0.099 | 0.904 | 0: 0.354 |
| Three Interval Oddity (Tone accuracy)<br>*aptitude* by *test-session* by HV+ HVB > LV | 0.048 | 0.05 | 0.345 | 0.371 | 0: 0.354 |
| Word Repetition (Tone accuracy)<br>*aptitude* by *test-session* by HV+ HVB > LV | 0.091 | 0.082 | 0.379 | 0.654 | 0: 0.758 |
| Training<br>*aptitude* by HV > LV | -0.037 | 0.119 | 0.129 | 0.572 | 0 : 0.253 |
| Training<br>*aptitude* by HVB > LV | 0.026 | 0.101 | 0.129 | 0.732 | 0 : 0.354 |