# Accepted Manuscript

Successful non-native speech perception is linked to frequency following response phase consistency

Akihiro Omote, Kyle Jasmin, Adam Tierney

Please cite this article as: Omote A, Jasmin K, Tierney A, Successful non-native speech perception is linked to frequency following response phase consistency, *CORTEX* (2017), doi: 10.1016/ j.cortex.2017.05.005.

**Title:** Successful non-native speech perception is linked to frequency following response phase consistency

**Authors:**  Akihiro Omote[1], Kyle Jasmin[1], Adam Tierney[1†]

## Affiliations

[1]Department of Psychological Sciences, Birkbeck, University of London, London, UK

[†]Corresponding author
Adam Tierney, Ph.D.
Birkbeck, University of London
Malet Street, London, WC1E 7HX
United Kingdom
Email:  a.tierney@bbk.ac.uk
Phone:  +44-020-7631-6368

1 **Abstract**

2 Some people who attempt to learn a second language in adulthood meet with greater success
3 than others. The causes driving these individual differences in second language learning skill
4 continue to be debated. In particular, it remains controversial whether robust auditory perception
5 can provide an advantage for non-native speech perception. Here, we tested English speech
6 perception in native Japanese speakers through the use of frequency following responses, the
7 evoked gamma band response, and behavioral measurements. Participants whose neural
8 responses featured less timing jitter from trial to trial performed better on perception of English
9 consonants than participants with more variable neural timing. Moreover, this neural metric
10 predicted consonant perception to a greater extent than did age of arrival and length of residence
11 in the UK, and neural jitter predicted independent variance in consonant perception after these
12 demographic variables were accounted for. Thus, difficulties with auditory perception may be one
13 source of problems learning second languages in adulthood.

14 **Keywords:** auditory, English, FFR, Japanese, speech

15 **1. Introduction**

16 Speaking and understanding a second language is a vital skill in an increasingly globalized world.
17 However, learning a second language poses difficulties that surpass those experienced in learning
18 a first language. Native Japanese speakers, for example, struggle to discriminate English /l/ and
19 /r/ (Goto, 1971; Miyawaki et al., 1975). Nevertheless, the difficulties which non-native speech
20 perception presents can be overcome. Native Japanese speakers, for example, through
21 experience (MacKain et al., 1981; Flege et al., 1995; Ingvalson et al. 2011) and training (Logan et
22 al., 1990; Lively et al. 1993, 1994; Bradlow et al., 1996, 1999; McCandliss et al., 2002; Iverson et
23 al., 2005; Lim and Holt, 2011) can learn to perceive and produce the distinction between /l/ and
24 /r/ with near-native accuracy. However, there are large individual differences in the degree to
25 which non-native speech sound categories can be successfully acquired: some people achieve
26 approximately native perception and production, while others produce heavily accented speech
27 and struggle to perceive non-native speech even after extensive training (Wong and Perrachione,
28 2007; Golestani and Zatorre, 2009; Perrachione et al., 2011; Hanulíková et al., 2012; Kempe et al.,
29 2012, 2015). Understanding the source of these individual differences would be an important step
30 towards the development of tools to boost non-native speech perception.

31 Learning a non-native speech sound category requires highly precise perception of durational,
32 pitch, and spectral information. One possible source of difficulties with non-native speech
33 perception, therefore, is imprecise auditory perception. Supporting this theory, individual
34 differences in non-native speech perception have been linked to non-verbal auditory perception
35 skills, including amplitude envelope discrimination (Kempe et al., 2012), frequency discrimination
36 (Lengeris and Hazan, 2010), pitch perception (Wong and Perrachione, 2007; Perrachione et al.,
37 2011), and spectral discrimination (Kempe et al., 2015). However, electrophysiology research has
38 supported a speech-specific source for non-native speech perception difficulties. Díaz et al. (2008,
39 2015), for example, found that non-native speech perception ability was linked to neural
40 discrimination of speech sounds but not non-verbal sounds differing in duration or frequency.

41 This link between speech sound discrimination and individual differences in non-native speech
42 perception has been replicated across languages (Garcia-Sierra et al., 2011; Jakoby et al., 2011;
43 Zhang et al., 2009).

44 Here we examine the link between non-native speech sound perception and auditory processing
45 in Japanese adults learning English as a second language using frequency-following responses
46 (FFRs), an electrophysiological response which reproduces the frequencies present in the evoking
47 sound and reflects early auditory processing in the brainstem and cortex (Coffey et al., 2016). The
48 FFR features high test-retest reliability (Hornickel et al. 2012) and reflects neural origins in the
49 brainstem and cortex (Coffey et al., 2016), making it an excellent measure of the robustness of
50 early auditory processing. The precision of FFRs has been linked to individual differences in the
51 development of language skills in children (Hornickel and Kraus, 2013; White-Schwoch et al.,
52 2015), but it remains unknown how FFR precision relates to second language acquisition.
53 Recently, Krizman et al., (2012) reported that bilingual FFRs more robustly encoded the
54 fundamental frequency (F0) of synthesized speech. Here, therefore, we predicted that non-native
55 speech perception ability would relate to F0 phase-locking. Given that impaired gamma-rate
56 phase-locking has also been shown to characterize children with language impairment (Heim et
57 al., 2011), we additionally investigated relationships between gamma phase-locking and non-
58 native speech perception.

59 **2. Methods**

60 *2.1 Participants*

61 Participants were 25 native Japanese speakers (13 female, aged 19 to 35 (M = 29.3, SD = 4.5))
62 with English learning experience at secondary school level or above in Japan. Participants were
63 required to have arrived in the UK after the age of 18 and to have been resident there for at least
64 1 month at the time of testing. Secondary inclusion criteria included normal audiometric
65 thresholds (<= 25 dB HL for octaves from 250 to 8000Hz) and lack of diagnosis of a language
66 impairment. Participants received a mean (sd) score of 7.6 (4.1) on the Musical Experience
67 portion of the Goldsmiths Musical Sophistication Index (Müllensiefen et al., 2014), indicating low
68 levels of musical training. Mean age of arrival in the UK was 27.8 (4.9) years, and mean duration
69 of residence in the UK was 2.6 (3.1) years. The Ethics Committee in the Department of
70 Psychological Sciences at Birkbeck, University of London approved all experimental procedures.
71 Informed consent was obtained from all participants. Participants were compensated £14 for
72 their participation.

73 *2.2 Behavioral Measures*

74 English speech perception was tested using the Receptive Phonology Test (Slevc and Miyake,
75 2006). Each question in this test is designed to assess a phonological contrast in English with
76 which Japanese subjects have difficulty. The test contains three main sections. In the *word* sub-
77 test, participants see a list of 26 word pairs which differ in a single speech sound (e.g., "late/rate").
78 Participants then hear a list of words and are asked to indicate which of the two words they heard.
79 In the *sentence* sub-test, participants see a list of 25 sentences, with one of the words replaced
80 with a word differing in a single speech sound (e.g., "My sister loves to play with crowns/clowns.")

81  Participants then hear a list of sentences and are asked to circle the word that they heard. Finally,
82  participants listen to a short story and are given a written version of the story that includes 42
83  underlined words. Participants are asked to circle any of the underlined words that are
84  mispronounced.

85  Because the original version of the Receptive Phonology Test featured a speaker of American
86  English, test materials were re-recorded by a native speaker of British English (Received
87  Pronunciation) in soundproof room with a RODE NT1-A Condenser Microphone. 3 of the items
88  from the original test were removed, as they feature speech sound contrasts which do not exist in
89  British Received Pronunciation. Audio recordings were presented to participants using Sennheiser
90  HD 25-1-II headphones. See Table 1 for a list of all of the speech sound contrasts included in the
91  test.

92  **2.3 Electrophysiology**

93  *2.3.1 Stimuli*

94  Participants were presented with two 170-ms synthesized speech sounds [la] and [ra]. These
95  syllables were synthesized using a Klatt synthesizer, as implemented in Praat (Boersma and
96  Weenink, 2016). The two syllables differed only during the first 70 ms, during which each had a
97  unique frequency trajectory for the third formant (F3). For [la], F3 was steady at 3400 Hz from 0
98  to 30 ms, then decreased linearly to 2530 Hz by 70 ms. For [ra], F3 was steady at 1601 Hz from 0
99  to 30 ms, then increased to 2530 Hz by 70 ms. All other stimulus characteristics were identical
100  across stimuli. F1 was steady at 478 Hz from 0 to 30 ms then increased to 705 Hz by 70 ms. F2 was
101  steady at 1088 Hz from 0 to 30 ms then decreased to 1035 Hz by 70 ms. From 70 to 170 ms F1, F2,
102  and F3 were steady at 705, 1035, and 2530 Hz, respectively. F0 and F4 were constant throughout
103  the stimulus at 100 Hz and 3850 Hz. A cosine off ramp with a duration of 20 ms was used to avoid
104  transients. Figure 1 displays waveforms and spectrograms for the two stimuli.

105  *2.3.2 Recording parameters*

106  During electrophysiological testing participants sat in a comfortable chair in a soundproof booth
107  with negligible ambient noise and read a book of their choice. Stimuli were presented through
108  Etymotic ER-II earphones in alternating polarity at 80 +/- 1 dB SPL to both ears with an inter-
109  onset interval of 251 ms. 6300 trials were collected for each stimulus, and stimuli were presented
110  in blocks (i.e. all [ra] trials were collected in a single block).  Electrophysiological data were
111  recorded in LabView 2.0 (National Instruments, Austin, TX) using a BioSEMI Active2 system via the
112  ActiABR module with a sample rate of 16,384 Hz and an online bandpass filter (100-3000 Hz, 20
113  dB/decade). The active electrode was placed at Cz, the grounding electrodes CMS and DRL were
114  placed on the forehead at FP1 and FP2, and the reference electrodes were placed on the earlobes.
115  Earlobe references were not electrically linked during data collection. Offset voltage for all
116  electrodes was kept below 50 mV.

117  *2.3.3 Data reduction*

118  Electrophysiological data reduction was conducted in Matlab R2016a. Offline amplification was
119  applied in the frequency domain for 3 decades below 100 Hz with a 20 dB rolloff per decade. The

120 data was organized into epochs 40 ms before through 210 ms after the onset of the stimulus and
121 baseline corrected. To ensure against contamination by electrical noise a second-order IIR notch
122 filter with a Q-factor of 100 was used with center frequencies of 50, 150, 250, 350, 450, and 550
123 Hz. A bandpass filter (0.1–2000 Hz, 12dB/oct) was then applied to the continuous EEG recording,
124 and epochs exceeding +/-100 μV were rejected as artifacts. The first 2,500 artifact-free responses
125 to each stimulus polarity then were selected for further analysis.

126 *2.3.4 Data analysis (>70 Hz)*

127 To investigate the precision of neural sound encoding we calculated inter-trial phase-locking. This
128 measure involves calculating the phase consistency at a particular frequency across trials and,
129 therefore, no averaging is necessary. This procedure provides information similar to spectral
130 analysis of average waveforms, but with a higher signal-to-noise ratio and less susceptibility to
131 artifact (Zhu et al. 2013).

132 All electrophysiological data analysis was conducted in Matlab 2016a. Parameters for FFR analysis
133 were used for frequencies >70 Hz, in accordance with the standards of previous research on
134 speech FFRs (Bidelman and Krishnan, 2009, Parbery-Clark et al., 2009). For FFR analysis (> 70 Hz),
135 phase-locking was calculated within 40-ms windows that were applied repeatedly across the
136 epoch with a 1 ms step size.  First, for each trial, a Hanning windowed fast Fourier transform was
137 calculated. Second, for each frequency, the resulting vector was transformed into a unit vector.
138 Third, all of the unit vectors were averaged. The length of the resulting vector—ranging from 0
139 (no phase consistency) to 1 (perfect phase consistency)—was then calculated as a measure of
140 cross-trial phase consistency. Phase locking factors for [la] and [ra] were averaged together to
141 form a global estimate of an individual's inter-trial phase locking.

142 This time-frequency data was then averaged in the following manner. First, data were collapsed
143 across the entire response (10 to 170 ms). Phase-locking at the fundamental frequency (100 Hz)
144 and the second through sixth harmonics was measured by extracting the maximum phase-locking
145 value in a 40-Hz bin centered on each frequency. (Harmonics above 600 Hz were not consistently
146 represented in every single participant and were therefore excluded.) Phase-locking at the
147 harmonics was averaged together to form a general measurement of harmonic encoding. In
148 addition, phase-locking was measured separately in the response to the consonant (10 to 80 ms)
149 and the response to the vowel (80-170 ms).

150 *2.3.5 Data analysis (<70 Hz)*

151 For lower-frequency analysis (< 70 Hz), phase-locking was calculated within 80-ms windows with a
152 1 ms step size. Visual inspection of the cross-subject average (see Figure 2) revealed an increase in
153 phase-locking over baseline between 0 and 60 ms. Gamma phase-locking was quantified,
154 therefore, as the average phase-locking within a window reaching from 0 to 60 ms and between
155 30 and 70 Hz.

156 *2.3.6 Statistical analyses*

157 Linear models of the behavioral and neural data were constructed using the lm() function with the
158 software package 'R', and model comparisons were performed with the anova() function. For

159  comparisons of correlations that shared one variable in common (Steiger, J.H. 1980), the r.test()
160  function in the 'psych' package from 'R' was used.

161  **3. Results**

162  First we tested whether the ability to discriminate English consonants was related to our neural
163  measures. Better performance (greater proportion correct items) on the consonant discrimination
164  items of the Phonology Test was associated with greater phase-locking to F0 ($R^2$=.379,
165  $F(1,23)$=14.03, p=.001) and with greater phase-locking within the gamma band ($R^2$=.21,
166  $F(1,23)$=6.11, p=.021). Vowel errors were not associated with F0 phase locking ($R^2$=.053,
167  $F(1,23)$=1.30, p=.27) or gamma phase locking ($R^2$=.000), and phase-locking to the harmonics (H2-
168  H6) did not correlate with performance on consonant items ($R^2$=.03, $F(1,23)$=.78, p=.34) or vowel
169  items ($R^2$=.025, $F(1,23)$=.059, p=.45). The correlation between phase-locking at F0 and consonant
170  perception was significantly greater than the correlation with vowel perception (T=2.76, p=.011);
171  similarly, the correlation between gamma phase-locking and consonant performance was
172  significantly greater than the correlation with vowel perception (T=2.95, p =.007). The correlation
173  between consonant perception and phase-locking at F0 was significantly greater than the
174  correlation with phase-locking at the higher harmonics (T=2.81, p =.01).  Figure 2 displays phase-
175  locking for the cortical evoked response and FFR across all subjects. Figure 3 displays cortical and
176  FFR phase-locking for good and poor perceivers of English consonants (top-bottom split). Figure 4
177  is a scatterplot displaying FFR phase-locking and cortical phase-locking versus consonant
178  perception performance.

179  One possible explanation for this relationship between English speech perception and F0 phase-
180  locking is that greater familiarity with English speech leads to enhanced encoding of neural
181  responses to English speech sounds. If so, one would expect the relationship between English
182  consonant perception and F0 phase-locking to be limited to the response to the consonant, which
183  did not overlap with any Japanese speech sound. On the other hand, if our results reflect a more
184  general relationship between precise auditory encoding and non-native speech perception, then
185  English consonant perception should also relate to F0-phase-locking in the response to the vowel,
186  which contained formant frequencies appropriate for a Japanese [a] (Nishi et al., 2008). We found
187  that F0 phase-locking in the response to the consonant (10-80 ms) correlated with performance
188  on consonant items ($R^2$ = 0.426, p = 0.001). F0 phase-locking in the response to the vowel (80-
189  170) also correlated with performance on consonant items ($R^2$ = 0.260, p = 0.009). Moreover, the
190  relationship between consonant perception and F0 phase-locking did not significantly differ
191  between these two portions of the response (T=0.97, p = 0.34).

192  To further test whether confounding effects of language experience could explain our results,
193  "Age Arrived in UK" and "Years in UK" were used to assess the extent of participants' experience
194  with English. "Years in UK" was cube root-transformed to bring its distribution closer to normality
195  (Shapiro-Wilk W=.89, p>.01 after transformation). Subjects who were older when they arrived in
196  the UK made more consonant errors, although the correlation was only marginally significant
197  ($R^2$=0.15, $F(1,23)$=4.02, p=.057). Age Arrived in UK also correlated negatively with F0 phase locking
198  ($R^2$=.17, $F(1,23)$=4.77, p=.039), as well as gamma phase locking ($R^2$=.25, $F(1,23)$=7.71, p=.01). The
199  number of years subjects had spent in the UK prior to testing was correlated with F0 phase

200  locking ($R^2$=.31, F(1,23)=7.51, p=.004), but not with gamma phase locking ($R^2$=.014, F(1,23)=.337,
201  p=.57).

202  To assess whether our neural measures predicted variance in phonological competence that
203  could not be simply explained by experience, we fit two linear models: one with age of arrival in
204  the UK and years residence in the UK predicting consonant performance (the "Experience Only"
205  model), and another which also included the consistency of the neural response (F0 phase
206  locking; the "Experience plus Neural model"). The two predictors in the Experience Only model
207  together accounted for 25% of the variance on consonant performance. The Experience plus
208  Neural model with F0 phase locking as a predictor performed significantly better than the
209  Experience Only model (F(1,21)=5.43, p=.030), with the F0 phase-locking predictor accounting for
210  an additional 15% of the variance for consonant performance. Including gamma phase locking as
211  an additional predictor only accounted for an additional 1.5% of the variance, and this reduction
212  in error was not significant (p=.50).

213  Finally, to investigate links between individual differences in low-frequency and high-frequency
214  phase-locking, we compared phase-locking in the gamma band to phase-locking in the FFR at F0
215  and the harmonics. Gamma phase-locking was correlated with phase-locking at both F0 ($R^2$=.31,
216  p=.004) and the harmonics ($R^2$=0.17, p=.039).

217  **4. Discussion**

218  Here we examined English speech perception and neural sound encoding in twenty-five native
219  speakers of Japanese who moved to the United Kingdom as adults. We found that English
220  consonant perception was linked to the degree of phase-locking to the fundamental frequency of
221  the frequency-following response (FFR) to sound and to phase-locking within the gamma band.
222  Vowel perception, however, did not relate to neural phase-locking. The relationship between
223  these neural metrics and English speech perception ability remained significant even after time in
224  the UK and age of arrival were controlled for.

225  That FFR phase-locking relates to second language speech perception suggests that difficulties
226  with auditory perception can interfere with the acquisition of non-native speech sound
227  categories. On the other hand, we found that non-native vowel perception was not linked to FFR
228  phase-locking, suggesting that vowel perception may depend less on the precision of auditory
229  processing. These findings support previous behavioral research demonstrating relationships
230  between non-native speech perception and auditory abilities including amplitude envelope
231  discrimination (Kempe et al., 2012), frequency discrimination (Lengeris and Hazan, 2010), and
232  spectral discrimination (Kempe et al., 2015). However, language learning is a complex process,
233  and there are likely many ways in which foreign language learning can be disrupted. Only a
234  portion of children with reading impairment, for example, display problems with auditory
235  perception (Ramus et al., 2003), and the causes of adult language learning difficulty are likely to
236  be similarly heterogenous. FFR phase-locking may be a useful metric to help identify people
237  whose difficulties with non-native language perception stem from auditory impairments.

238  These findings support and extend previous work demonstrating links between the precision of
239  neural sound encoding, language skill, and language experience. Krizman et al. (2015), for

240 example, found that in Spanish-English bilinguals degree of bilingual experience was linked to the
241 strength of fundamental frequency (F0) encoding in the FFR. Here we replicate this relationship in
242 native speakers of Japanese learning English as a second language, and extend this finding by
243 showing that this same neural metric can also explain individual differences in non-native speech
244 perception, even after language experience is accounted for. Hornickel and Kraus (2013)
245 demonstrated that the inter-trial consistency of the FFR is linked to individual differences in
246 language skills in school-age children; here we show that precise neural encoding of sound is
247 linked to successful adult language learning as well. Chandrasekaran et al. (2012) showed that the
248 robustness of FFR pitch encoding can predict subsequent short-term learning of lexical tones;
249 here we show that FFR phase-locking is linked to long-term language learning of non-tonal speech
250 sounds.

251 What is the mechanism underlying this relationship between FFR phase-locking and non-native
252 speech perception ability? One possibility is that FFR phase-locking reflects the precision of
253 temporal perception. FFR phase-locking has been linked to the ability to precisely synchronize
254 movements with sound onsets (Tierney and Kraus, 2013, 2016; Woodruff Carr et al., 2016). This
255 suggests that precise tracking of sound timing relies upon consistent auditory neural timing, as
256 synchronization places stringent demands upon the precision of auditory time perception (on the
257 order of a few milliseconds; Repp, 2000). The ability to track sound timing is also vital for speech
258 perception, as the temporal information contained in the speech envelope contains information
259 relevant to speech sound discrimination (Rosen, 1992); in fact, discrimination of speech sounds is
260 possible even if spectral information is greatly reduced (Shannon et al., 1995). Moreover, non-
261 native speech perception may rely more upon temporal information than does native speech
262 perception. For example, Japanese adults have a strong bias towards the use of temporal
263 information such as closure duration and formant transition duration when distinguishing [la] and
264 [ra], whereas native English speakers rely more heavily upon the frequency of the third formant
265 (Iverson et al., 2005).

266 We replicate the finding of Krizman et al. (2012) that F0 encoding in the FFR is related to degree
267 of bilingual experience but encoding of the harmonics is not. Moreover, we show that phase-
268 locking at the F0 but not the harmonics is also linked to non-native speech perception ability. The
269 specificity of this relationship was predicted based on these previous findings, but the underlying
270 mechanism remains unclear. One possibility is that this result reflects a relationship between non-
271 native speech perception ability and cortical auditory encoding. There is strong evidence that
272 frequency-following responses at 250 Hz and above are generated within the auditory brainstem,
273 as cooling the inferior colliculus in cats abolishes the scalp-recorded FFR (Smith et al. 1975) and
274 patients with inferior colliculus lesions do not display an FFR (Sohmer et al. 1977). However, both
275 of these studies included no stimuli below 250 Hz, and recent work has suggested that the FFR at
276 100 Hz is generated within multiple sources, including both cortical and subcortical regions
277 (Coffey et al., 2016). Thus, the higher frequencies of the FFR may reflect a greater contribution
278 from more peripheral areas such as the inferior colliculus, as generally the upper limit of phase-
279 locking to sound is lower in more central structures (Joris et al., 2004). Our finding of a
280 relationship between non-native speech perception ability and phase-locking within both the low-
281 frequency FFR and the gamma band, therefore, may indicate that learning a second language in
282 adulthood relies upon precise cortical but not subcortical auditory processing. This hypothesis

283 cannot be properly evaluated by the current study; however, it could be tested by future work
284 examining FFR phase-locking and non-native speech perception using MEG.

285 Previous work (Nagarajan et al., 1999; Heim et al. 2011) has demonstrated that children with
286 language learning difficulties have less phase-locked gamma band onset responses to sounds
287 presented with a short inter-stimulus interval (ISI). Here we find that degree of gamma phase-
288 locking is linked to non-native speech perception. Given that our stimuli were presented with a
289 short ISI, this could reflect an impaired ability to process rapidly presented sounds on the part of
290 the participants who struggled to learn to perceive English. Future work could examine this
291 hypothesis by examining links between non-native speech perception and gamma phase-locking
292 to stimuli presented at different ISIs. This enhanced gamma phase-locking in participants better
293 able to perceive English may also reflect greater recruitment of speech processing resources in
294 response to synthesized English speech sounds in these participants, as gamma phase-locking has
295 been shown to be greater for speech stimuli as compared to non-speech stimuli (Palva et al.,
296 2002). This would be consistent with fMRI evidence showing that subjects who are better at
297 learning novel speech sounds display more STG activity when passively listening to speech sounds
298 (Archila-Suerte et al., 2016). Finally, gamma phase-locking has also been hypothesized to be an
299 important component of speech perception in multi-time resolution models (Poeppel et al.,
300 2008), in which phonetic information is carried within the gamma band and prosodic information
301 is carried within the delta and theta bands. Greater gamma phase-locking in the participants who
302 were better able to perceive English speech may, therefore, indicate more precise neural
303 encoding of the timing of the speech envelope. This interpretation is supported by our finding
304 that gamma phase-locking was correlated with FFR phase-locking.

305 One limitation of this work is that it is difficult to rule out the possibility that the link between
306 neural sound encoding and non-native speech perceptual ability is driven by experiential factors.
307 Time spent in the United Kingdom, for example, was linked to both F0 phase-locking and English
308 perception, a relationship which is likely contributing to the link between F0 phase-locking and
309 speech perception performance. However, the relationship between neural sound encoding and
310 non-native speech perception held even after time in the UK and age of arrival were controlled
311 for, suggesting that this relationship partially reflects the dependence of successful non-native
312 language learning on auditory skills. Moreover, the relationship between non-native speech
313 perception and F0 phase-locking held both for the neural response to the consonant, which did
314 not overlap with any Japanese speech sound category, and the response to the vowel, which
315 contained formant frequencies similar to those of the Japanese [a] (Nishi et al. 2008).
316 Nevertheless, in a retrospective study it is difficult to account for all possible confounding
317 experiential factors. This limitation could be addressed in future work in which participants are
318 tested prior to beginning study of a foreign language for the first time or through the use of very
319 short-term training paradigms (Lim and Holt, 2011).

320

321

322 **References**
323

324 Archila-Suerte P, Bunta F, Hernandez A (2016) Speech sound learning depends on individuals'
325 ability, not just experience. International Journal of Bilingualism 20: 231-253.

326 Bidelman G, Krishnan A (2009) Neural correlates of consonance, dissonance, and the hierarchy of
327 musical pitch in the human brainstem. Journal of Neuroscience 29: 13165-13171.

328 Boersma P, Weenink D (2016) Praat: doing phonetics by computer [Computer program]. Version
329 6.0.22, retrieved from http://www.praat.org/

330 Bradlow A, Pisoni D, Akahane-Yamada R, Tohkura Y (1996) Training Japanese listeners to identify
331 English /r/ and /l/: IV. Some effects of perceptual learning on speech production. JASA 101: 2299-
332 2310.

333 Bradlow A, Akahane-Yamada R, Pisoni D, Tohkura Y (1999) Training Japanese listeners to identify
334 English /r/ and /l/: Long-term retention of learning in perception and production. Perception and
335 Psychophysics 61: 977-985.

336 Chandrasekaran B, Kraus N, Wong P (2012) Human inferior colliculus activity relates to individual
337 differences in spoken language learning. J Neurophysiol 107: 1325-1336.

338 Coffey E, Herholz S, Chepesiuk A, Baillet S, Zatorre R (2016) Cortical contributions to the auditory
339 frequency-following response revealed by MEG. Nature Communications 7: 11070.

340 Díaz B, Baus C, Escera C, Costa A, Sebastián-Gallés N (2008) Brain potentials to native phoneme
341 discrimination reveal the origin of individual differences in learning the sounds of a second
342 language. PNAS 105: 16083-16088.

343 Díaz B, Mitterer H, Broersma M, Escera C, Sebastián-Gallés N (2015) Variability in L2 phonemic
344 learning originates from speech-specific capabilities: an MMN study on late bilinguals.
345 Bilingualism: Language and Cognition 19: 955-970.

346 Flege J, Takagi N, Mann V (1995) Lexical familiarity and English-language experience affect
347 Japanese adults' perception of /r/ and /l/. JASA 99: 1161-1173.

348 Garcia-Sierra A, Rivera-Gaxiola M, Percaccio C, Conboy B, Romo H, Klarman L, Ortiz S, Kuhl P
349 (2011) Bilingual language learning: an ERP study relating early brain responses to speech,
350 language input, and later word production. Journal of Phonetics 39: 546-557.

351 Golestani N, Zatorre R (2009) Individual differences in the acquisition of second language
352 phonology. Brain and Language 109: 55-67.

353 Goto H (1971) Auditory perception by normal Japanese adults of the sounds "L" and "R".
354 Neuropsychologia 9: 317-323.

355 Hanulíková A, Dediu D, Fang Z, Basnaková J, Huettig F (2012) Individual differences in the
356 acquisition of a complex L2 phonology: a training study. Language Learning 62: 79-109.

357 Heim S, Friedman J, Keil A, Benasich A (2011) Reduced oscillatory activity during rapid auditory
358 processing as a correlate of language-learning impairment. Journal of Neurolinguistics 24: 538-
359 555.

360 Hornickel J, Kraus N (2013) Unstable representation of sound: a biological marker of dyslexia.
361 Journal of Neuroscience 33: 3500-3504.

362 Ingvalson E, McClelland J, Holt L (2011) Predicting native English-like performance by native
363 Japanese speakers. J Phon 39: 571-584.

364 Iverson P, Hazan V, Bannister K (2005) Phonetic training with acoustic cue manipulations: a
365 comparison of methods for teaching English /r/-/l/ to Japanese adults. JASA 118: 3267-3278.

366 Jakoby H, Goldstein A, Faust M (2011) Electrophysiological correlates of speech perception
367 mechanisms and individual differences in second language attainment. Psychophysiology 48:
368 1517-1531.

369 Joris P, Schreiner C, Rees A (2004) Neural processing of modulated sounds. Physiol Rev 84: 541-
370 577.

371 Kempe V, Thoresen J, Kirk N, Schaeffler F, Brooks P (2012) Individual differences in the
372 discrimination of novel speech sounds: effects of sex, temporal processing, musical and cognitive
373 abilities. PLoS ONE 7: e48623.

374 Kempe V, Bublitz D, Brooks P (2015) Musical ability and non-native speech-sound processing are
375 linked through sensitivity to pitch and spectral information. British Journal of Psychology 106:
376 349-366.

377 Krizman J, Marian V, Shook A, Skoe E, Kraus N (2012) Subcortical encoding of sound is enhanced in
378 bilinguals and relates to executive function advantages. Proceedings of the National Academy of
379 Sciences 109: 7877-7881.

380 Krizman J, Slater J, Skoe E, Marian V, Kraus N (2015) Neural processing of speech in children is
381 influenced by extent of bilingual experience. Neuroscience Letters 585: 48-53.

382 Lengeris A, Hazan V (2010) The effect of native vowel processing ability and frequency
383 discrimination acuity on the phonetic training of English vowels for native speakers of Greek. JASA
384 128: 3757-3568.

385 Lim S, Holt L (2011) Learning foreign sounds in an alien world: videogame training improves non-
386 native speech categorization. Cognitive Science 35: 1390-1405.

387 Lively S, Logan J, Pisoni D (1993) Training Japanese listeners to identify English /r/ and /l/. II: The
388 role of phonetic environment and talker variability in learning new perceptual categories. JASA
389 94: 1242-1255.

390 Lively S, Pisoni D, Yamada R, Tohkura Y, Yamada T (1994). Training Japanese listeners to identify
391 English /r/ and /l/. III. Long-term retention of new phonetic categories. JASA 96: 2076-2087.

392  Logan J, Lively S, Pisoni D (1990) Training Japanese listeners to identify English /r/ and /l/: a first
393  report. JASA 89: 875-886.

394  MacKain K, Best C, Strange W (1981) Categorical perception of English /r/ and /l/ by Japanese
395  bilinguals. Applied Psycholinguistics 2: 369-390.

396  Marian V, Blumenfeld H, Kaushanskaya M (2007) The language experience and proficiency
397  questionnaire (LEAP-Q): assessing language profiles in bilinguals and multilinguals. Journal of
398  Speech, Language, and Hearing Research 50: 940-967.

399  McCandliss B, Fiez J, Protopapas A, Conway M, McClelland J (2002) Success and failure in teaching
400  the [r]-[r] contrast to Japanese adults: tests of a Hebbian model of plasticity and stabilization in
401  spoken language perception. Cognitive, Affective, and Behavioral Neuroscience 2: 89-108.

402  Miyawaki K, Strange W, Verbrugge R, Liberman A, Jenkins J (1975) An effect of linguistic
403  experience: the discrimination of [r] and [l] by native speakers of Japanese and English.
404  Perception and Psychophysics 18: 331-340.

405  Nagarajan S, Mahncke H, Salz T, Tallal P, Roberts T, Merzenich M (1999) Cortical auditory signal
406  processing in poor readers. PNAS 96: 6483-6488.

407  Nishi K, Strange W, Akahane-Yamada R, Kubo R, Trent-Brown S (2008) Acoustic and perceptual
408  similarity of Japanese and American English vowels. JASA 124: 576-588.

409  Palva S, Palva J, Shtyrov Y, Kujala T, Ilmoniemi R, Kaila K, Naatanen R (2002) Distinct gamma-band
410  evoked responses to speech and non-speech sounds in humans. Journal of Neuroscience 22: 1-5.

411  Parbery-Clark A, Skoe E, Kraus N (2009) Musical experience limits the degradative effects of
412  background noise on the neural processing of sound. Journal of Neuroscience 29: 14100-14107.

413  Perrachione T, Lee J, Ha L, Wong P (2011) Learning a novel phonological contrast depends on
414  interactions between individual differences and training paradigm design. JASA 130: 461-472.

415  Poeppel D, Idsardi W, van Wassenhove V (2008) Speech perception at the interface of
416  neurobiology and linguistics. Philosophical Transactions of the Royal Society B 363: 1071-1086.

417  Ramus F, Rosen S, Dakin S, Day B, Castellote J, White S, Frith U (2003) Theories of developmental
418  dyslexia: insights from a multiple case study of dyslexic adults. Brain 126: 841-865.

419  Repp B (2000) Compensation for subliminal timing perturbations in perceptual-motor
420  synchronization. Psychol Res 63: 106-128.

421  Rosen S (1992) Temporal information in speech: acoustic, auditory and linguistic aspects.
422  Philosophical Transactions: Biological Sciences 336: 367-373.

423  Shannon R, Zeng F, Kamath V, Wygonski J, Ekelid M (1995) Speech recognition with primarily
424  temporal cues. Science 270: 303-304.

425  Skoe E, Kraus N (2010) Auditory brainstem response to complex sounds: a tutorial. Ear and
426  Hearing 31: 302-324.

427  Slevc R, Miyake A (2006) Individual differences in second-language proficiency: does musical
428  ability matter? Psychological Science 17: 675-681.

429  Steiger J (1980) Tests for comparing elements of a correlation matrix. Psychological Bulletin 87:
430  245-251.

431  Tierney A, Kraus N (2013) The ability to move to a beat is linked to the consistency of neural
432  responses to sound. Journal of Neuroscience 33: 14981-14988.

433  Tierney A, Kraus N (2016) Getting back on the beat: links between auditory-motor integration and
434  precise auditory processing at fast time scales. European Journal of Neuroscience 43: 782-791.

435  White-Schwoch T, Woodruff Carr K, Thompson EC, Anderson S, Nicol T, Bradlow AR, Zecker SG,
436  Kraus N (2015) Auditory processing in noise: a preschool biomarker for literacy. PLoS Biology 13:
437  e1002196.

438  Wong P, Perrachione T (2007) Learning pitch patterns in lexical identification by native English-
439  speaking adults. Applied Psycholinguistics 28: 565-585.

440  Woodruff Carr K, White-Schwoch T, Tierney A, Strait DL, Kraus N (2014) Beat synchronization
441  predicts neural speech encoding and reading readiness in preschoolers. Proceedings of the
442  National Academy of Sciences 111: 14559-14564.

443  Zhang Y, Kuhl P, Imada T, Iverson P, Pruitt J, Stevens E, Kawakatsu M, Tohkura Y, Nemoto I (2009)
444  Neural signatures of phonetic learning in adulthood: a magnetoencephalography study.
445  NeuroImage 46: 226-240.

446  Zhu L, Bharadwaj H, Xia J, Shinn-Cunningham B (2013) A comparison of spectral magnitude and
447  phase-locking value analyses of the frequency-following response to complex tones. JASA 134:
448  384-395.

449

450

451

452

453

454

455

456

| Speech sound contrast | Number of items |
|---|---|
| consonants | 38 |
| b-v | 4 |
| f-h | 6 |
| l-r | 14 |
| n-ŋ | 3 |
| s-ʃ | 3 |
| s-θ | 8 |
| vowels | 32 |
| æ-ɛ | 4 |
| æ-ʌ | 6 |
| ɑː-ʌ | 1 |
| ɒ-ʊ | 1 |
| ɒ-ʌ | 2 |
| ʊ-ɔː | 5 |
| ɜː-ɑː | 5 |
| iː-ɪ | 4 |
| ɪ-ɛ | 4 |

457

458    **Table 1.** Speech sound contrasts included in the receptive phonology test.

459

460

461

462

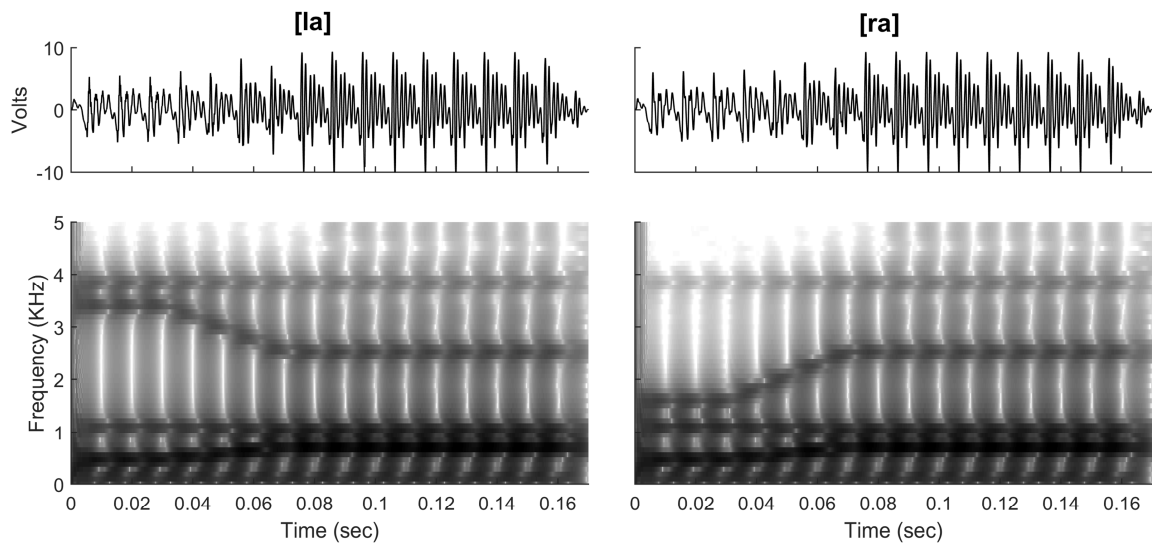463

464

465

466

467

468

469

470

471     **Figure 1**. Waveforms (top) and spectrograms (bottom) of synthesized speech stimuli. The [la] and
472     [ra] stimuli differed only in the first 70 ms, and were identical thereafter.

473     **Figure 2**. (Left) Time-frequency plot of inter-trial phase locking across all subjects for the
474     frequency following response (71-600 Hz). (Right) Time-frequency plot of inter-trial phase locking
475     across all subjects for the cortical response (8-70 Hz).

476     **Figure 3**. (Left, top) Time-frequency plot of inter-trial phase locking for the frequency following
477     response for participants with good versus poor perception of English consonants. Participants
478     were divided into top and bottom halves based on performance on the consonant portions of the
479     receptive phonology test. (Right, top) Time-frequency plot of inter-trial phase locking for the
480     cortical response for good versus poor consonant perceivers. (Left, bottom) Inter-trial phase
481     locking in the frequency following response as a function of frequency across the entire response
482     (10-170 ms) for good (red) versus poor (blue) consonant perceivers. Error bars are one standard
483     error of the mean. (Right, bottom) Inter-trial phase locking in the frequency following response as
484     a function of frequency across the first 60 ms of the response for good versus poor consonant
485     perceivers.

486     **Figure 4**. (Left) Scatterplot displaying performance on the consonant portions of the receptive
487     phonology test (displayed as portion correct) versus inter-trial phaselocking at the fundamental
488     frequency during the entirety of the frequency following response. (Right) Scatterplot displaying
489     consonant perception versus inter-trial phaselocking within the gamma band (31-70 Hz) during
490     the first 60 ms of the cortical response. R-values and p-values are derived from Pearson
491     correlations.

492

**Frequency following response**

**Cortical and subcortical responses (8-70 Hz)**

Frequency following response

Good consonant perception

Poor consonant perception

Cortical response

Good consonant perception

Poor consonant perception

Inter-trial phase locking