

1 Relationships between native and non-native speech perception

2 Pamela Fuhrmeister<sup>1</sup>, Matthew C. Phillips<sup>1</sup>, D. Betsy McCoach<sup>2</sup>, & Emily B. Myers<sup>1</sup>

3 <sup>1</sup> University of Connecticut, Department of Speech, Language, and Hearing Sciences

4 <sup>2</sup> University of Connecticut, Neag School of Education

5 Author Note

6 This material is based upon work supported in part by the National Science  
7 Foundation under grant DGE-1747486 to the University of Connecticut and NSF BCS  
8 1554510 to EBM. Any opinions, findings, and conclusions or recommendations expressed in  
9 this material are those of the author(s) and do not necessarily reflect the views of the  
10 National Science Foundation. The authors are grateful to Hannah Mechtenberg for help  
11 with data collection. The data and analysis code for the experiments in this paper can be  
12 found at <https://osf.io/2zg6c/>.

13 Correspondence concerning this article should be addressed to Pamela Fuhrmeister.

14 E-mail: [pamela.fuhrmeister@uconn.edu](mailto:pamela.fuhrmeister@uconn.edu)

## Abstract

15

16 Individuals differ in their ability to perceive and learn unfamiliar speech sounds, but we  
17 lack a comprehensive theoretical account that predicts individual differences in this skill.  
18 Predominant theories largely attribute difficulties of non-native speech perception to the  
19 relationships between non-native speech sounds/contrasts and native-language categories.  
20 The goal of the current study was to test whether the predictions made by these theories  
21 can be extended to predict individual differences in naive perception of non-native speech  
22 sounds or learning of these sounds. Specifically, we hypothesized that the internal structure  
23 of native-language speech categories is the cause of difficulty in perception of unfamiliar  
24 sounds such that learners who show more graded (i.e., less categorical) perception of  
25 sounds in their native language would have an advantage for perceiving non-native speech  
26 sounds because they would be less likely to assimilate unfamiliar speech tokens to their  
27 native-language categories. We tested this prediction in two experiments in which listeners  
28 categorized speech continua in their native language and performed tasks of discrimination  
29 or identification of difficult non-native speech sound contrasts. Overall, results did not  
30 support the hypothesis that individual differences in categorical perception of  
31 *native-language* speech sounds is responsible for variability in sensitivity to *non-native*  
32 speech sounds. However, participants who responded more consistently on a speech  
33 categorization task showed more accurate perception of non-native speech sounds. This  
34 suggests that individual differences in non-native speech perception are more related to the  
35 stability of phonetic processing abilities than to individual differences in phonetic category  
36 structure.

37

*Keywords:* speech perception, categorical perception, individual differences,

38

non-native speech sound learning, perceptual assimilation

## Relationships between native and non-native speech perception

39

40 Many adults struggle to perceive and produce speech sounds in a second language  
41 (e.g., Bradlow, Akahane-Yamada, Pisoni, & Tohkura, 1999). Some obvious factors  
42 contribute to more successful learning of speech sounds, including the age at which  
43 someone starts learning the language or the amount of time they spend using their second  
44 language (e.g., Flege, Yeni-Komshian, & Liu, 1999; Piske, MacKay, & Flege, 2001).  
45 However, even among learners with similar experience with a second language, we still see  
46 a large range of individual variability in perception and production of non-native speech  
47 sounds: Some learners meet with excellent success in distinguishing difficult non-native  
48 contrasts, where others never manage to grasp the distinctions (e.g., Bradlow et al., 1999;  
49 Golestani & Zatorre, 2009; Lim & Holt, 2011; Luthra et al., 2019). What, then, accounts  
50 for this variability? A collection of findings suggests that perception or production of  
51 non-native speech sounds may be related to an individual's cue-weighting strategies or  
52 abilities (Schertz, Cho, Lotto, & Warner, 2015), auditory processing skills (Kachlicka,  
53 Saito, & Tierney, 2019), musical training/pitch perception (especially for learning of tonal  
54 contrasts, Bowles, Chang, & Karuzis, 2016; Kempe, Bublitz, & Brooks, 2015; Perrachione,  
55 Lee, Ha, & Wong, 2011; Slevc & Miyake, 2006), phonological skills (Earle & Arthur, 2017;  
56 Fuhrmeister, Schlemmer, & Myers, 2020; Perrachione et al., 2011), or brain structure  
57 (Golestani, Molko, Dehaene, LeBihan, & Pallier, 2007; Golestani, Paus, & Zatorre, 2002).

58 In the current study, we test a prediction that falls out of dominant theories of  
59 non-native speech perception: namely that the challenge of non-native speech perception is  
60 the result of the difficulties that adults have in detecting differences between items within  
61 native language speech categories, a phenomenon known as “categorical perception.” As  
62 will be described in detail in the following paragraphs, popular theories of non-native  
63 speech perception predict that the relationship between native and non-native speech  
64 sounds explains the difficulty in perceiving non-native sounds (e.g., Best, McRoberts, &

65 Goodell, 2001; Best & Tyler, 2007; Flege, 1995; Iverson et al., 2003; Kuhl, 1994; Kuhl et  
66 al., 2008). These theories reliably predict *which* speech sounds or speech sound contrasts  
67 will be relatively easier or relatively difficult to perceive for a listener of a given  
68 native-language based on the relationship between native and non-native speech sounds  
69 (e.g., Best et al., 2001; Best, McRoberts, & Sithole, 1988). However, it is unknown whether  
70 these theories can be extended to predict which *individuals* will be more successful  
71 non-native speech sound perceivers or learners. Silbert et al. (2015) tested whether  
72 individuals have a general aptitude for perceiving non-native speech sound contrasts, and  
73 instead of a general aptitude, their results indicate that listeners tended to show different  
74 patterns of perception for segmental vs. tonal contrasts. This suggests that we may need to  
75 consider relationships between various skills and perception of segmental and  
76 suprasegmental contrasts separately, and hints that difficulty perceiving non-native  
77 contrasts might be tied to more specific relationships between an individual non-native  
78 contrast and corresponding native category representations. The goal of the current study  
79 is to test the hypothesis that individual differences in native-language speech category  
80 representations (specifically how categorically individuals perceive sounds) can predict  
81 success in perceiving or learning new speech categories. In the following paragraphs, we  
82 discuss categorical perception of speech, predictions from theories of non-native speech  
83 perception, and our reasoning for bringing these literatures together to generate new  
84 predictions for individual differences in non-native speech perception.

### 85 **Individual differences in categorical perception of native-language speech** 86 **sounds**

87 For some time, the dominant view in the field has been that listeners perceive speech  
88 sounds categorically, i.e., acoustic differences between speech categories are highly  
89 detectable, while differences within speech categories are often poorly detected (Liberman,  
90 Cooper, Shankweiler, & Studdert-Kennedy, 1967; e.g., Liberman, Harris, Hoffman, &

91 Griffith, 1957). Since this seminal finding, much research in the field has expanded on this  
92 phenomenon but has only recently begun to examine whether individuals differ  
93 systematically in whether their perception of phonetic categories is more categorical or  
94 more graded (e.g., Kapnoula, Winn, Kong, Edwards, & McMurray, 2017; Kong & Edwards,  
95 2016).

96 One potential limitation of previous studies is that two-alternative forced choice  
97 identification tasks that are often used to measure categorical perception do not allow  
98 listeners to sufficiently demonstrate gradedness in their perception of speech sounds. For  
99 instance, when only given the choice between /d/ and /t/, many listeners who could  
100 otherwise detect subtle acoustic differences between two exemplars in the /d/ category  
101 would still identify all or most /d/ tokens as belonging to the /d/ category because those  
102 tokens would still be better exemplars of /d/ than /t/. Other tasks for measuring  
103 categorical perception offer the listener more opportunities to demonstrate graded  
104 perception of speech sounds. Some of these methods include eye-tracking (e.g., Clayards,  
105 Tanenhaus, Aslin, & Jacobs, 2008; McMurray, Danelz, Rigler, & Seedorff, 2018; McMurray,  
106 Tanenhaus, & Aslin, 2002), goodness judgments (Drouin, Theodore, & Myers, 2016; Miller,  
107 1997), or visual analogue scales, in which a listener moves a visual slider on a screen  
108 between two alternatives (Kapnoula, Edwards, & McMurray, 2021; Kapnoula et al., 2017;  
109 Kong & Edwards, 2016). These studies suggest that when the task affords an opportunity  
110 to demonstrate sensitivity to variation within the phonetic category, listeners show  
111 sensitivity to the graded internal structure of speech categories (Clayards et al., 2008;  
112 Drouin et al., 2016; McMurray et al., 2018, 2002; Miller, 1997). Furthermore, considerable  
113 variability exists even among typically developing individuals in how graded or  
114 categorically speech categories are represented when measured in this way (Kapnoula et  
115 al., 2021, 2017; Kong & Edwards, 2016). Taken together, these findings suggest that while  
116 most listeners do continue to show a classic “categorical perception” pattern, with lessened  
117 sensitivity to sound contrasts falling within a category, these listeners differ systematically

118 in the degree to which they show a more categorical vs. graded perceptual pattern. Far  
119 from having uniform perceptual abilities, adults differ in the degree to which they can  
120 detect fine-grained phonetic differences within a native-language speech category. Given  
121 that theoretical accounts of non-native speech perception (discussed below) hinge on  
122 assimilation within the native category, this makes *native-language* categoricity/gradience  
123 an attractive explanation for individual differences in the difficulty of *non-native* speech  
124 perception.

### 125 **Theoretical predictions from non-native speech sound learning**

126 Several theories attribute difficulties in non-native speech sound learning to  
127 perceptual similarity of non-native contrasts to native language speech categories (e.g.,  
128 Best & Tyler, 2007; Flege, 1995; Iverson et al., 2003; Kuhl et al., 2008). For example, the  
129 native-language magnet model discusses this in terms of category prototypes:  
130 native-language categories will act as a “magnet” so that unfamiliar speech sounds that are  
131 acoustically or perceptually similar to native-language sounds will often be perceived as  
132 exemplars of the native-language category that is more familiar to the listener (Kuhl, 1994;  
133 Kuhl et al., 2008). Similarly, Iverson et al. (2003) provide a perceptual interference  
134 account of non-native speech sound learning. They suggest that experience with a (native)  
135 language changes perception and perhaps also lower-level auditory processing for certain  
136 acoustic dimensions, which in turn interferes with learning new sounds that are contrasted  
137 by that dimension.

138 The perceptual assimilation model predicts that non-native speech contrasts will be  
139 perceptually assimilated to native-language sounds, and these assimilation patterns  
140 determine how difficult it will be for a listener to differentiate these non-native speech  
141 sounds (Best et al., 2001; Best & Tyler, 2007). Specifically, the distinctions between two  
142 non-native speech sounds that assimilate to one native-language speech category will be  
143 particularly difficult to perceive, whereas contrasts that assimilate to two different native

144 speech categories will be easier (Best et al., 2001). For example, the voiced dental and  
145 retroflex stop consonants in Hindi are often difficult for native speakers of many varieties of  
146 English to differentiate because they are allophones of the alveolar /d/ category (e.g., in  
147 “width” or “address,” Polka, 1991). These sounds then get assimilated to the alveolar  
148 category, making it difficult to perceive them as separate sounds. Mayr and Escudero  
149 (2010) found tentative support for the idea that individual differences in assimilation  
150 patterns between native and non-native speech sounds predicted more accurate perception  
151 of non-native contrasts; however, Hattori and Iverson (2009) tested a similar question but  
152 did not find much support for the notion that the degree of assimilation between a difficult  
153 non-native contrast and the perceptually similar native-language category was related to  
154 perception or production accuracy of the non-native contrast.

155 The speech learning model also assumes that the (dis)similarity between phonemes in  
156 the first and second language will influence the likelihood of establishing new phonological  
157 categories in a second language (Flege, 1995). This model focuses on experienced  
158 second-language learners rather than naive listeners and therefore a longer time span of  
159 learning; however, many of the core ideas and predictions are similar to the other models  
160 discussed here.

161 Although these theories are highly accurate in predicting *which* non-native speech  
162 sounds will be most difficult to acquire for learners of a given language background, they  
163 do not explain the variability commonly observed among learners with the *same* native  
164 language. Notably, these models do not make explicit predictions about individual  
165 differences in non-native speech perception, but implicit in the models’ assumptions is that  
166 variation in the perception of native-language speech category structure should predict how  
167 easy or difficult acquiring perceptual sensitivity to a new, non-native speech sound contrast  
168 will be for an individual. Therefore, these theories may be able to generate predictions on  
169 an individual level: All of them would likely predict a relationship between subtle variation

170 in how native-language categories are perceived and how well someone can perceive or  
171 learn non-native speech sounds. For instance, if an individual shows perception of  
172 native-language categories that is more graded and less categorical (i.e., that person can  
173 distinguish subtle within-category differences), they may be less likely to assimilate  
174 similar-sounding non-native speech sounds to existing native categories. In other words,  
175 those with more graded native-language categories may be better equipped to detect  
176 acoustic variability that defines non-native contrasts.

### 177 **Current study**

178       The primary aim of the current study is to test whether individual differences in  
179 categorical perception of native-language speech sounds predict success on non-native  
180 speech sound learning tasks. We used a visual analogue scale task similar to those used in  
181 Kapnoula et al. (2017) and Kong and Edwards (2016) to derive a continuous measure of  
182 how categorically or graded participants perceive native-language speech sounds, which we  
183 refer to as *categoricity*. A finding in which categoricity negatively predicts non-native  
184 speech sound learning outcomes (i.e., participants who are less categorical are more  
185 successful on non-native tasks) would lend support to the idea that theories of non-native  
186 speech sound learning, such as the native-language magnet model (Kuhl, 1994; Kuhl et al.,  
187 2008), perceptual interference model (Iverson et al., 2003), perceptual assimilation model  
188 (Best et al., 2001; Best & Tyler, 2007), or speech learning model (Flege, 1995) can be  
189 extended to account for individual differences in non-native speech sound learning. In the  
190 current study we used a discrete visual analogue scale that had seven discrete points on the  
191 line (one for each continuum point). This allowed us to measure not only categoricity, but  
192 also how consistently participants rated the stimuli each time they were presented (i.e.,  
193 how often they rated the first step on the continuum as “1” when it was presented). Some  
194 recent evidence suggests that perception of native-language speech sounds becomes more  
195 graded throughout adolescence (McMurray et al., 2018) and that adult speech category

196 representations are more graded than might be reflected by a two-alternative forced choice  
197 task McMurray et al. (2002). McMurray et al. (2018) argue that shallower categorization  
198 slopes in two-alternative forced choice phonetic identification tasks are a result of noisy  
199 representations, rather than the older view that they are due to more graded  
200 representations (e.g., Burnham, Earnshaw, & Clark, 1991; Hazan & Barrett, 2000).  
201 Therefore, we also collected a measure of how precise (or noisy) a participant's responses  
202 on the task were (*response consistency*). If more precise representations of speech sounds  
203 allow for the integration and formation of new speech categories, we would expect to see  
204 that more consistent responses on this task predict non-native speech sound learning.

## 205 Experiment 1

206 In Experiment 1, we tested whether categoricity and response consistency of two  
207 different native-language speech contrasts (a stop and a fricative contrast) predicted naive  
208 discrimination of a difficult non-native speech contrast or learning/retention of this  
209 contrast after a brief training period.

## 210 Method

### 211 Participants

212 Fifty-eight participants were recruited from the University of Connecticut  
213 (Connecticut, United States) community. All participants were native speakers of North  
214 American English and had typical hearing and no reading or language disorders. One  
215 participant could not complete the experiment due to an equipment error; data from the  
216 remaining 57 participants are included in the analyses. All participants gave informed  
217 consent according to the University of Connecticut Institutional Review Board  
218 requirements and were paid \$10 per hour for participation.

## 219 Stimuli

220 We tested participants' discrimination and identification of a difficult non-native  
221 speech sound contrast for our participants (native speakers of North American English),  
222 the voiced dental and retroflex stop consonants found in Hindi. We recorded five exemplars  
223 each of the minimal pair non-words /d̪ug/ and /ɖug/. Stimuli were recorded by a female  
224 native speaker of Hindi; stimuli were recorded in a soundproof booth and were scaled to a  
225 mean amplitude of 65 dB sound pressure level in Praat (Boersma & Weenink, 2013).  
226 Fribbles served as novel objects to pair new words with (stimulus images courtesy of  
227 Michael J. Tarr, Center for the Neural Basis of Cognition and Department of Psychology,  
228 Carnegie Mellon University, <http://www.tarrlab.org/>). Non-native tasks were presented on  
229 a desktop computer with the experiment presentation software OpenSesame (Mathôt,  
230 Schreij, & Theeuwes, 2012).

231 To measure perception of native-language speech categories, participants rated tokens  
232 from two seven-step continua on the visual analogue scale. One continuum consisted of a  
233 fricative contrast embedded in real words (sign-shine) and one consisted of a synthetic stop  
234 contrast of consonant-vowel syllables (/bɑ/-/dɑ/)<sup>1</sup>. The /bɑ/-/dɑ/ continuum was created  
235 using a Klatt synthesizer at Haskins Laboratories. Stimuli for the sign-shine continuum  
236 were recorded by a female, native speaker of English, and waveform averaging in Praat  
237 (Boersma & Weenink, 2013) was used to create blends from 20% /s/ to 80% /ʃ/ in 10%  
238 steps. Native-language visual analogue scale tasks were presented using E-Prime 3.0  
239 (Psychology Software Tools, Pittsburgh, PA). We chose these two contrasts because we did  
240 not want the acoustic properties of the stimuli to be too different between the native and  
241 non-native speech tasks. For instance, dental and retroflex stop consonants are

---

<sup>1</sup> We did not have strong theoretical reasons for including contrasts embedded in both real words and syllables other than that we know that these are well validated continua that elicit typical categorical patterns of perception.

242 differentiated by the frequency of the burst and the third and fourth formant trajectories,  
243 i.e., primarily spectral cues (Stevens & Blumstein, 1975). We therefore chose two  
244 native-language contrasts that are also mainly differentiated by spectral cues.

## 245 **Procedure**

246 Participants came to the lab twice. The first session took place in the evening hours  
247 (between 5 and 9 PM) and included consent and non-native phonetic training and  
248 assessments in this order: discrimination pretest<sup>2</sup>, identification training, identification  
249 posttest, and discrimination posttest. The second session took place the following morning  
250 between 8 and 10 AM and included (in this order) a test of retention of non-native speech  
251 sounds (identification posttest followed by the discrimination posttest) and after that, the  
252 visual analogue scale tasks to measure categorical perception of native-language speech  
253 sounds (the order of the continua was counterbalanced)<sup>3</sup>. The training and testing schedule  
254 was motivated by previous work showing that relationships between phonological skills in  
255 the native-language and a non-native speech sound learning task emerged after a period of  
256 sleep (Earle & Arthur, 2017). We therefore wanted to test the possibility that relationships  
257 between categoricity or consistency and non-native speech perception tasks would emerge  
258 only after a period of offline consolidation.

## 259 **Non-native speech sound learning tasks.**

260 ***Discrimination.*** The discrimination task was an AX task (i.e., same/different  
261 judgment). On each trial (64 trials total), participants heard two of the minimal pair  
262 non-words with an inter-stimulus interval of 1 s and were asked to indicate whether the  
263 non-words started with the same or different sounds. Four practice trials with feedback

---

<sup>2</sup> We chose to only administer a discrimination pretest and not an identification pretest. In similar study designs, we have never done a pretest for an identification task because that requires at least some amount of familiarization with the stimuli. A discrimination task, on the other hand, can be done naively, as participants simply indicate whether they think the sounds they heard were the same or different.

<sup>3</sup> We additionally collected a battery of tests of cognitive skills; these will be reported elsewhere.

264 were administered, and these are not included in the analyses. No feedback was given on  
 265 any subsequent trials.

266 **Identification.** First, participants were familiarized with the non-word-object  
 267 pairings by seeing one visual stimulus (Fribble) on the screen at a time while hearing all  
 268 five auditory exemplars of that non-word. Next, participants learned the sounds via a  
 269 two-alternative forced choice identification task. For this task, they saw both objects on  
 270 the screen and heard one auditory token and then were asked to choose the corresponding  
 271 picture. Visual feedback (“Correct”/“Incorrect”) was given on each trial. Participants  
 272 completed a total of 400 trials with a two-minute break halfway through. After training,  
 273 participants completed an identification assessment, which included 50 trials just like in  
 274 training but with no feedback.

275 **Native-language speech perception measures.** Participants completed two  
 276 tasks measuring categoricity of native-language speech sounds. In each task, auditory  
 277 tokens taken from one of two continua were presented (sign-shine or ba-da), and the order  
 278 of the two continua was counterbalanced. Participants indicated their response on a visual  
 279 analogue scale. In this task, participants moved a slider to discrete (numbered) points on  
 280 the line between two stimuli (in this case, words/syllables) to indicate how, for example,  
 281 sign-like or shine-like each stimulus token sounded (see Figure 1 for an illustration of the  
 282 task and Figure 2 for examples of a graded and categorical participant).

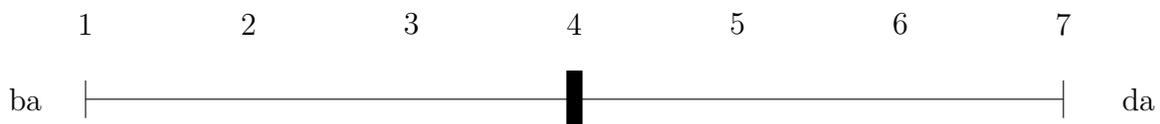


Figure 1. Illustration of the discrete visual analogue scale task.

### 283 Analysis approach

284 All analyses were performed in R (R Core Team, 2021). Categoricity and non-native  
 285 discrimination data were analyzed with generalized linear mixed effects models with the  
 286 lme4 package (Bates, Mächler, Bolker, & Walker, 2015). Unless otherwise noted, the full

287 random effects structure was used unless the model did not converge, in which case,  
288 random effects were removed one by one until the model converged without warnings  
289 (Barr, Levy, Scheepers, & Tily, 2013). All data and code to reproduce the analyses  
290 reported in this paper can be found at our OSF repository: <https://osf.io/2zg6c/>.

291 **Non-native discrimination.** A common way to analyze data from a  
292 discrimination task such as the one used here is to compute a d prime ( $d'$ ) score to account  
293 for response bias (Macmillan & Creelman, 2004). For instance, when indicating whether  
294 stimuli sound the same or different, a participant who responds “different” on every trial in  
295 the experiment would score 100% for all the different trials. However, this does not mean  
296 that the participant can actually distinguish the sounds because they also answered  
297 incorrectly for all the same trials. Using  $d'$  scores can overcome this bias by subtracting the  
298 standardized “false alarm” score (“same” trials on which the participant answered  
299 “different”) from the standardized “hit” score (“different” trials on which the participant  
300 correctly answered “different”), i.e.,  $z(\text{hits}) - z(\text{false alarms})$ . In short, this score represents  
301 the distance between the standardized distributions of responses when the stimuli are  
302 different and when they are the same. However, this approach requires averaging over  
303 trials, and we would rather incorporate this approach into a hierarchical model to take trial  
304 and item variability into account while still accounting for response bias. This can be done  
305 with a generalized linear mixed effects model with the probit link because the probit link  
306 transformation is the inverse of the cumulative distribution function of the standard  
307 normal distribution (Razzaghi, 2013; Wright, Horry, & Skagerberg, 2009).

308 To test whether categoricity or response consistency predicted accuracy on the  
309 non-native discrimination task, we fit a generalized linear mixed effects model with the  
310 probit link. The model predicted whether the response on each trial was different (yes = 1,  
311 no = 0). Fixed effects included condition (deviation coded: different = .5, same = -.5), the  
312 interaction of condition and time (backwards difference coded to test contrasts for learning:  
313 immediate posttest-pretest, and retention: next-day posttest-immediate posttest), and the

314 interaction of condition and each predictor variable (ba-da categoricity, ba-da consistency,  
315 s-sh categoricity, and s-sh consistency, all mean centered) was nested within time. Nesting  
316 the predictors within the factor of time tests for simple effects of the predictors at each  
317 time point separately (Schad, Vasishth, Hohenstein, & Kliegl, 2020, e.g., whether  
318 categoricity predicted discrimination accuracy at the pretest, immediate posttest, or  
319 next-day posttest separately). In this model, the intercept represents overall bias, or  
320 whether there was more of a tendency to answer “same” or “different” regardless of  
321 condition. The coefficient for condition represents  $d'$ , or sensitivity considering the actual  
322 condition (same/different). The interactions are what we are most interested in for the  
323 present question, which can be interpreted as the change in sensitivity ( $d'$ ) for each unit  
324 increase in categoricity or response consistency. We did not estimate main effects of  
325 categoricity or response consistency because they would not be informative for our question  
326 without interpreting them in the context of the interaction with condition.

327 **Non-native identification.** With a two-alternative forced choice identification  
328 task without feedback, we sometimes see that participants score well below chance  
329 performance (for example scoring 2% accuracy on the task). This makes sense if  
330 participants can indeed distinguish the sounds but they simply switched the labels. To  
331 address this, we recoded correct and incorrect responses for identification data (i.e., 0 was  
332 recoded as 1 and 1 was recoded as 0) if participants' mean accuracy was significantly below  
333 chance as measured by a binomial test (below 38% accuracy for 50 trials). No participants  
334 scored below chance, so no data was recoded in this experiment.

335 We analyzed identification performance using a generalized linear mixed effects model  
336 with a logit link that predicted accuracy (0 or 1). Fixed effects included ba-da categoricity,  
337 ba-da consistency, s-sh categoricity, and s-sh consistency (all mean centered), which were  
338 nested within time (deviation coded: immediate posttest = -.5, next-day = .5).

339 **Native-language speech perception measures.** To measure categoricity of  
340 native-language speech sounds, we fit responses from the visual analogue scale (for each

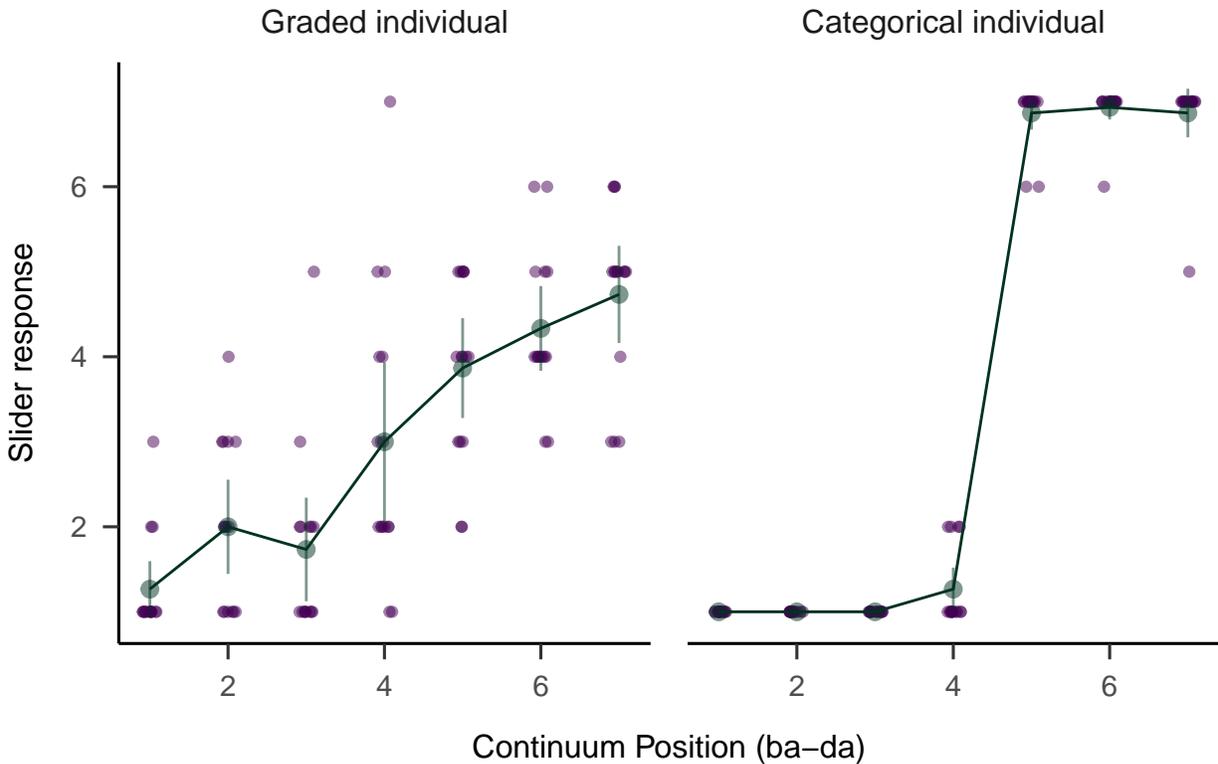


Figure 2. Examples of participants who showed more graded (left) and categorical (right) patterns of perception on the task with the visual analogue scale.

341 continuum separately) to a three-parameter logistic function with a non-linear mixed effects  
 342 regression model. The three-parameter model estimates the maximum asymptote, the  
 343 inflection point (i.e., category boundary), and the slope of the function (more categorical  
 344 response patterns have higher slope values). Non-linear regression models were fit with the  
 345 nlme package (Pinheiro, Bates, DebRoy, Sarkar, & R Core Team, 2021). Each participant's  
 346 categoricity score (the population mean plus the participant's individual random effect  
 347 adjustment) was extracted from the model. An individual measure of response consistency  
 348 on each phonetic continuum was obtained by calculating the mean of the squared residuals  
 349 from the model for each participant separately (residuals were squared to avoid negative  
 350 values). Because larger values represented less consistent responses, we changed the sign of  
 351 the response consistency measure for ease of interpretation, such that *larger* values  
 352 represent *more* consistent response patterns. The measures of slope and consistency were

353 entered into further analyses described below. We additionally computed all the possible  
 354 correlations between these measures. The correlation between the categoricity measures  
 355 from the two continua were moderately correlated,  $r = 0.31$ ,  $p = 0.02$ ; all other correlations  
 356 were non-significant: consistency for both continua,  $r = 0.16$ ,  $p = 0.22$ , ba-da categoricity  
 357 and consistency:  $r = -0.07$ ,  $p = 0.60$ , s-sh categoricity and consistency:  $r = -0.05$ ,  $p = 0.74$ .

Table 1

*Results of the discrimination task in Experiment 1. Time1 = pretest, Time2 = immediate posttest, and Time3 = next-day posttest.*

<i>Predictors</i>	<i>Estimate</i>	<i>SE</i>	<i>z</i>	<i>p</i>
Intercept	-0.07	0.05	-1.52	0.13
Condition ( $d'$ )	1.26	0.15	8.39	<.001
Time2-1	0.02	0.03	0.65	0.51
Time3-2	0.00	0.03	0.01	0.99
Time2-1 x Condition	0.64	0.07	9.84	<.001
Time3-2 x Condition	0.19	0.07	2.85	<.001
Time1:Condition x Categoricity ba-da	-0.17	0.43	-0.39	0.7
Time2:Condition x Categoricity ba-da	0.41	0.43	0.95	0.34
Time3:Condition x Categoricity ba-da	0.15	0.43	0.34	0.73
Time1:Condition x Consistency ba-da	0.35	0.35	0.99	0.32
Time2:Condition x Consistency ba-da	0.00	0.36	-0.01	1
Time3:Condition x Consistency ba-da	-0.11	0.36	-0.30	0.77
Time1:Condition x Categoricity s-sh	0.11	1.04	0.10	0.92
Time2:Condition x Categoricity s-sh	-1.28	1.05	-1.22	0.22
Time3:Condition x Categoricity s-sh	-1.23	1.05	-1.17	0.24
Time1:Condition x Consistency s-sh	0.58	0.33	1.73	0.08
Time2:Condition x Consistency s-sh	0.42	0.34	1.25	0.21

---

Time3:Condition x Consistency s-sh                      0.31   0.34   0.92   0.36

---

Table 2

*Results of the identification task in Experiment 1. Time1 = immediate posttest, and Time2 = next-day posttest.*

<i>Predictors</i>	<i>Estimate</i>	<i>SE</i>	<i>z</i>	<i>p</i>
Intercept	2.30	0.26	8.80	<.001
Time	0.13	0.12	1.10	0.27
Time1:Categoricity ba-da	1.07	0.62	1.74	0.08
Time2:Categoricity ba-da	0.17	0.67	0.25	0.81
Time1:Categoricity s-sh	-2.76	1.47	-1.88	0.06
Time2:Categoricity s-sh	-2.31	1.64	-1.41	0.16
Time1:Consistency ba-da	-0.36	0.50	-0.72	0.47
Time2:Consistency ba-da	-0.33	0.56	-0.60	0.55
Time1:Consistency s-sh	0.92	0.47	1.97	0.05
Time2:Consistency s-sh	0.73	0.52	1.39	0.16

358

## Results

### 359 Discrimination

360 Results of the discrimination analysis are summarized in Table 1. We found a  
 361 significant main effect of condition, meaning participants answered “different” more often  
 362 when the stimulus pair was different. There were significant interactions of condition with  
 363 each of the time contrasts, which indicates that participants improved in their  
 364 discrimination accuracy of the Hindi sounds at each time point. None of the measures of  
 365 categoricity or consistency significantly predicted discrimination accuracy.

## 366 Identification

367 Results of the identification analysis are summarized in Table 2. The intercept was  
368 significantly different from zero, indicating participants identified the sounds above chance,  
369 but there were no other significant effects<sup>4</sup>. The relationship between consistency on the  
370 fricative continuum and identification accuracy approached significance at the immediate  
371 posttest. This is also in the expected direction, in that more consistent responses predicted  
372 higher accuracy; however, we did not see the same pattern with the stop continuum, so we  
373 will err on the side of caution and not interpret this further.

## 374 Discussion

375 Neither categoricity nor consistency significantly predicted discrimination or  
376 identification accuracy at any of the time points we tested. One obvious explanation is that  
377 we lacked statistical power in this design. We address this issue in Experiment 2.  
378 Additionally, if independent variables in a regression model are highly correlated, we may  
379 not see any significant effects because no one variable explains unique variance. We only  
380 saw weak correlations (if any) among predictor variables of categoricity and consistency;  
381 however, to ensure that our null results were not due to issues of multicollinearity, we  
382 calculated variance inflation factors (a diagnostic tool to detect multicollinearity in  
383 multiple regression models) for all predictors in the model using the car package in R (Fox  
384 & Weisberg, 2019). Variance inflation factors of over 5 (or over 10 according to some  
385 authors) are considered problematic (e.g., Thompson, Kim, Aloe, & Becker, 2017). Most  
386 variance inflation factors were around 1 with only a couple closer to 2 (values for each

---

<sup>4</sup> We additionally tested the interaction between each predictor and time and found that the interaction between ba-da categoricity and time was significant,  $\beta = -0.91$ ,  $SE = 0.32$ ,  $z = -2.83$ ,  $p = 0.005$ , suggesting that *more* categorical perception of the ba-da continuum predicted more accurate identification at both time points, but more strongly on the first day. Because the simple effects were not statistically significant and this pattern is not consistent with any theoretical predictions, we think this is likely a Type I error and advise caution in any interpretation of this finding.

387 predictor can be found in our supplementary analysis document on OSF,  
388 <https://osf.io/2zg6c/>), indicating no issues with multicollinearity. Finally, we acknowledge  
389 that we only included a very brief training period in this experiment design. Even though  
390 we saw significant improvement on the discrimination task after training and after the  
391 overnight interval, it is possible that predictors of real-world or longer-term learning of  
392 non-native speech contrasts differ from those found for short-term, laboratory-based  
393 studies. The results presented here cannot speak to possible relationships between  
394 native-language speech perception and long-term learning of non-native speech sounds.

395 Another possibility is that the relationship between categoricity and perception of  
396 unfamiliar speech sounds might be cue-specific or it might only hold for sounds that are  
397 perceptually similar in the native and non-native languages (for more on cue-specificity, see  
398 e.g., Kapnoula et al., 2021; Saito et al., 2022). This could at least explain why we did not  
399 find a significant relationship between s-sh categoricity and perception of the Hindi  
400 contrast. Though the sounds /b/ and /d/ are more similar to the Hindi sounds, some work  
401 has found that listeners show more categorical patterns of perception for stop consonants  
402 than other types of sounds (Eimas, 1963; Healy & Repp, 1982; Repp, 1981). Therefore, we  
403 may find a wider range of variability in categoricity of the s-sh continuum, especially when  
404 testing a larger sample. This could allow us to detect a relationship between categoricity  
405 and perception of a non-native contrast that is similar to the s-sh sounds.

## 406 Experiment 2

407 The purpose of Experiment 2 is twofold: First, to address the issue of statistical  
408 power raised in Experiment 1, we simplified the design and collected a much larger sample  
409 size. Second, we tested the question of whether the relationship between categoricity and  
410 discrimination of unfamiliar speech sounds is specific to perceptually similar contrasts. In  
411 this experiment, we tested participants' categoricity with the same phonetic continua  
412 (ba-da and s-sh) as in Experiment 1. We also tested participants' discrimination of the

413 Hindi dental and retroflex stop consonants and their discrimination of an additional  
414 unfamiliar speech contrast, namely the voiceless alveolo-palatal /ç/ and retroflex /ʂ/  
415 fricatives found in Polish. Prior to explicit instruction, these speech sounds are not well  
416 discriminated by native English listeners (Lisker, 2001; McGuire, 2007; Żygis & Padgett,  
417 2010), due to their similarity to the /f/ category in English (Best et al., 2001; Kuhl et al.,  
418 2008). We therefore tested whether categoricity or consistency of native-language sounds  
419 predict discrimination of perceptually similar non-native sounds (i.e., whether perceptual  
420 measures of ba-da predicts discrimination of the Hindi sounds and s-sh predicts  
421 discrimination of the Polish sounds). We did not include a training task in this experiment,  
422 so this experiment focused on the relationship between categoricity/consistency and *naïve*  
423 perception of non-native speech sound contrasts. We note, however, that the predictions  
424 from the perceptual assimilation model, for example, apply to naïve perception (Best et al.,  
425 2001). In addition, in our previous work, we have seen that performance on a  
426 discrimination pretest typically predicts performance on the posttests after training  
427 (Fuhrmeister & Myers, 2020); therefore, it is possible that any findings from this  
428 experiment could be relevant for learning, but we did not test that explicitly here.

## 429 Method

### 430 Participants

431 A sample of 269 participants was recruited from the online data collection platform  
432 Prolific ([www.prolific.co](http://www.prolific.co)). Our goal was to have a sample size of at least 200 participants,  
433 and data collection was stopped after our sample size reached at least 200 after eliminating  
434 participants for the reasons described below. All participants included in data analysis  
435 were native speakers of North American English and had typical hearing. A total of 67  
436 participants were excluded from data analysis for the following reasons: 32 for failing to  
437 use headphones during the experiment, 13 for learning a second language in which they  
438 were fluent before the age of 13 years, 6 for self-reported hearing difficulty, and 16 for

439 exposure to a language with a retroflex sound whether from family members, a caregiver,  
440 in a school environment, or through the media. These languages include unspecified  
441 Chinese ( $n = 4$ ), Mandarin ( $n = 3$ ), Vietnamese ( $n = 3$ ), Hindi ( $n = 2$ ), Russian ( $n = 2$ ),  
442 Serbian ( $n = 1$ ), and Sicilian ( $n = 1$ ). The remaining sample was comprised of 202  
443 participants. All participants gave informed consent according to the University of  
444 Connecticut Institutional Review Board and were paid \$10 per hour for participation.

## 445 **Stimuli**

446 Perception of non-native speech sounds was assessed using Hindi and Polish  
447 phonemic contrasts. For the Hindi contrast, the same dental and retroflex stimuli from  
448 Experiment 1 were used (minimal pair non-words /ɖʊg/ and /ɖ̠ʊg/). Polish stimuli  
449 consisted of minimal pair, consonant-vowel syllables that differed in whether they began  
450 with the voiceless alveolo-palatal /çɔ/ or voiceless retroflex /ʂɔ/) consonant. Stimuli were  
451 recorded by a male, native speaker of Polish in a soundproof booth and were rescaled to a  
452 mean amplitude of 65dB sound pressure level in Praat (Boersma & Weenink, 2013). All  
453 tasks in Experiment 2 were presented to participants remotely on a laptop or desktop  
454 computer through Gorilla Experiment Builder ([www.gorilla.sc](http://www.gorilla.sc), Anwyl-Irvine, Massonnié,  
455 Flitton, Kirkham, & Evershed, 2018).

456 The stimuli used to measure the perception of native-language speech categories were  
457 the same seven-step fricative (sign-shine) and stop consonant (ba-da) continua from  
458 Experiment 1.

## 459 **Procedure**

460 Participants were routed from Prolific to Gorilla to complete the experiment remotely  
461 between the dates of July 4, 2021 and October 20, 2021. After providing consent,  
462 participants completed a headphone check in which they listened to three pure tones in  
463 stereo and were asked to indicate which tone was the quietest (Woods, Siegel, Traer, &

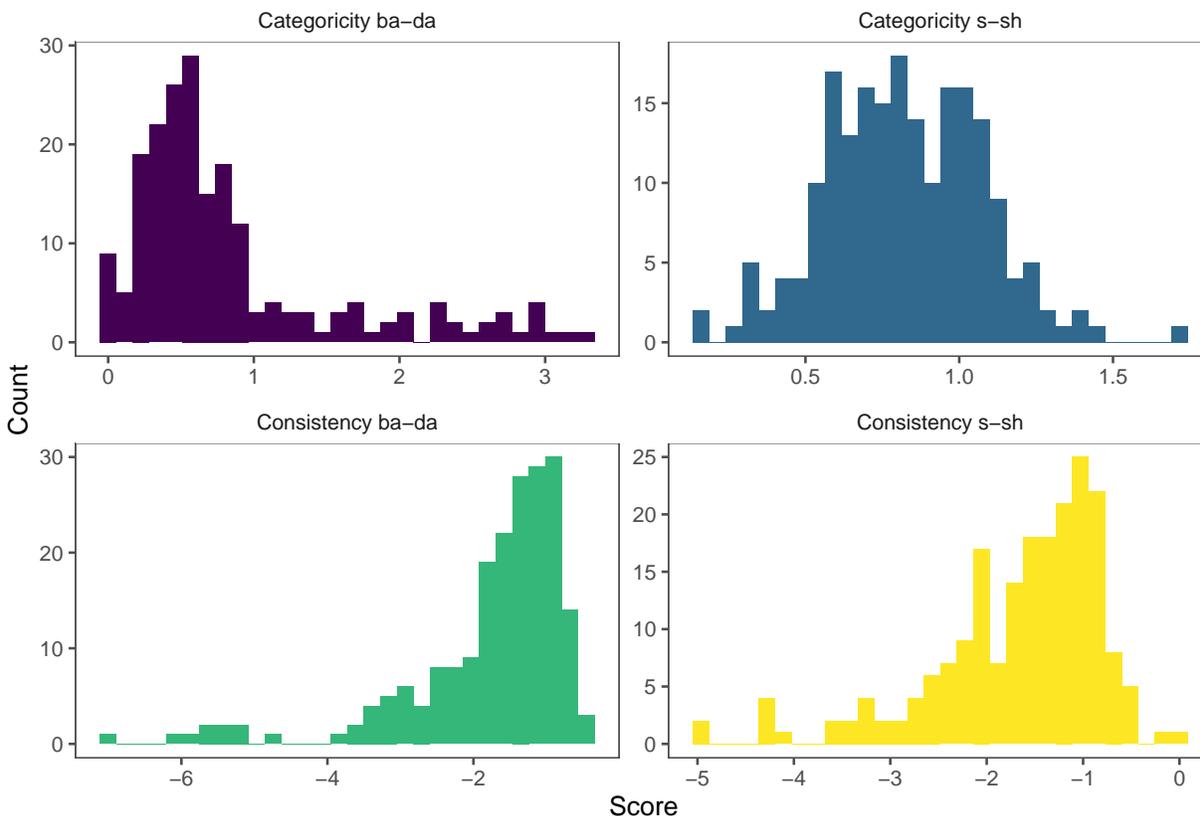


Figure 3. Histograms of categoricity and consistency measures on each continuum. Note the differences in axes among the plots.

464 McDermott, 2017). This task utilizes phase cancellation in the stimuli and can determine  
 465 whether the participant is listening to the stimuli through headphones or over loudspeakers.  
 466 If participants passed the headphone check they continued with the experiment, and if they  
 467 failed, they were asked to redo the headphone check a second time after being reminded to  
 468 wear headphones. If participants failed the headphone check the second time, they  
 469 completed the experiment although their data was excluded from analysis. Participants  
 470 then completed two questionnaires collecting information about demographics and language  
 471 experience and completed the four phonetic tasks (Hindi AX discrimination, Polish AX  
 472 discrimination, ba-da categorization, and sign-shine categorization). To address potential  
 473 concerns that the order of tasks in Experiment 1 may have influenced performance on the  
 474 visual analogue scale tasks (specifically that hearing the Hindi sounds perceived as unusual

475 exemplars of /d/ would make participants more gradient in their perception of the ba-da  
476 continuum), we set the order of tasks in this experiment so that no participant would hear  
477 similar native and non-native stimuli in two tasks that were completed consecutively.  
478 Specifically, participants completed the tasks in one of the following orders, to which they  
479 were randomly assigned: (1) Hindi AX, Polish AX, ba-da categorization, sign-shine  
480 categorization; (2) Polish AX, Hindi AX, sign-shine categorization, ba-da categorization;  
481 (3) ba-da categorization, sign-shine categorization, Hindi AX, Polish AX; (4) sign-shine  
482 categorization, ba-da categorization, Polish AX, Hindi AX. We opted for counterbalancing  
483 task order rather than having participants perform the native phonetic tasks first because  
484 we think it is equally possible that the task with the visual analogue scale would influence  
485 discrimination of the non-native contrasts. For example, we can imagine that drawing  
486 listeners' attention to subtle acoustic differences in sounds in their native language might  
487 provide them with a strategy for discriminating non-native sounds. Using a Bayesian *t*-test  
488 and Bayes factor test, we did not find evidence that categorization slopes differed on the  
489 ba-da continuum between participants depending on whether native or non-native  
490 perceptual tasks were completed first ( $BF_{10} = 1.77$ ), and we found moderate evidence for  
491 *no* difference on the s-sh continuum ( $BF_{10} = 0.17$ , see for example Schönbrodt &  
492 Wagenmakers, 2018; Wetzels & Wagenmakers, 2012, for Bayes factor interpretations).

493 **Non-native speech sound learning tasks.** Two AX discrimination tasks just  
494 like the one described in the previous experiments were administered to test participants'  
495 discrimination accuracy of the Hindi and Polish sounds<sup>5</sup>. Each task consisted of 64 trials.

496 **Native-language speech perception measures.** To measure categoricity and  
497 consistency of native-language speech sounds, we administered a similar task to the visual

---

<sup>5</sup> The inter-stimulus interval for the discrimination task of the Polish sounds was 500 ms, which differed from the Hindi task. Although the interstimulus interval can influence performance on a discrimination task, the inter-stimulus intervals from both tasks were in the range of peak discrimination performance discussed in Gerrits and Schouten (2004).

498 analogue scale described in Experiment 1. Due to constraints from the online experiment  
 499 presentation software, we did not have a slider for participants to move along a line; rather,  
 500 participants saw a bar with numbers 1-7 (the number of points on the continuum) and the  
 501 corresponding words/syllables on each side of the bar. They then clicked on a number from  
 502 1-7 to represent where between the two phonemes they thought the token belonged. This  
 503 task was similar to the visual analogue scale in Experiment 1 in that participants saw the  
 504 numbers 1-7 on a visual scale between the two phonemes.

### 505 **Analysis approach**

506 The slopes of the logistic curve to measure categoricity and the measure of  
 507 consistency in each continuum were calculated as described in Experiment 1. Figure 3  
 508 shows the distributions of the categoricity and consistency measures. We again computed  
 509 all possible correlations of these measures: categoricity between the two continua,  $r = 0.32$ ,  
 510  $p = <.001$ ; consistency between the two continua,  $r = 0.34$ ,  $p = <.001$ ; ba-da categoricity  
 511 and consistency,  $r = 0.29$ ,  $p = <.001$ ; and s-sh categoricity and consistency,  $r = 0.13$ ,  $p =$   
 512  $0.07$ .

513 We followed the same analysis procedure as in Experiment 1 to predict discrimination  
 514 accuracy of the Hindi and Polish contrasts. We fit separate models for each contrast. For  
 515 each model, the dependent variable was whether participants answered “different” (0 or 1),  
 516 and fixed effects included a deviation coded contrast for condition (same = -.5, different =  
 517 .5), the interaction of condition and categoricity (mean centered), and the interaction of  
 518 condition and consistency (mean centered).

Table 3

*Results of the Hindi discrimination task in Experiment 2.*

<i>Predictors</i>	<i>Estimate</i>	<i>SE</i>	<i>z</i>	<i>p</i>
Intercept	-0.14	0.04	-3.43	<.001

Condition (d')	0.58	0.08	6.90	<.001
Condition x Categoricity	-0.02	0.07	-0.24	0.81
Condition x Consistency	0.10	0.05	2.15	0.03

Table 4

*Results of the Polish discrimination task in Experiment 2.*

<i>Predictors</i>	<i>Estimate</i>	<i>SE</i>	<i>z</i>	<i>p</i>
Intercept	-0.20	0.09	-2.24	0.03
Condition (d')	1.17	0.17	6.70	<.001
Condition x Categoricity	0.02	0.20	0.12	0.91
Condition x Consistency	0.30	0.06	5.19	<.001

519

## Results

520

521

522

523

524

525

526

527

528

Results of the model predicting discrimination accuracy of the Hindi contrast are summarized in Table 3. Similar to the previous experiments, we see that the intercept is significantly different from zero with participants answering “same” more often, and there is a main effect of condition, suggesting participants answered “different” more often on different trials. More interestingly for the current purposes, we observe a significant interaction of condition and consistency, suggesting that more consistent responders on the ba-da task discriminated the Hindi sounds more accurately. The interaction of condition and categoricity did not reach significance. Variance inflation factors were all close to 1, suggesting no issues with multicollinearity.

529

530

531

532

Results of the model predicting discrimination accuracy of the Polish sounds are summarized in Table 4. We again see a significant intercept, main effect of condition, and interaction of condition and consistency. This suggests that participants who responded more consistently on the s-sh task also discriminated the Polish sounds more accurately.

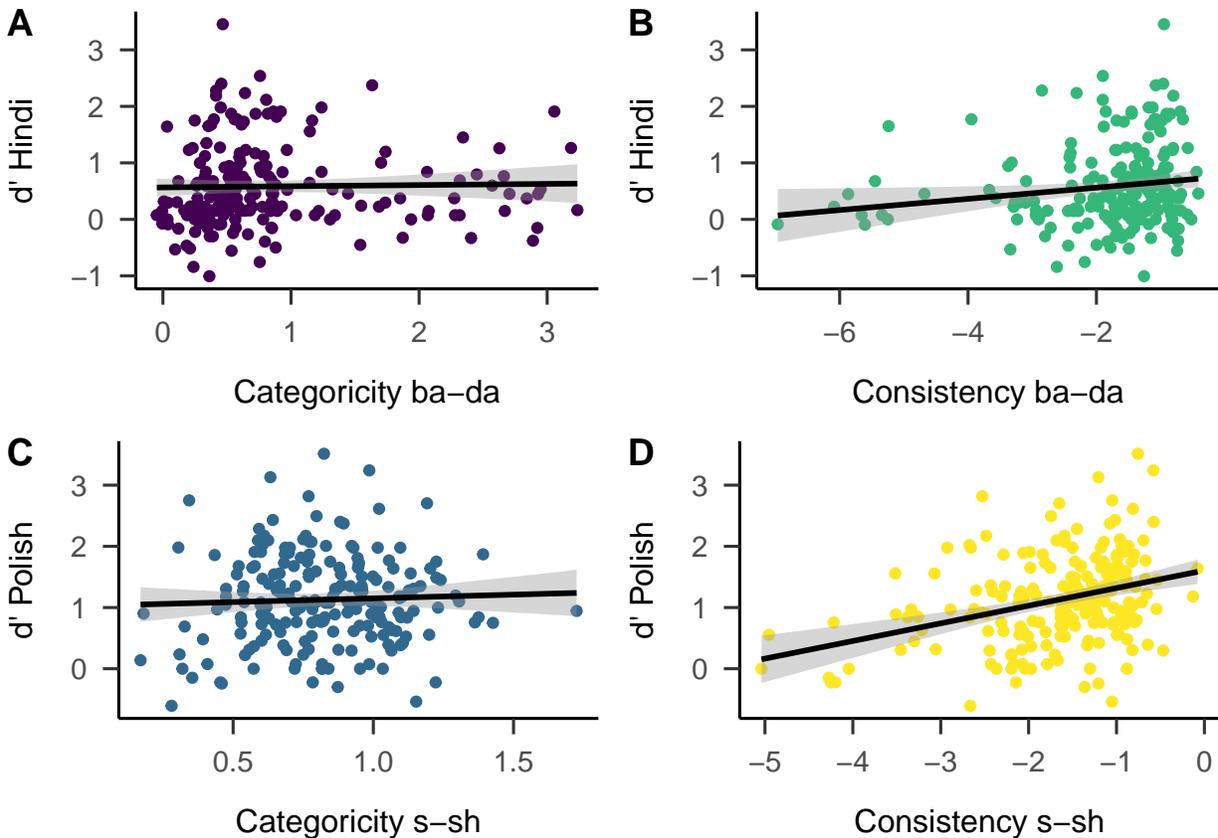


Figure 4. Relationships between discrimination accuracy of the Hindi contrast and A. ba-da categoricity and B. ba-da consistency. Relationships between discrimination accuracy of the Polish contrast and C. s-sh categoricity and D. s-sh consistency. We used  $d'$  scores for visualization purposes only; however, the statistical models predicted trial-level accuracy.

533 The interaction between condition and categoricity was not significant. Variance inflation  
 534 factors were again all close to 1 in this model, suggesting no issues with multicollinearity.  
 535 Results of Experiment 2 are displayed in Figure 4.

536

## Discussion

537 Experiment 2 was a partial replication and extension of Experiment 1. With a much  
 538 larger sample size than Experiment 1, we saw that consistency of responses on a  
 539 native-language visual analogue scale task predicted better discrimination accuracy of a  
 540 perceptually similar non-native speech contrast. This pattern of results was found for both

541 non-native contrasts we tested. This suggests that Experiment 1 was likely underpowered  
542 to detect this relationship. We return to this finding in the General Discussion.

543 Despite the larger sample in Experiment 2, categoricity still did not significantly  
544 predict discrimination accuracy of either non-native contrast. Once again, variance  
545 inflation factors for all predictors in each model were close to 1, suggesting the null result  
546 of categoricity is not simply due to multicollinearity. However, we cannot conclude there is  
547 no effect with a null result in the frequentist framework. For this reason, we did a  
548 small-scale Bayesian meta-analysis to synthesize the estimates in the two experiments  
549 (categoricity predicting naive discrimination accuracy in Experiment 1 and categoricity  
550 predicting discrimination accuracy for both sets of stimuli in Experiment 2; three estimates  
551 in total). Details on this analysis can be found in the appendix. To summarize briefly here,  
552 we performed Bayes factor tests and a sensitivity analysis (as recommended by several  
553 authors) to determine the relative evidence for the alternative hypothesis (an effect of  
554 categoricity on discrimination accuracy) over the null hypothesis under different  
555 assumptions (e.g., Hoijtink, Mulder, Lissa, & Gu, 2019; Kass & Raftery, 1995; Schad,  
556 Nicenboim, Bürkner, Betancourt, & Vasishth, 2021). We mostly found anecdotal evidence  
557 for the null hypothesis and only found moderate evidence for the null model when we use a  
558 very wide prior. Thus, we have evidence against a very large effect; however, if there is an  
559 effect of categoricity on non-native discrimination, it is possible that it is so small that we  
560 were still unable to find evidence for it even when synthesizing results across both  
561 experiments.

## 562 General discussion

563 Studies of non-native speech perception often find a wide range of individual  
564 differences in how accurately listeners perceive unfamiliar speech sounds. The goal of the  
565 current study was to test whether current theories of non-native speech perception can be  
566 extended to account for this variability in a particular way. We predicted that listeners

567 who show more graded perceptual patterns of speech sounds in their native language would  
568 discriminate difficult non-native speech sounds more accurately than those who show more  
569 categorical patterns of perception. In two different experiments, we did not find that  
570 individual differences predicted discrimination accuracy of non-native contrasts. Instead,  
571 we found that *consistency* of responses on a speech categorization task predicted  
572 discrimination accuracy.

### 573 **Individual differences in categorical perception**

574 In general, the current findings do not lend support to the idea that categoricity of  
575 native-language speech sounds is what makes differentiating non-native speech sounds hard  
576 at the *individual* level. On the basis of several theories of non-native speech sound learning  
577 (e.g., Best et al., 2001; Best & Tyler, 2007; Flege, 1995; Iverson et al., 2003; Kuhl, 1994;  
578 Kuhl et al., 2008), we reasoned that individual differences in categoricity would predict the  
579 degree to which individuals assimilate difficult non-native speech contrasts to a  
580 perceptually similar native-language speech sound. It is intuitive to assume that perceptual  
581 reorganization of speech in infancy (Kuhl et al., 2006; Werker & Tees, 1984) is related to  
582 the development of categorical perception. Specifically, if we lose perceptual sensitivity to  
583 certain non-native speech sounds and those non-native speech sounds are then assimilated  
584 to native-language speech categories during perception, it seems logical that categorical  
585 perception could be responsible for the assimilation of perceptually similar non-native  
586 sounds to native-language categories. However, categorical perception and perceptual  
587 assimilation may not be as related as is often assumed. This is consistent with the more  
588 modern view that speech perception follows a more protracted developmental pattern, in  
589 which perceptual representations continue to be refined (this includes becoming more  
590 graded) at least throughout adolescence (Hazan & Barrett, 2000; Idemaru & Holt, 2013;  
591 McMurray, 2022; McMurray et al., 2018; Nitttrouer, 2004; Zevin, 2012). Thus, graded  
592 perception of speech seems to be the norm (McMurray, 2022; McMurray et al., 2018;

593 Toscano, McMurray, Dennhardt, & Luck, 2010), though individuals differ in the extent of  
594 this gradiency (Fuhrmeister & Myers, 2021; Kapnoula et al., 2021, 2017; Kong & Edwards,  
595 2016). However, these differences may not have much of an impact on how accurately  
596 listeners can perceive non-native speech sounds.

597         It is possible, however, that the behavioral measures of categoricity that we used in  
598 the experiments presented here were not sensitive enough for measuring individual  
599 variability. For example, we used discrete points on the visual analogue scale, whereas  
600 previous studies have used continuous scales (Kapnoula et al., 2021, 2017; Kong &  
601 Edwards, 2016). Furthermore, in Experiment 2, we used a task more akin to a Likert scale,  
602 and we also found a correlation of 0.29 between the consistency and categoricity measures  
603 for the ba-da continuum. This could indicate that at least some participants were treating  
604 the task more like a two-alternative forced choice task, as consistency and categoricity  
605 measures are typically more highly correlated for two-alternative forced choice tasks (see  
606 e.g., Kapnoula et al., 2017; McMurray et al., 2018). Future work could use a more sensitive  
607 measure such as a visual analogue scale without discrete options or an online measure such  
608 as eye tracking to measure categoricity and test its relationship with non-native speech  
609 perception. Nonetheless, we did see a great deal of individual variability in the categoricity  
610 measure from the visual analogue scale that we used (see Figure 3), suggesting a decent  
611 amount of sensitivity for capturing individual differences. The distribution of scores does  
612 suggest that a handful of participants may have responded in a very categorical manner on  
613 the ba-da task; however, most participants did indeed show more graded patterns of  
614 perception, as the bulk of the distribution of scores falls between 0 and 1. It is also notable  
615 that we only saw this pattern with the stop continuum and not the fricative continuum,  
616 suggesting that this pattern of responses might have more to do with the nature of the  
617 sounds (e.g., Eimas, 1963; Healy & Repp, 1982; Repp, 1981) rather than the task itself.

**618 Perceptual consistency**

619 Results of the current study suggest that the most sensitive perceivers of novel  
620 non-native contrasts were those who had the most consistent categorization responses or  
621 ratings of native-language sounds. The fact that native-language response consistency was  
622 related to discrimination accuracy of non-native contrasts supports the hypothesis that  
623 trial-to-trial stability in the auditory response leads to stronger speech and language skills.  
624 This interpretation is consistent with findings in the field that general auditory and speech  
625 skills themselves support perception of both native and non-native speech sounds (and  
626 possibly also learning of new speech sounds). For example, Díaz, Baus, Escera, Costa, and  
627 Sebastián-Gallés (2008) reported electrophysiological differences in good vs. poor perceivers  
628 in response to both native and non-native speech sounds, which they interpreted as  
629 evidence of a speech-specific skill. Additionally, work from our lab and others has found  
630 that measures of phonological skills predict learning of a non-native speech sound contrast  
631 (Earle & Arthur, 2017; Fuhrmeister et al., 2020; see also Perrachione et al., 2011 for similar  
632 findings for learning of a non-native tonal contrast). More evidence for this idea comes  
633 from Heffner and Myers (2021), who found relationships between measures of non-native  
634 speech sound learning and adaptation to different types of speech in listeners' native  
635 language (rate and accent adaptation).

636 The relationship between consistency and discrimination accuracy might suggest that  
637 fidelity or stability of encoding of the stimulus is beneficial for accurate discrimination of  
638 non-native speech sounds. This parallels the findings from studies using the auditory  
639 brainstem response, which is a measure of how accurately the brainstem encodes an  
640 auditory stimulus (e.g., Skoe & Kraus, 2010). Fidelity of neural responses such as the  
641 auditory brainstem response has been linked to individual differences in speech perception  
642 abilities (Bidelman, Moreno, & Alain, 2013), non-native speech sound learning (Kachlicka  
643 et al., 2019; Omote, Jasmin, & Tierney, 2017), and individual differences in reading

644 abilities among unimpaired readers (Skoe, Brody, & Theodore, 2017). It is possible that we  
645 are seeing a similar phenomenon behaviorally: Listeners who encode familiar speech sounds  
646 more accurately or with higher fidelity may be able to leverage this same skill to  
647 distinguish subtle differences in unfamiliar sounds.

648         The findings presented here are also reminiscent of some findings of consistency in  
649 speech production. For example, several studies have shown that the precision or  
650 compactness of pronunciations of native-language vowels predicts perception (Kartushina  
651 & Frauenfelder, 2013) and production (Kartushina & Frauenfelder, 2014; Kartushina,  
652 Hervais-Adelman, Frauenfelder, & Golestani, 2016) accuracy of a non-native vowel  
653 contrast. In fact, Flege and Bohn (2021) recently published a revised version of the speech  
654 learning model, and based partially on these findings that production precision predicts  
655 learning of second-language speech categories, this version adds the “L1 category precision  
656 hypothesis” to the theory. This hypothesis asserts that the more precise the first-language  
657 categories are when a learner is exposed to a second language, the more likely the learner  
658 will be to establish new phonetic categories in the second language because the first and  
659 second-language sounds will be more discernible when the first-language categories are  
660 more compact/precise. Although the revised speech learning model defines category  
661 precision in terms of speech *production*, our findings suggest that consistency in *perception*  
662 may also be relevant.

663         According to the revised speech learning model, speech category precision may be  
664 related to auditory processing and acuity (Flege & Bohn, 2021). This is consistent with  
665 other studies that have found links between auditory abilities and non-native speech sound  
666 learning. For instance, Kidd, Watson, and Gygi (2007) found support for a general  
667 auditory ability for recognizing speech and non-speech sounds. Lengeris and Hazan (2010)  
668 observed that a non-speech frequency discrimination task correlated with perception of  
669 first- and second-language speech contrasts, which they take to argue that auditory

670 processing can explain individual differences in speech perception. Finally, Saito et al.  
671 (2022) found that more precise discrimination accuracy of acoustic cues that differentiate  
672 non-native speech sounds resulted in more accurate perception of those sounds.

673 The literature on the auditory brainstem response, speech production, and now our  
674 results from a speech perception task all seem to converge to suggest that consistency or  
675 precision of speech representation at multiple levels partially explains individual differences  
676 we observe in production and perception of non-native speech sounds. In sum, the category  
677 precision hypothesis in the revised speech learning model (and related accounts that  
678 prioritize the stability and precision of phonetic category representations) holds promise for  
679 the future of individual differences research in non-native speech sound learning.

680

### Conclusions

681 The findings presented here do not support the assumption that categorical  
682 perception of native-language speech sounds is what is responsible for individual differences  
683 in perception of non-native speech sounds. In contrast, our findings do suggest at least  
684 some degree of overlap of native and non-native perceptual systems, specifically with  
685 respect to consistency or precision of representations.

## References

- 686  
687 Anwyl-Irvine, A., Massonnié, J., Flitton, A., Kirkham, N., & Evershed, J. (2018). Gorillas  
688 in our midst: Gorilla. Sc, a new web-based experiment builder. *bioRxiv*, 438242.
- 689 Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for  
690 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,  
691 68(3), 255–278.
- 692 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects  
693 models using lme4. *Journal of Statistical Software*, 67(1), 1–48.  
694 <https://doi.org/10.18637/jss.v067.i01>
- 695 Best, C. T., McRoberts, G. W., & Goodell, E. (2001). Discrimination of non-native  
696 consonant contrasts varying in perceptual assimilation to the listener's native  
697 phonological system. *The Journal of the Acoustical Society of America*, 109(2), 775–794.
- 698 Best, C. T., McRoberts, G. W., & Sithole, N. M. (1988). Examination of perceptual  
699 reorganization for nonnative speech contrasts: Zulu click discrimination by  
700 english-speaking adults and infants. *Journal of Experimental Psychology: Human*  
701 *Perception and Performance*, 14(3), 345.
- 702 Best, C. T., & Tyler, M. D. (2007). Nonnative and second-language speech perception:  
703 Commonalities and complementarities. *Language Experience in Second Language*  
704 *Speech Learning: In Honor of James Emil Flege*, 1334, 1–47.
- 705 Bidelman, G. M., Moreno, S., & Alain, C. (2013). Tracing the emergence of categorical  
706 speech perception in the human auditory system. *Neuroimage*, 79, 201–212.
- 707 Boersma, P., & Weenink, D. (2013). Praat: Doing phonetics by computer [computer  
708 program]. Version 5.3. 51. *Online: Http://Www. Praat.org*, 2.
- 709 Bowles, A. R., Chang, C. B., & Karuzis, V. P. (2016). Pitch ability as an aptitude for tone  
710 learning. *Language Learning*, 66(4), 774–808.
- 711 Bradlow, A. R., Akahane-Yamada, R., Pisoni, D. B., & Tohkura, Y. (1999). Training  
712 japanese listeners to identify english/r/and/l: Long-term retention of learning in

- 713 perception and production. *Perception & Psychophysics*, *61*(5), 977–985.
- 714 Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan.  
715 *Journal of Statistical Software*, *80*(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 716 Burnham, D. K., Earnshaw, L. J., & Clark, J. E. (1991). Development of categorical  
717 identification of native and non-native bilabial stops: Infants, children and adults.  
718 *Journal of Child Language*, *18*(2), 231–260.
- 719 Clayards, M., Tanenhaus, M. K., Aslin, R. N., & Jacobs, R. A. (2008). Perception of  
720 speech reflects optimal use of probabilistic speech cues. *Cognition*, *108*(3), 804–809.
- 721 Díaz, B., Baus, C., Escera, C., Costa, A., & Sebastián-Gallés, N. (2008). Brain potentials  
722 to native phoneme discrimination reveal the origin of individual differences in learning  
723 the sounds of a second language. *Proceedings of the National Academy of Sciences*,  
724 *105*(42), 16083–16088.
- 725 Drouin, J. R., Theodore, R. M., & Myers, E. B. (2016). Lexically guided perceptual tuning  
726 of internal phonetic category structure. *The Journal of the Acoustical Society of*  
727 *America*, *140*(4), EL307–EL313.
- 728 Earle, F. S., & Arthur, D. T. (2017). Native phonological processing abilities predict  
729 post-consolidation nonnative contrast learning in adults. *The Journal of the Acoustical*  
730 *Society of America*, *142*(6), EL525–EL531.
- 731 Eimas, P. D. (1963). The relation between identification and discrimination along speech  
732 and non-speech continua. *Language and Speech*, *6*(4), 206–217.
- 733 Flege, J. E. (1995). Second language speech learning: Theory, findings, and problems.  
734 *Speech Perception and Linguistic Experience: Issues in Cross-Language Research*, *92*,  
735 233–277.
- 736 Flege, J. E., & Bohn, O.-S. (2021). The revised speech learning model (SLM-r). *Second*  
737 *Language Speech Learning: Theoretical and Empirical Progress*, 3–83.
- 738 Flege, J. E., Yeni-Komshian, G. H., & Liu, S. (1999). Age constraints on second-language  
739 acquisition. *Journal of Memory and Language*, *41*(1), 78–104.

- 740 Fox, J., & Weisberg, S. (2019). *An R companion to applied regression* (Third). Thousand  
741 Oaks CA: Sage.
- 742 Fuhrmeister, P., & Myers, E. B. (2020). Desirable and undesirable difficulties: Influences of  
743 variability, training schedule, and aptitude on nonnative phonetic learning. *Attention,*  
744 *Perception, & Psychophysics*, *82*(4), 2049–2065.
- 745 Fuhrmeister, P., & Myers, E. B. (2021). Structural neural correlates of individual  
746 differences in categorical perception. *Brain and Language*, *215*, 104919.
- 747 Fuhrmeister, P., Phillips, M. C., McCoach, D. B., & Myers, E. B. (2022 November 25/2022  
748 November 25). Relationships between native and non-native speech perception.  
749 Retrieved from <https://osf.io/2zg6c/>
- 750 Fuhrmeister, P., Schlemmer, B., & Myers, E. B. (2020). Adults show initial advantages  
751 over children learning difficult non-native speech sounds. *Journal of Speech, Language,*  
752 *and Hearing Research*, *63*(8), 2667–2679.
- 753 Gerrits, E., & Schouten, M. E. (2004). Categorical perception depends on the  
754 discrimination task. *Perception & Psychophysics*, *66*(3), 363–376.
- 755 Goh, J. X., Hall, J. A., & Rosenthal, R. (2016). Mini meta-analysis of your own studies:  
756 Some arguments on why and a primer on how. *Social and Personality Psychology*  
757 *Compass*, *10*(10), 535–549.
- 758 Golestani, N., Molko, N., Dehaene, S., LeBihan, D., & Pallier, C. (2007). Brain structure  
759 predicts the learning of foreign speech sounds. *Cerebral Cortex*, *17*(3), 575–582.
- 760 Golestani, N., Paus, T., & Zatorre, R. J. (2002). Anatomical correlates of learning novel  
761 speech sounds. *Neuron*, *35*(5), 997–1010.
- 762 Golestani, N., & Zatorre, R. J. (2009). Individual differences in the acquisition of second  
763 language phonology. *Brain and Language*, *109*(2-3), 55–67.
- 764 Hattori, K., & Iverson, P. (2009). English /r/-/l/ category assimilation by Japanese adults:  
765 Individual differences and the link to identification accuracy. *The Journal of the*  
766 *Acoustical Society of America*, *125*(1), 469–479.

- 767 Hazan, V., & Barrett, S. (2000). The development of phonemic categorization in children  
768 aged 6–12. *Journal of Phonetics*, 28(4), 377–396.
- 769 Healy, A. F., & Repp, B. H. (1982). Context independence and phonetic mediation in  
770 categorical perception. *Journal of Experimental Psychology: Human Perception and*  
771 *Performance*, 8(1), 68.
- 772 Heffner, C. C., & Myers, E. B. (2021). Individual differences in phonetic plasticity across  
773 native and nonnative contexts. *Journal of Speech, Language, and Hearing Research*,  
774 64(10), 3720–3733.
- 775 Hoijsink, H., Mulder, J., Lissa, C. van, & Gu, X. (2019). A tutorial on testing hypotheses  
776 using the bayes factor. *Psychological Methods*, 24(5), 539.
- 777 Idemaru, K., & Holt, L. L. (2013). The developmental trajectory of children’s perception  
778 and production of english/r/-/l. *The Journal of the Acoustical Society of America*,  
779 133(6), 4232–4246.
- 780 Iverson, P., Kuhl, P. K., Akahane-Yamada, R., Diesch, E., Kettermann, A., Siebert, C., et  
781 al. (2003). A perceptual interference account of acquisition difficulties for non-native  
782 phonemes. *Cognition*, 87(1), B47–B57.
- 783 Kachlicka, M., Saito, K., & Tierney, A. (2019). Successful second language learning is tied  
784 to robust domain-general auditory processing and stable neural representation of sound.  
785 *Brain and Language*, 192, 15–24.
- 786 Kapnoula, E. C., Edwards, J., & McMurray, B. (2021). Gradient activation of speech  
787 categories facilitates listeners’ recovery from lexical garden paths, but not perception of  
788 speech-in-noise. *Journal of Experimental Psychology: Human Perception and*  
789 *Performance*, 47(4), 578.
- 790 Kapnoula, E. C., Winn, M. B., Kong, E. J., Edwards, J., & McMurray, B. (2017).  
791 Evaluating the sources and functions of gradiency in phoneme categorization: An  
792 individual differences approach. *Journal of Experimental Psychology: Human*  
793 *Perception and Performance*, 43(9), 1594.

- 794 Kartushina, N., & Frauenfelder, U. H. (2013). On the role of L1 speech production in L2  
795 perception: Evidence from spanish learners of french. *Proceedings of the 14th*  
796 *Interspeech Conference*, 2118–2122.
- 797 Kartushina, N., & Frauenfelder, U. H. (2014). On the effects of L2 perception and of  
798 individual differences in L1 production on L2 pronunciation. *Frontiers in Psychology*, 5,  
799 1246.
- 800 Kartushina, N., Hervais-Adelman, A., Frauenfelder, U. H., & Golestani, N. (2016). Mutual  
801 influences between native and non-native vowels in production: Evidence from  
802 short-term visual articulatory feedback training. *Journal of Phonetics*, 57, 21–39.
- 803 Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical*  
804 *Association*, 90(430), 773–795.
- 805 Kempe, V., Bublitz, D., & Brooks, P. J. (2015). Musical ability and non-native  
806 speech-sound processing are linked through sensitivity to pitch and spectral  
807 information. *British Journal of Psychology*, 106(2), 349–366.
- 808 Kidd, G. R., Watson, C. S., & Gygi, B. (2007). Individual differences in auditory abilities.  
809 *The Journal of the Acoustical Society of America*, 122(1), 418–435.
- 810 Kong, E. J., & Edwards, J. (2016). Individual differences in categorical perception of  
811 speech: Cue weighting and executive function. *Journal of Phonetics*, 59, 40–57.
- 812 Kuhl, P. K. (1994). Learning and representation in speech and language. *Current Opinion*  
813 *in Neurobiology*, 4(6), 812–822.
- 814 Kuhl, P. K., Conboy, B. T., Coffey-Corina, S., Padden, D., Rivera-Gaxiola, M., & Nelson,  
815 T. (2008). Phonetic learning as a pathway to language: New data and native language  
816 magnet theory expanded (NLM-e). *Philosophical Transactions of the Royal Society B:*  
817 *Biological Sciences*, 363(1493), 979–1000.
- 818 Kuhl, P. K., Stevens, E., Hayashi, A., Deguchi, T., Kiritani, S., & Iverson, P. (2006).  
819 Infants show a facilitation effect for native language phonetic perception between 6 and  
820 12 months. *Developmental Science*, 9(2), F13–F21.

- 821 Lengeris, A., & Hazan, V. (2010). The effect of native vowel processing ability and  
822 frequency discrimination acuity on the phonetic training of english vowels for native  
823 speakers of greek. *The Journal of the Acoustical Society of America*, 128(6), 3757–3768.
- 824 Liberman, A. M., Cooper, F. S., Shankweiler, D. P., & Studdert-Kennedy, M. (1967).  
825 Perception of the speech code. *Psychological Review*, 74(6), 431.
- 826 Liberman, A. M., Harris, K. S., Hoffman, H. S., & Griffith, B. C. (1957). The  
827 discrimination of speech sounds within and across phoneme boundaries. *Journal of*  
828 *Experimental Psychology*, 54(5), 358.
- 829 Lim, S., & Holt, L. L. (2011). Learning foreign sounds in an alien world: Videogame  
830 training improves non-native speech categorization. *Cognitive Science*, 35(7),  
831 1390–1405.
- 832 Lisker, L. (2001). Hearing the polish sibilants [s š ś]: Phonetic and auditory judgements.  
833 *Travaux Du Cercle Linguistique de Copenhague XXXI. To Honour Eli*  
834 *Fischer-Jørgensen*, 226–238.
- 835 Luthra, S., Fuhrmeister, P., Molfese, P. J., Guediche, S., Blumstein, S. E., & Myers, E. B.  
836 (2019). Brain-behavior relationships in incidental learning of non-native phonetic  
837 categories. *Brain and Language*, 198, 104692.
- 838 Macmillan, N. A., & Creelman, C. D. (2004). *Detection theory: A user's guide*. Psychology  
839 press.
- 840 Mathôt, S., Schreij, D., & Theeuwes, J. (2012). OpenSesame: An open-source, graphical  
841 experiment builder for the social sciences. *Behavior Research Methods*, 44(2), 314–324.
- 842 Mayr, R., & Escudero, P. (2010). Explaining individual variation in L2 perception:  
843 Rounded vowels in english learners of german. *Bilingualism: Language and Cognition*,  
844 13(3), 279–297.
- 845 McGuire, G. (2007). English listeners' perception of polish alveopalatal and retroflex  
846 voiceless sibilants: A pilot study. *UC Berkeley PhonLab Annual Report*, 3(3).
- 847 McMurray, B. (2022). *The acquisition of speech categories: Beyond perceptual narrowing*,

- 848 *beyond unsupervised learning and beyond infancy.*
- 849 McMurray, B., Danelz, A., Rigler, H., & Seedorff, M. (2018). Speech categorization  
850 develops slowly through adolescence. *Developmental Psychology, 54*(8), 1472.
- 851 McMurray, B., Tanenhaus, M. K., & Aslin, R. N. (2002). Gradient effects of  
852 within-category phonetic variation on lexical access. *Cognition, 86*(2), B33–B42.
- 853 Miller, J. L. (1997). Internal structure of phonetic categories. *Language and Cognitive*  
854 *Processes, 12*(5-6), 865–870.
- 855 Nittrouer, S. (2004). The role of temporal and dynamic signal components in the  
856 perception of syllable-final stop voicing by children and adults. *The Journal of the*  
857 *Acoustical Society of America, 115*(4), 1777–1790.
- 858 Omote, A., Jasmin, K., & Tierney, A. (2017). Successful non-native speech perception is  
859 linked to frequency following response phase consistency. *Cortex, 93*, 146–154.
- 860 Perrachione, T. K., Lee, J., Ha, L. Y., & Wong, P. C. (2011). Learning a novel phonological  
861 contrast depends on interactions between individual differences and training paradigm  
862 design. *The Journal of the Acoustical Society of America, 130*(1), 461–472.
- 863 Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D., & R Core Team. (2021). *nlme: Linear and*  
864 *nonlinear mixed effects models*. Retrieved from  
865 <https://CRAN.R-project.org/package=nlme>
- 866 Piske, T., MacKay, I. R., & Flege, J. E. (2001). Factors affecting degree of foreign accent  
867 in an L2: A review. *Journal of Phonetics, 29*(2), 191–215.
- 868 Polka, L. (1991). Cross-language speech perception in adults: Phonemic, phonetic, and  
869 acoustic contributions. *The Journal of the Acoustical Society of America, 89*(6),  
870 2961–2977.
- 871 R Core Team. (2021). *R: A language and environment for statistical computing*. Vienna,  
872 Austria: R Foundation for Statistical Computing. Retrieved from  
873 <https://www.R-project.org/>
- 874 Razzaghi, M. (2013). The probit link function in generalized linear models for data mining

- 875 applications. *Journal of Modern Applied Statistical Methods*, 12(1), 19.
- 876 Repp, B. H. (1981). Two strategies in fricative discrimination. *Perception &*  
877 *Psychophysics*, 30(3), 217–227.
- 878 Saito, K., Kachlicka, M., Suzukida, Y., Petrova, K., Lee, B. J., & Tierney, A. (2022).  
879 Auditory precision hypothesis-L2: Dimension-specific relationships between auditory  
880 processing and second language segmental learning. *Cognition*, 229, 105236.
- 881 Schad, D. J., Nicenboim, B., Bürkner, P.-C., Betancourt, M., & Vasishth, S. (2021).  
882 Workflow techniques for the robust use of bayes factors. *arXiv Preprint*  
883 *arXiv:2103.08744*.
- 884 Schad, D. J., Vasishth, S., Hohenstein, S., & Kliegl, R. (2020). How to capitalize on a  
885 priori contrasts in linear (mixed) models: A tutorial. *Journal of Memory and Language*,  
886 110, 104038.
- 887 Schertz, J., Cho, T., Lotto, A., & Warner, N. (2015). Individual differences in phonetic cue  
888 use in production and perception of a non-native sound contrast. *Journal of Phonetics*,  
889 52, 183–204.
- 890 Schönbrodt, F. D., & Wagenmakers, E.-J. (2018). Bayes factor design analysis: Planning  
891 for compelling evidence. *Psychonomic Bulletin & Review*, 25(1), 128–142.
- 892 Silbert, N. H., Smith, B. K., Jackson, S. R., Campbell, S. G., Hughes, M. M., & Tare, M.  
893 (2015). Non-native phonemic discrimination, phonological short term memory, and  
894 word learning. *Journal of Phonetics*, 50, 99–119.
- 895 Skoe, E., Brody, L., & Theodore, R. M. (2017). Reading ability reflects individual  
896 differences in auditory brainstem function, even into adulthood. *Brain and Language*,  
897 164, 25–31.
- 898 Skoe, E., & Kraus, N. (2010). Auditory brainstem response to complex sounds: A tutorial.  
899 *Ear and Hearing*, 31(3), 302.
- 900 Slevc, L. R., & Miyake, A. (2006). Individual differences in second-language proficiency:  
901 Does musical ability matter? *Psychological Science*, 17(8), 675–681.

- 902 Stevens, K. N., & Blumstein, S. E. (1975). Quantal aspects of consonant production and  
903 perception: A study of retroflex stop consonants. *Journal of Phonetics*, *3*(4), 215–233.
- 904 Sutton, A. J., & Abrams, K. R. (2001). Bayesian methods in meta-analysis and evidence  
905 synthesis. *Statistical Methods in Medical Research*, *10*(4), 277–303.
- 906 Thompson, C. G., Kim, R. S., Aloe, A. M., & Becker, B. J. (2017). Extracting the variance  
907 inflation factor and other multicollinearity diagnostics from typical regression results.  
908 *Basic and Applied Social Psychology*, *39*(2), 81–90.
- 909 Toscano, J. C., McMurray, B., Dennhardt, J., & Luck, S. J. (2010). Continuous perception  
910 and graded categorization: Electrophysiological evidence for a linear relationship  
911 between the acoustic signal and perceptual encoding of speech. *Psychological Science*,  
912 *21*(10), 1532–1540.
- 913 Werker, J. F., & Tees, R. C. (1984). Cross-language speech perception: Evidence for  
914 perceptual reorganization during the first year of life. *Infant Behavior and*  
915 *Development*, *7*(1), 49–63.
- 916 Wetzels, R., & Wagenmakers, E.-J. (2012). A default bayesian hypothesis test for  
917 correlations and partial correlations. *Psychonomic Bulletin & Review*, *19*(6), 1057–1064.
- 918 Woods, K. J., Siegel, M. H., Traer, J., & McDermott, J. H. (2017). Headphone screening to  
919 facilitate web-based auditory experiments. *Attention, Perception, & Psychophysics*,  
920 *79*(7), 2064–2072.
- 921 Wright, D. B., Horry, R., & Skagerberg, E. M. (2009). Functions for traditional and  
922 multilevel approaches to signal detection theory. *Behavior Research Methods*, *41*(2),  
923 257–267.
- 924 Zevin, J. D. (2012). A sensitive period for shibboleths: The long tail and changing goals of  
925 speech perception over the course of development. *Developmental Psychobiology*, *54*(6),  
926 632–642.
- 927 Żygis, M., & Padgett, J. (2010). A perceptual study of polish fricatives, and its  
928 implications for historical sound change. *Journal of Phonetics*, *38*(2), 207–226.

## Appendix

## Meta-analysis of estimates in the two experiments

929 In two different experiments, we did not see that categoricity significantly predicted  
930 discrimination accuracy of a non-native speech contrast. However, with a null result in the  
931 frequentist framework, we cannot claim to have *no* effect. Furthermore, we saw a fair  
932 amount of between-experiment variability, as evidenced by the range in standard errors of  
933 the effect of interest. We therefore conducted a Bayesian meta-analysis of the effect  
934 (interaction between condition and categoricity) that included the experiments presented  
935 in current study. Meta-analyses are useful for synthesizing results of several experiments,  
936 and they estimate the effect of interest and the precision of that effect using the estimates  
937 and standard errors from the individual studies. We acknowledge that we have a limited  
938 number of data sets for a meta-analysis; however, Goh, Hall, and Rosenthal (2016)  
939 advocate specifically for what they call a “mini meta-analysis,” in which authors perform a  
940 meta-analysis on similar experiments reported in a paper. We further chose to perform the  
941 meta-analysis in the Bayesian framework in order to perform Bayes factor tests to test  
942 whether we have relative evidence for the alternative hypothesis (an effect of categoricity  
943 on discrimination accuracy) over the null hypothesis.

944 **Data sets and extraction of estimates**

945 Meta-analyses synthesize data across several experiments using the estimate of the  
946 effect in question (here the interaction of categoricity and condition) and its standard  
947 error. We included three estimates and their standard errors in this meta-analysis. For  
948 Experiment 1, we fit a model using only the data from the discrimination pretest because  
949 that is common to both experiments. We additionally only used the measure of  
950 categoricity from the ba-da continuum from Experiment 1 to be more like Experiment 2.  
951 For each data set (Experiment 1, Experiment 2 Hindi contrast, Experiment 2 Polish  
952 contrast), we fit a generalized linear mixed effects model with a probit link that predicted

953 whether the participant answered “different” (0 = no, 1 = yes). Fixed effects included  
954 condition (deviation coded, different = .5, same = -.5) and the interaction of condition and  
955 categoricity (mean centered). When possible, the full random effects structure was  
956 included and random effects were removed one by one in case the model did not converge  
957 without warnings. The estimates and standard errors from the interaction effect were  
958 extracted from each model and entered into the meta-analysis.

### 959 **Meta-analysis**

960 We conducted a random effects meta-analysis (as opposed to a fixed effects  
961 meta-analysis), which assumes a different true effect for each data set (Sutton & Abrams,  
962 2001). To do this, we fit a Bayesian generalized linear mixed effects model that estimates  
963 an intercept from the estimates of the effects from each data set weighted by their standard  
964 error with a random intercept for data set. This model was fit using the brms package  
965 (Bürkner, 2017). To test whether we have evidence for this effect (the interaction of  
966 categoricity and condition), we did a Bayes factor test using bridge sampling, which  
967 informs us on the relative evidence for one model over another (in this case the alternative  
968 vs. a null model). Bayes factors are supposed to be taken as graded evidence: A Bayes  
969 factor of 1 indicates no evidence for one model over another, Bayes factors between 1 and 3  
970 or between 1/3 and 1 are considered anecdotal or inconclusive evidence for one model over  
971 another, and Bayes factors over 3 or less than 1/3 are typically considered to indicate at  
972 least moderate evidence for one model over the other (Schönbrodt & Wagenmakers, 2018;  
973 Wetzels & Wagenmakers, 2012).

974 The Bayesian framework allows us to incorporate prior beliefs into the statistical  
975 model. Bayes factors are sensitive to the priors used in the model, so we did a sensitivity  
976 analysis, in which we report Bayes factors for a range of priors, as recommended by several  
977 authors (e.g., Hoijsink et al., 2019; Kass & Raftery, 1995; Schad et al., 2021). The effect of  
978 interest here is the interaction effect (interaction of categoricity and condition), and for

979 that effect, we chose a normal prior with a mean of 0 and a range of standard deviations  
 980 from .05 to .5. For the standard deviation, we chose regularizing priors from a normal  
 981 distribution with a mean of 0 and standard deviation of 1.

Table A1

*Results of the meta-analysis and sensitivity analyses. We report the prior, estimate, credible interval (CrI), and bayes factor in favor of the alternative hypothesis (BF10).*

<i>Prior</i>	<i>Estimate</i>	<i>CrI</i>	<i>BF10</i>
Normal (0,0.05)	0.005	[-0.09, 0.09]	0.90
Normal (0,0.1)	0.011	[-0.15, 0.16]	0.76
Normal (0,0.2)	0.017	[-0.24, 0.26]	0.53
Normal (0,0.3)	0.017	[-0.31, 0.33]	0.40
Normal (0,0.4)	0.015	[-0.37, 0.36]	0.32
Normal (0,0.5)	0.013	[-0.4, 0.39]	0.26

## 982 Results and discussion

983 Results of the meta-analysis including sensitivity analyses are summarized in Table  
 984 A1. We report the prior for the effect of interest used in each Bayes factor test, the  
 985 posterior mean (i.e., meta-analytic estimate), the 95% credible interval, and the Bayes  
 986 factor in favor of the alternative hypothesis.

987 Most of the evidence is inconclusive; however, when we assume a larger effect size  
 988 (prior width of 0.5), we see moderate evidence for the null hypothesis. In no case did we  
 989 see evidence for the alternative hypothesis. At the very least, we can conclude that we have  
 990 evidence against a large effect. It is nonetheless possible that there is a true *small* effect of  
 991 categoricity on discrimination accuracy of a non-native contrast; however, we may need  
 992 even more data than we have collected in the two experiments here to obtain evidence for

993 such a small effect.