

¹ The effects of variability and similarity on talker-independent adaptation to
² Spanish-accented English

³ Holly A. Zaharchuk^{a,b,*}, Janet G. van Hell^a

^aDepartment of Psychology, The Pennsylvania State University,

^bDepartment of Speech, Language, and Hearing Sciences, University of Connecticut,

⁴ **Abstract**

to-do

⁵ **Keywords:**

to-do

⁷ **1. Introduction**

⁸ The production of a single speech sound not only varies diachronically within an individual talker, but
⁹ also varies synchronically across groups of talkers. Listeners have to accommodate multiple sources of
¹⁰ variation in the speech signal in order to achieve successful comprehension. Perceptual adaptation is the
¹¹ process of detecting the *patterns* of variation in the speech signal and applying them during comprehension.
¹² This learning process allows listeners to map unfamiliar, variable, or secondary acoustic cues onto their
¹³ stable phonetic categories. Previous research has shown that perceptual adaptation to specific cue-category
¹⁴ mappings can transfer from one talker to another (e.g., Kraljic and Samuel, 2006, 2007; Reinisch and Holt,
¹⁵ 2014). In addition, the benefits of exposure to an L2 accent can transfer to a new talker with the same accent.
¹⁶ However, there is disagreement in the literature about the conditions under which adaptation generalizes
¹⁷ across talkers. Two different factors have been argued to facilitate generalization: (1) *variability* during
¹⁸ exposure (Baese-Berk et al., 2013; Bradlow and Bent, 2008) and (2) *similarity* between the exposure and test
¹⁹ talkers (Xie and Myers, 2017; Xie et al., 2021). The present study features three experiments that directly
²⁰ compare variability and similarity during exposure to Spanish-accented English. The overarching goal is to
²¹ understand how listeners use acoustic-phonetic information to accommodate linguistic variation.

²² **1.1. L2-accented speech**

²³ During L2 acquisition, a speaker's L1 phonetic inventory shapes the production of L2 speech sounds
²⁴ (Flege and Bohn, 2021; for a recent review, see Nagle and Baese-Berk, 2022). This cross-language influence in
²⁵ bilingual speech production is a product of the relations between L1 and L2 phonetic categories (Flege, 2007;
²⁶ Flege et al., 2003). These relations determine how an L2 phone is categorized within the L1 phonological
²⁷ system (Best et al., 2001; Best and Tyler, 2007). In the experiments reported here, we focused on two-category
²⁸ assimilation (Best, 1994; Tyler et al., 2014), where two contrasting L2 phonemes map onto two contrasting L1
²⁹ phonemes. Specifically, we focused on three contrasting phoneme pairs found in both English and Spanish:
³⁰ /p/-/b/ (*park* vs. *bark*), /t/-/d/ (*tune* vs. *dune*), and /k/-/g/ (*coal* vs. *goal*). These consonant pairs all
³¹ have the same manner of articulation (stop) but different places of articulation (bilabial, alveolar, and
³² velar, respectively). Critically, within each pair, one member differs from the other in terms of voicing
³³ (voiceless-voiced, respectively). Voicing refers to the status of the glottis, which comprises two folds of tissue

*Corresponding author

Email address: holly.zaharchuk@uconn.edu (Holly A. Zaharchuk)

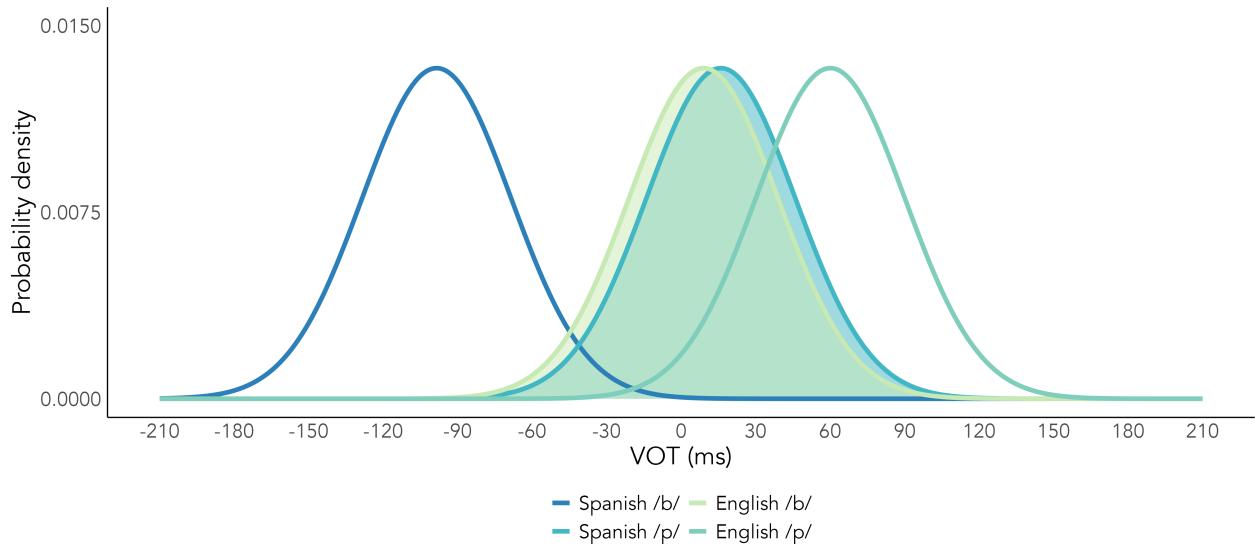


Figure 1: Overlapping VOT distributions for Spanish voiceless stops and English voiced stops.

located in the larynx at the top of the neck. Constricting the vocal folds causes them to vibrate, resulting in voiced speech sounds. By contrast, opening the vocal folds limits vibration, resulting in voiceless speech sounds. The consonants /p/, /t/, and /k/ are voiceless, while the consonants /b/, /d/, and /g/ are voiced. This small articulatory distinction plays an outsized role in the phonological systems of many languages, including English and Spanish. However, it can be difficult for listeners to detect this contrast in continuous speech. How, then, can a listener tell the difference between a talker's /p/ and /b/?

Throughout language acquisition, listeners (and speakers themselves) learn which acoustic-phonetic features are most strongly associated with this fine-grained phonological voicing contrast. There are several interrelated cues that listeners use to distinguish voiced from voiceless stops in word-initial position, including onset fundamental frequency (Whalen et al., 1993), formant transitions (Cooper, 1974; Benkí, 2001), and vowel duration (Port and Dalby, 1982; Viswanathan et al., 2020). However, the primary cue to word-initial stop consonant voicing is voice onset time (VOT), which is the interval between the release of the closure of the articulators, called the burst, and the beginning of voicing (Lisker and Abramson, 1964; Cho and Ladefoged, 1999; Chodroff et al., 2019). Listeners are highly sensitive to the probabilistic distributions of VOT across voiced and voiceless English stops (Clayards et al., 2008).

The associations between VOT and stop consonants vary cross-linguistically (Chodroff et al., 2019; Lisker and Abramson, 1970). There are three general types of VOT: lead, where voicing begins before the release, resulting in negative values; short lag, where voicing begins immediately or relatively quickly after the release; and long lag, where voicing begins relatively slowly after the release (Abramson and Whalen, 2017). In Spanish, lead VOTs correspond to voiced stops and short lag VOTs correspond to voiceless stops (Vicente, 1986; Williams, 1977). For example, the /b/ in Spanish *barco* (boat) would be approximately -98 ms on average, while the /p/ in Spanish *parque* (park) would be approximately 16 ms on average (values calculated from Chodroff et al., 2019). By contrast, in English, short lag VOTs correspond to voiced stops and long lag VOTs correspond to voiceless stops (Chodroff and Wilson, 2017). For example, the /b/ in English *bark* would be approximately 9 ms on average, while the /p/ in English *park* would be approximately 60 ms on average (values calculated from Chodroff et al., 2019). In other words, the VOT for Spanish /p/ is a better fit for the contrasting English /b/ phoneme category than for the corresponding English /p/ phoneme category. The overlap in VOT between Spanish /p/ and English /b/ is illustrated in Figure 1, with probability densities constructed from the data presented in Chodroff et al. (2019).

To return to the notion of two-category assimilation, consider an L1 Spanish speaker who is learning English as an L2. She already has an established /p/-/b/ contrast in her L1 Spanish that is associated

65 with VOT. The /p/-/b/ contrast in her L2 English is associated with the same cue in the same direction,
66 such that VOTs for /p/ are longer than those for /b/. As a result, the English /p/-/b/ contrast will be
67 relatively easy to integrate into her existing phonological system compared to, for example, certain vowel
68 contrasts in English (Baigorri et al., 2019). The same will be true for both the /t/-/d/ and /k/-/g/ contrasts.
69 However, assimilating the L2 English /p/ to the L1 Spanish /p/ means that the English phoneme will be
70 produced like the Spanish phoneme. That is, L2 English /p/ will be produced with a short lag VOT (similar
71 to Spanish /p/). This will create a perceptual problem from a listener's perspective, since short lag VOTs
72 are associated with English /b/, not /p/; without supporting contextual information, listeners may confuse
73 Spanish-accented English /p/ with /b/ (e.g., *park* perceived as *bark*). To summarize, cross-language transfer
74 in Spanish-accented English can create ambiguity between voiced and voiceless stops (Flege and Eefting,
75 1987). Thankfully, listeners are adept at resolving such ambiguities through perceptual adaptation.

76 1.2. *Adaptation to L2-accented speech*

77 A formal way of understanding perceptual adaptation is provided by the ideal adapter framework
78 (Kleinschmidt and Jaeger, 2015; Kleinschmidt, 2019). This model is built on three core assumptions: (1) the
79 relations between acoustic cues and phonetic categories are probabilistic; (2) novel cue-category relations are
80 incorporated into the perceptual system through statistical learning; and (3) cue distributions covary with
81 socio-indexical categories. If a listener is an “ideal adapter,” then she will be sensitive to the probability
82 distributions of cue-category mappings (1). Figure 1 provides examples of the probability distributions
83 for the acoustic cue VOT over the phonetic categories /p/ and /b/. An ideal adapter is able to form new
84 representations of cue-category distributions or update her existing ones through statistical learning (2). The
85 specificity of these cue-category representations, called generative models, depends on two factors: how much
86 information is gained about a cue distribution and how well a category predicts a cue by using a particular
87 grouping (Kleinschmidt, 2019). An ideal adapter will structure her generative models according to the
88 most “informative” and “useful” socio-indexical groupings (3). We will return to the structure of a listener’s
89 generative models in relation to L2-accented talkers in Section 1.3 after reviewing the relevant literature.

90 Previous research shows that listeners adapt quickly to L2-accented speech (Bent and Baese-Berk, 2021).
91 A typical experimental design features an exposure phase, in which participants gain experience with a
92 speech pattern, followed by a test phase. Talker-specific adaptation is the result of hearing the same talker
93 during exposure and test. For example, participants who train on Talker A during exposure tend to perform
94 better on Talker A during test than participants who train on Talker B (Bradlow and Bent, 2008; Clarke and
95 Garrett, 2004; Xie et al., 2021; Xie et al., 2017; Xie et al., 2018; cf. Bradlow et al., 2023). Talker-independent
96 adaptation, or generalization, is the result of exposure to a type of talker rather than exposure to a specific
97 talker. For example, participants who train on Talkers A, B, and C with Accent X tend to perform better on
98 Talker D with Accent X than participants who trained on Talkers E, F, and G with Accent Y (Alexander
99 and Nygaard, 2019; Bradlow and Bent, 2008; Sidaras et al., 2009; Tzeng et al., 2016; Xie et al., 2021, 2018).
100 Overall, experience with L2-accented talkers improves perception of novel L2-accented talkers (Witteman
101 et al., 2013; Tzeng et al., 2024; Reinisch and Holt, 2014; Bieber and Gordon-Salant, 2022; Baese-Berk et al.,
102 2013).

103 1.2.1. *Talker-independent adaptation to L2-accented speech*

104 The pattern of results for talker-independent adaptation becomes more complicated when considering
105 different types of exposure to an accent. Studies often compare single-talker exposure, in which one talker
106 produces all of the stimuli for a condition, to multi-talker exposure, in which several different talkers produce
107 the stimuli for a condition. For example, participants in a single-talker exposure condition listen to Talker
108 A with Accent X and participants in a multi-talker exposure condition listen to Talkers B, C, and D with
109 Accent X. During test, both groups listen to Talker E with Accent X. Differences in task performance with
110 Talker E index generalization from exposure (Baese-Berk et al., 2013; Bradlow and Bent, 2008; Xie and
111 Myers, 2017; Xie et al., 2021).

112 Experiment 2 of Bradlow and Bent (2008) illustrates how exposure generalizes across L2-accented talkers
113 with the same L1. In this experiment, participants completed a sentence transcription task during both

114 exposure and test, and performance was measured in terms of sentence recognition accuracy. There were four
115 exposure conditions: talker-specific, multi-talker, single-talker, and control. Participants in the talker-specific
116 condition trained directly on the Mandarin-accented English test talker, while those in the multi-talker
117 condition were exposed to five novel Mandarin-accented English talkers. Participants in the single-talker
118 condition were exposed to just one of the Mandarin-accented English talkers from the multi-talker condition.
119 Control training featured five L1-accented English talkers. Training was followed by two post-tests: one
120 featured a novel talker with a familiar L2 accent (Mandarin) and the other featured a novel talker with an
121 unfamiliar L2 accent (Slovakian). Performance on the Mandarin-accented English test talker was higher in
122 the talker-specific and multi-talker conditions than in the single-talker and control conditions, which did not
123 differ from one another. By contrast, performance on the Slovakian-accented English test talker did not differ
124 between conditions. These results suggest that exposure to multiple talkers with the same L1 highlights the
125 features that characterize a specific L2 accent.

126 Baese-Berk et al. (2013) provided evidence that exposure to multiple L2-accented talkers with *different*
127 L1s also facilitates generalization. They used the same stimulus materials and procedures as Experiment 2 in
128 Bradlow and Bent (2008) in order to compare a multi-accent exposure group directly to the multi-talker
129 exposure and control groups from the previous study. During multi-accent training, participants were exposed
130 to five L2-accented English talkers, each of whom had a different L1: Thai, Korean, Hindi, Romanian, and
131 Mandarin. Multi-accent and multi-talker exposure yielded higher performance on the Mandarin-accented test
132 talker than control exposure. Multi-accent exposure also generalized to the Slovakian-accented test talker,
133 while multi-talker and control exposure did not. Taken together with the results of Bradlow and Bent (2008),
134 these results suggest that exposure to multiple L2-accented talkers facilitates talker-independent adaptation
135 to a greater extent than exposure to a single L2-accented talker. A multi-talker versus single-talker advantage
136 was found even when the L2-accented talkers in the multi-talker condition did not share an L1 with the test
137 talker. By contrast, Xie et al. (2021) and Xie and Myers (2017) both observed comparable talker-independent
138 adaptation effects for multi- and single-talker conditions.

139 Specifically, Xie et al. (2021) sought to replicate Bradlow and Bent (2008) with one key change: counter-
140 balancing the test and exposure talkers. In the original study, different (partially overlapping) sets of
141 Mandarin-accented English talkers were used in the multi-talker and single-talker exposure conditions. Thus,
142 the lack of generalization from single-talker training may have been the result of the specific talkers, rather
143 than the superiority of multi-talker exposure per se. The specific combination of exposure and test talkers
144 is likely to have affected performance under both exemplar-based (e.g., Goldinger, 1998; Johnson, 2006),
145 and hybrid (e.g., Pierrehumbert, 2016; Kleinschmidt and Jaeger, 2015) models of speech perception (see
146 Introduction to Xie et al., 2021). To address this potential confound, Xie et al. (2021) used a single set of
147 Mandarin-accented talkers for both multi-talker and single-talker exposure; in addition, these talkers were
148 rotated through the exposure talker and test talker roles. The results replicated Bradlow and Bent (2008)'s
149 finding that multi-talker exposure facilitates generalization more than control exposure. However, contrary
150 to Bradlow and Bent (2008), they found that single-talker exposure also facilitated generalization more than
151 control exposure. When comparing multi- and single-talker exposure, the difference in performance depended
152 on the particular combination of exposure and test talkers. Overall, these results suggest that generalization
153 is strongly influenced by talker-specific features.

154 Xie and Myers (2017) investigated the acoustic-phonetic features of the exposure talkers that facilitate
155 generalization. In this study, participants were exposed to Mandarin-accented English through an auditory
156 lexical decision task. Experimental exposure included multisyllabic real words that biased listeners toward
157 perceiving ambiguous word-final stops as /d/ rather than as /t/ (e.g., *overload*). Control exposure did
158 not include these critical items. Both types of exposure featured five Mandarin-accented talkers; the only
159 difference was in the presence of disambiguating lexical contexts for learning ambiguous word-final /d/.
160 During test, participants performed a primed cross-modal lexical decision task with a novel Mandarin-accented
161 talker. Previous exposure to Mandarin-accented word-final English /d/ in disambiguating lexical contexts
162 generalized to the novel talker. Specifically, lexical activation for the /d/-final member of minimal pairs
163 like *seed* versus *seat* was increased in the experimental versus the control group. The results suggest that
164 listeners did not simply expand their phonetic category boundaries with exposure; instead, they used the
165 lexical contexts from training to re-tune their categories for particular accented features.

166 Two follow-up experiments with single-talker exposure revealed the importance of exposure-test similarity.
167 Of the five talkers included in the multi-talker exposure, two were selected. The first talker was the dissimilar
168 talker, who differed from the test talker on the three key acoustic measures associated with the word-final
169 /t/-/d/ contrast. The second talker was the similar talker, who did not differ from the test talker in the
170 means of these measures. Exposure to the similar talker generalized to the test talker, with experimental
171 exposure decreasing lexical competition between *seed*-*seat* minimal pairs compared to control exposure. By
172 contrast, exposure to the dissimilar talker did not increase performance relative to control exposure. Moreover,
173 generalization from the similar talker was as strong as generalization from multi-talker exposure (which
174 included the similar talker). Together, the results of Xie and Myers (2017) show that the correspondence
175 between exposure and test talkers at the acoustic-phonetic level is critical for understanding the observed
176 patterns of generalization across studies.

177 *1.2.2. Competing hypotheses for generalization*

178 There are two competing explanations for why multi-talker exposure may or may not provide additional
179 benefits for generalization over single-talker exposure: the exposure-to-variability hypothesis and the similarity-
180 based hypothesis.

181 The exposure-to-variability hypothesis posits that L2-accented talkers exhibit similarities in production
182 that differ from L1-accented norms (Baese-Berk et al., 2013). The exposure-to-variability hypothesis also
183 posits that among L2-accented talkers with the *same L1*, cross-language influence shifts the relations between
184 acoustic cues and phonetic categories similarly across talkers. Among L2-accented talkers with *different L1s*,
185 typological features that are unique to the L2 lead to similarly accented realizations regardless of the
186 L1. Multi-talker exposure thus allows listeners to abstract away from the peculiarities of any given talker
187 and home in on these commonalities. For example, a listener encountering one unfamiliar Spanish-accented
188 talker may not know whether their short lag VOTs are specific to that talker or characteristic of the L2
189 accent. By contrast, a listener encountering multiple unfamiliar Spanish-accented talkers at once would not
190 only see that short lag VOTs are common across the talkers, but also that other accent features (e.g., vowel
191 height) exhibit covariation with voicing (Clayards, 2017). In Bradlow and Bent (2008), multi-talker exposure
192 outperformed single-talker exposure, suggesting that training on multiple L2-accented talkers (with the same
193 L1) allowed listeners to separate the characteristic features of an accent from the idiosyncratic features of a
194 talker. The results of Bradlow and Bent (2008) support the exposure-to-variability hypothesis, but do not
195 align with the effects observed by Xie and Myers (2017) and Xie et al. (2021).

196 By contrast, the similarity-based hypothesis posits that acoustic-phonetic overlap between the exposure
197 and test talkers, rather than variability during exposure, facilitates generalization (Xie et al., 2021). In Xie
198 and Myers (2017), listeners exhibited comparable talker-independent adaptation effects after both single- and
199 multi-talker exposure to Mandarin-accented realizations of word-final /d/, which is perceptually confusable
200 with /t/ (e.g., *seed* vs. *seat*). Critically, both exposure conditions contained a Mandarin-accented talker
201 with similar word-final /d/ acoustics to the test talker. Xie et al. (2021) also demonstrated equivalent
202 generalization effects from single- and multi-talker exposure to Mandarin-accented speech. They argue that
203 exposure to multiple talkers merely increases the likelihood that listeners will encounter a cue distribution
204 that is relevant for adapting to the test talker. For example, multi-talker training on Spanish-accented speech
205 is more likely to include at least one talker with short lag VOTs than single-talker training. Bradlow and
206 Bent (2008)'s multi-talker exposure condition may have included talkers who were more similar to the test
207 talker than their single-talker exposure condition, which would confound talker similarity and exposure to
208 variability. Xie et al. (2021) addressed this confound by counterbalancing the combinations of exposure and
209 test talkers across participants. Thus, the similarity-based hypothesis may account for the differential effects
210 of single- and multi-talker exposure in perceptual adaptation studies. The present study was designed to
211 distinguish between the roles of variability and similarity in talker-independent adaptation. To summarize,
212 the similarity-based hypothesis focuses on specific cue-category mappings, while the exposure-to-variability
213 hypothesis focuses on the covariation between cues.

214 1.3. Comparing variability and similarity

215 As discussed above, the results of Bradlow and Bent (2008) support the exposure-to-variability hypothesis,
216 while those of Xie and Myers (2017) and Xie et al. (2021) support the similarity-based hypothesis. We argue
217 that the socio-indexical structure of cue-category mappings in the ideal adapter framework provides a link
218 between these two hypotheses and can account for these conflicting findings (Kleinschmidt, 2019). According
219 to the ideal adapter framework, listeners represent each cue-category mapping according to informative
220 and useful groupings. For L2-accented speech, these groupings may be structured according to each talker
221 (talker-specific) or to the L2 accent shared by the talkers (talker-independent).

222 The exposure-to-variability hypothesis argues that increasing variability during exposure increases
223 systematic covariation among relevant cues, enabling listeners to separate the idiosyncratic features of a talker
224 from the common features of an accent. To restate this hypothesis in terms of the ideal adapter framework,
225 covariation among individual exposure talkers promotes the development of a robust talker-independent
226 model (or refinement of an existing one). In turn, this talker-independent model guides adaptation to the
227 novel test talker. This exposure-to-variability perspective assumes that talker-independent models are better
228 for generalization than talker-specific models.

229 By contrast, the similarity-based hypothesis argues that increasing the similarity between the cue-category
230 mappings of the exposure and test talkers facilitates generalization. Translating this hypothesis into the ideal
231 adapter framework, listeners develop robust talker-specific models during exposure. During test, listeners
232 use the model that provides the most information about the test talker's cue distributions and most readily
233 predicts their individual cue values to guide adaptation. This similarity-based perspective assumes that
234 talker-specific models are best for generalization because they provide precise information about cue-category
235 mappings.

236 Overall, the specificity of a listener's set of generative models may explain the effects of variability and
237 similarity on generalization. The present study tests these two hypotheses in order to probe the mechanisms
238 underlying adaptation to L2-accented speech. The overarching goal is to understand how the L1 speech
239 recognition system learns the patterns of cross-language influence in L2 speech production.

240 1.4. Present study

241 How do listeners generalize their experience with L2-accented speech? On the one hand, the *structure*
242 of exposure may be the primary driver of talker-independent adaptation. Variability has been the primary
243 focus of structural inquiries (Baese-Berk et al., 2013), with the number of talkers being the key manipulation
244 (Bradlow and Bent, 2008; Bieber and Gordon-Salant, 2022; Choi and Perrachione, 2019; Xie et al., 2021). On
245 the other hand, the *content* of exposure may be the key to effective generalization. At a high level, training
246 on the relevant L2 accent tends to facilitate talker-independent adaptation (Bradlow et al., 2023; Clarke and
247 Garrett, 2004; Xie et al., 2018; Alexander and Nygaard, 2019). At a detailed level, exposure to the relevant
248 acoustic-phonetic features of an L2 accent explains generalization beyond shared L1 (Xie and Myers, 2017;
249 Sidaras et al., 2009; Reinisch and Holt, 2014). Together, a shared L1 and common acoustic-phonetic features
250 create similarity between exposure and test talkers that facilitates adaptation. The present study builds on
251 this body of work investigating generalization of L2-accented speech. The key difference between our study
252 and earlier work is the operationalization of variability and similarity. Specifically, we increased the precision
253 with which variability was implemented (c.f. Baese-Berk et al., 2013) and expanded the scope of similarity (c.f.
254 Xie and Myers, 2017) in order to better delineate the exposure-to-variability and similarity-based hypotheses.
255 The goal was to clarify some of the inconsistent findings in the literature for talker-independent adaptation
256 to L2-accented speech (Bent and Baese-Berk, 2021).

257 Regarding variability, previous studies have almost exclusively investigated this factor as a comparison
258 between single-talker exposure and multi-talker exposure (c.f. Baese-Berk et al., 2013; Bradlow et al., 2023).
259 However, multi-talker exposure reduces the amount of experience with any given talker relative to single-talker
260 exposure. From the larger perceptual adaptation literature, we know that listeners are highly sensitive to
261 acoustic-phonetic mappings and use them to adapt to new talkers (e.g., Kraljic and Samuel, 2006). We
262 also know that there is a high degree of within- and between-talker variability in both L1- and L2-accented
263 speech production (Xie and Jaeger, 2020; Wade et al., 2007). Depending on the idiosyncrasies of a given

264 exposure talker, altering the amount of experience with this talker may benefit or impede generalization.
265 This asymmetry in talker-specific exposure may explain some of the inconsistent findings in the literature
266 (Xie and Myers, 2017; Xie et al., 2021). Moreover, exposure to multiple talkers increases the sources of
267 covariation to account for, which in turn increases the difficulty of accounting for them during exposure.
268 If we assume, as the exposure-to-variability hypothesis does, that experience with covariation is critical
269 for generalization, this lack of clarity obscures how and why variability might facilitate generalization. In
270 the present study, we used a single group of talkers with the same L1 to control the type and amount of
271 experience with each talker across levels of variability.

272 Regarding similarity, previous studies have either conceptualized this factor at the accent level or at
273 the acoustic-phonetic level. When it comes to operationalizing either definition of similarity, we run into
274 the same problem as we did with multi-talker versus single-talker exposure. That is, we end up comparing
275 two (or more) groups of exposure talkers between conditions. If we assume again that listeners are highly
276 sensitive to talker-specific acoustic-phonetic features, then this design confounds talker- and accent-specific
277 effects. In other words, we cannot know whether listeners are “learning a talker or learning an accent” or
278 learning both (Xie and Myers, 2017). Here, we manipulated the stimuli listeners encounter during exposure
279 rather than the talkers in order to maintain consistency across levels of similarity. Moreover, we directly
280 crossed the factors of similarity and variability in the same design, which, to our knowledge, has not yet been
281 done. Overall, this study allows us to address the ongoing debate between the exposure-to-variability and
282 similarity-based hypotheses for generalization.

283 2. Norming study

284 2.1. Participants

285 Prior to conducting the main experiments, we recruited a separate group of participants ($N = 688$) from
286 the Penn State subject pool to norm the auditory stimuli. Participants provided implied consent in line with
287 Penn State IRB policies and were compensated 0.5 class credits after completing the experiment, which took
288 15–30 minutes. Eligible participants were between the ages of 18 and 40 years, spoke English as their first
289 and only fluent language, had normal hearing, had normal or corrected-to-normal vision, and did not have a
290 history of language-related disorders. We removed ineligible participants ($N = 34$) and those with poor data
291 quality ($N = 83$; see Section 2.5) from further analysis. This left 571 participants.

292 2.2. Stimuli

293 Stimuli were grouped by onset phoneme (e.g., *park* has the onset phoneme /p/). There were three
294 **experimental onset** groups: critical, competitor, and control. Critical onsets were voiceless stops (/p/, /t/,
295 and /k/), competitor onsets were voiced stops (/b/, /d/, and /g/), and control onsets were voiceless fricatives
296 (/f/, /s/, and /ʃ/). Across the experimental onset groups, phonemes were also grouped (roughly) according
297 to their place of articulation: labial (/p/, /b/, and /f/), alveolar (/t/, /d/, and /s/), and postalveolar/velar
298 (/k/, /g/, and /ʃ/). These will be referred to as the **cross-experimental onset** groups. There was one
299 **filler onset** group: /m/, /n/, /l/, /ɹ/, /h/, and /w/. Filler onsets were also grouped by their place of
300 articulation: nasal (/m/ or /n/), alveolar (/l/ or /ɹ/), and other back consonants (/h/ or /w/). These will
301 be referred to as the **within-filler onset** groups. Finally, each cross-experimental onset group was paired
302 with one of the within-filler onset groups according to frontness/backness: front (/p/, /b/, /f/, /m/, and
303 /n/), mid (/t/, /d/, /s/, /l/, and /ɹ/), and back (/k/, /g/, /ʃ/, /h/, and /w/). These will be referred to as
304 the **cross-condition onset** groups.

305 Stimulus selection began by downloading real words with one to four syllables, three to eight letters,
306 and one to two morphemes from the English Lexicon Project (ELP) restricted lexicon (Balota et al., 2007).
307 Within this set of items, we limited our search to words with experimental or filler onsets followed directly
308 by a vowel (e.g., *peach* was considered but *preach* was not). We removed duplicate word stems with different
309 suffixes (e.g., *paints* and *painting* were removed but *paint* was kept). We also removed any inappropriate,
310 harmful, or distracting words and word stems. The remaining set of items will be referred to as the ELP
311 pool.

312 Multisyllabic stimulus selection began by drawing real words with two to four syllables, five to eight
313 letters, and one to two morphemes from the ELP pool. We removed words with minimal pairs between the
314 critical and competitor groups (e.g., *pocket* and *docket* were removed but *socket* and *locket* were kept). Once
315 the set of options was established, real words were selected and pseudowords were created. Pseudowords were
316 created by changing the onsets of real words with filler onsets that had not been selected as potential real
317 words for the study. For experimental pseudowords, onsets were assigned according to the cross-condition
318 onset groups (e.g., *machine* became **pachine*, **bachine*, and **fachine*). For filler pseudowords, one of the
319 other five filler onsets was substituted for the existing filler onset (e.g., *medicine* became **hedicine*). In total,
320 544 multisyllabic real words and 555 multisyllabic pseudowords were normed.

321 Monosyllabic stimulus selection began by drawing real words with one syllable, three to six letters,
322 and one morpheme from the ELP pool. From this set of options, we pulled all words with critical onsets
323 that had cross-experimental minimal pairs with competitor onsets (e.g., *park*-*bark*). Both members of each
324 critical-competitor minimal pair were normed. We also selected words with filler onsets that had a minimal
325 pair with a different filler onset (e.g., *mall*-*hall*). In total, 319 monosyllabic real words were normed.

326 Norming took place in three waves of testing. Overall, there were nine experimental lists of 180 to 540
327 items: three in the first wave, two in the second, and four in the third. Lists always contained an equal
328 number of real words and pseudowords, and items from each experimental onset groups were always presented
329 in separate lists.

330 2.3. Recording

331 One female L1-accented English talker from the US (Talker 0; the first author) recorded all items. Audio
332 was captured with a head-worn condenser microphone (Shure SM35-XLR) connected to an audio interface
333 (Sound Devices USBPre2) and recorded with Praat (Broersma and Weenink, 2021) in mono at 44.1 kHz in a
334 sound-attenuated booth. After annotating the recordings and extracting the individual sound files, stimuli
335 were normalized to 70 dB and had 50 ms of silence added to the beginning and end.

336 2.4. Task and procedure

337 The study was conducted online using Pavlovia. At the beginning of the study, participants responded to
338 a yes/no question for each eligibility criterion. Ineligible participants were not able to complete the study.

339 The study featured an auditory lexical decision task. Participants were randomly assigned to an
340 experimental list. On each trial, participants indicated whether an auditory stimulus was a real English
341 word or not by pressing the *d* or *k* key on their keyboard. One of the two response-key relations—real-*d* or
342 real-*k*—was assigned randomly to each participant. The total number of trials varied by list.

343 2.5. Analysis

344 Analyses were conducted with R version 4.2.2 using the *stats* package (R Core Team, 2022). Each wave of
345 testing was analyzed separately. The experimental lists within each wave were combined for analysis. We first
346 conducted one-sample t-tests comparing each participant's accuracy to chance (50%) to check for data quality.
347 Data from participants whose performance was indistinguishable from chance were removed from further
348 analysis (see Section 2.1). Next, we conducted one-sample t-tests comparing the accuracy on each item to
349 chance (50%). Items that did not yield above-chance accuracy were removed from further consideration.
350 All three variants of each experimental pseudoword and both members of each cross-experimental minimal
351 pair needed to have above-chance accuracy in order to be considered for final selection. The set of items
352 that passed the norming phase will be referred to as the selection pool. For details on the stimuli that were
353 selected for Experiment 1, see Section 3.1.5.

354 3. Experiment 1: Investigating differences in generalization from exposure to Spanish-accented 355 stops versus fricatives

356 XX It would be good to briefly describe the main research question and theory-guided predictions.

357 3.1. Methods

358 We used an exposure-test design to understand how different kinds of experience with Spanish-accented
359 speech change listeners' VOT-stop mappings. The exposure phase established the comparison between the
360 similarity-based and exposure-to-variability hypotheses of talker-independent adaptation. The test phase
361 assessed the effects of each type of exposure on perception.

362 3.1.1. Design

363 Participants were exposed to multisyllabic real words and pseudowords before being tested on monosyllabic
364 real words. Participants heard these items produced by four Spanish-accented talkers: three during exposure
365 and one during test. During exposure, multisyllabic items without onset competitors provided disambiguating
366 lexical contexts for categorizing Spanish-accented onsets. For example, consider the real word *pencil* and
367 the pseudoword **pachine*. In both cases, the onset may be interpreted as /p/ or /b/. This ambiguity
368 does not affect whether **pachine* is perceived as a real word or not, since both **pachine* and **bachine* are
369 pseudowords. However, resolving this ambiguity is necessary for distinguishing between the real word *pencil*
370 and pseudoword **bencil*. Critically, participants hearing ambiguous real words like *pencil* would never hear
371 unambiguous real words like *beehive*. Thus, in the context of the exposure task, participants should learn to
372 perceive the ambiguous short lag VOTs as /p/ rather than as /b/. During test, monosyllabic items with
373 onset competitors created ambiguous lexical contexts in which to assess learning from exposure.

374 **Exposure similarity** was operationalized as the relation between the experimental onsets encountered
375 during exposure and the critical onsets encountered during test. There were three levels of Similarity:
376 Direct, Indirect, and Control. Each level refers to the type of information participants received about
377 Spanish-accented voiceless stops. In Direct conditions, participants were exposed directly to critical onsets
378 (e.g., *peanut* and **pachine*). In Indirect conditions, participants were exposed to competitor onsets (e.g.,
379 *beehive* and **bachine*), thereby gaining experience with the shifted VOT continuum that they would encounter
380 during test. In Control conditions, participants were exposed to control onsets (e.g., *football* and **fachine*),
381 thereby gaining general experience with Spanish-accented speech but not with stop VOTs. This design
382 allowed the talkers to remain the same across the three levels of Similarity.

383 **Exposure variability** was operationalized as the relation between onset phonemes and exposure talkers.
384 There were two levels of Variability: Invariant and Variant. Each level refers to the type of experience
385 with each talker. This is illustrated in Figure 2. In Invariant conditions, listeners heard each of the three
386 exposure talkers produce one onset out of the three in an experimental group (e.g., Direct-Invariant: Talker A
387 produced *peanut*, *palace*, and *pianist*; Talker B produced *terminal*, *tiara*, and *textile*; and Talker C produced
388 *kingdom*, *counter*, and *kayak*). This means that all of the words with a given onset were produced by only
389 one talker. In Variant conditions, listeners heard each of the three exposure talkers produce all three onsets
390 in an experimental group (e.g., Direct-Variant: Talker A produced *peanut*, *terminal*, and *kingdom*; Talker B
391 produced *palace*, *tiara*, and *counter*; and Talker C produced *pianist*, *textile*, and *kayak*). This means that
392 one third of the words with a given onset were produced by each talker. This design allowed the items and
393 talkers to remain the same between levels of Variability.

394 The two exposure factors of Similarity (Direct, Indirect, Unrelated) and Variability (Variant, Invariant)
395 were manipulated between participants, so each participant was assigned to one of the six combinations
396 of Similarity and Variability. Target type was manipulated within participants. There were three levels of
397 Target: Identity, Competitor, and Unrelated (See Figure 2). Each level refers to the type of visual target
398 that followed each auditory prime. Identity targets exactly matched the auditory primes (e.g., *park-park*).
399 Competitor targets were the minimal pairs of the auditory primes (e.g., *park-bark*). Unrelated targets only
400 shared vowels with the auditory primes (e.g., *park-wand*). Differential performance on these three conditions
401 indexed learning from exposure.

402 3.1.2. Participants

403 We recruited 296 participants through the online platform Prolific. Participants provided implied consent
404 in line with Penn State IRB policies and were compensated \$6 after completing the experiment, which took
405 approximately 30 minutes. The experiment was available to individuals whose Prolific user profiles aligned

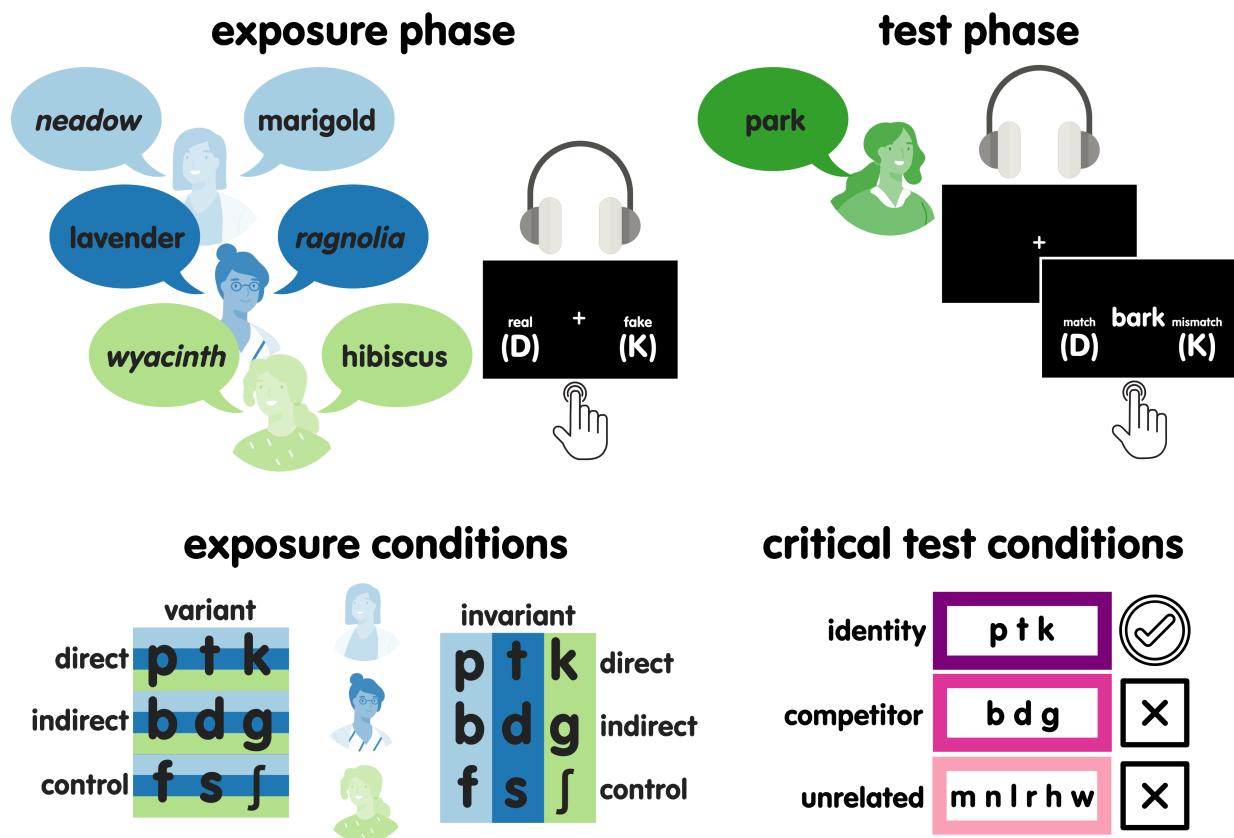


Figure 2: Experiment 1 design. The exposure phase panel illustrates the task with filler real words and pseudowords.

Table 1: Talker background information.

Talker	Region	Age of English acquisition	Age of arrival in US	Age at time of recording
1	Veracruz, Mexico	10	18	26
2	Quintana Roo, Mexico	13	26	28
3	Puebla, Mexico	5	24	27
4	Mexico City, Mexico	5	28	31

406 with the following eligibility criteria: between 18 and 40 years of age, located in the US at the time of
 407 the study, English as their first and only fluent language, normal hearing, normal or corrected-to-normal
 408 vision, and without a history of language-related disorders. The Prolific user profiles for each participant also
 409 included sex (two options, select one) and race/ethnicity information (four options, select all that apply).

410 The eligibility criteria were cross-checked with the responses to a post-experiment questionnaire (see
 411 Section 3.1.6), and any ineligible participants was removed from further analysis ($N = 9$). We also removed
 412 participants with any knowledge of Spanish, with self-rated proficiency greater than or equal to 3/5 in any
 413 languages other than English, with self-rated proficiency less than 5/5 in English, or whose place of origin
 414 was not the US ($N = 15$). This left 287 eligible participants.

415 Finally, we removed participants with poor data quality from further analysis ($N = 15$). Poor data quality
 416 was defined as: accuracy statistically indistinguishable from or significantly below chance on experimental
 417 real words in the exposure task, zero correct experimental trials in any of the three conditions in the test
 418 task, or fewer than 50% of experimental trials with correct responses and reaction times between 50 and
 419 2500 ms in either task (see Section 3.1.6 for details). This left 257 participants for analysis (Age: $M = 31$,
 420 $SD = 6$, Min = 18, Max = 40; Sex: Female = 121, Male = 135, Prefer not to say = 1; Race: Asian = 5,
 421 Black = 43, Multiple selected = 21, Other = 5, White = 182, Not provided = 1).

422 We recruited 46 additional participants through Prolific to complete the experiment without the exposure
 423 phase. Removing ineligible participants left 38 participants. Removing participants with low data quality
 424 left 37 participants for analysis (Age: $M = 31$, $SD = 5$, Min = 21, Max = 40; Sex: Female = 19, Male = 18;
 425 Race: Asian = 2, Multiple selected = 5, Other = 2, White = 28). All aspects of recruitment were the same
 426 for this group as for the main group of participants.

427 3.1.3. *Talkers*

428 Seven female Spanish-accented English talkers from Latin America and Spain were recruited to be talkers
 429 for the experiment. Talkers were compensated with one \$20 Amazon giftcard per recording session (30-60
 430 minutes each; 1-2 total). Recording followed the procedure described in Section 2.3.

431 Out of the seven talkers, we selected four from Mexico (Talkers 1-4) to control for country-level dialectal
 432 variation in the L1. Each talker was from a different region of Mexico and reported living in this region for the
 433 majority of their lives before moving to the US for college or graduate school. All four talkers reported that
 434 they grew up speaking Spanish with their caregivers and began acquiring English in traditional classroom
 435 settings. Language background information for each talker is provided in Table 1.

436 3.1.4. *Stimuli*

437 In total, there were 936 auditory items and 648 visual items across tasks. The exposure task featured
 438 multisyllabic real words and pseudowords with experimental and filler onsets from the selection pool. The
 439 final set of exposure items included 360 real words (24 per onset) and 360 pseudowords (24 per onset).
 440 The test task featured two types of stimuli: auditory primes and visual targets. Auditory primes were
 441 monosyllabic real words with critical or filler onsets from the selection pool. The final set of auditory primes
 442 included 216 real words (24 per onset). The final set of visual targets included 648 real words (3 per prime).

443 During auditory stimulus selection, we considered a number of parameters. From the ELP, we included
 444 word length, orthographic neighborhood density (OND), phonological neighborhood density (PND), and US
 445 Zipf frequency. From Brysbaert et al. (2019), we included percent known and prevalence. We also included
 446 the position of lexical stress, calculated from the pronunciation information provided by the ELP, and mean

Table 2: Mean and standard deviation of stimulus parameters by condition and task.

Phase	Word type	Condition	N	Length	LDT	Lexical stress	OND	PND	Frequency	Percent known	Prevalence
Exposure	Real word	Direct	72	6.85 (0.93)	0.89 (0.09)	1.11 (0.32)	0.56 (1.40)	1.12 (2.13)	3.38 (0.74)	0.98 (0.02)	2.12 (0.29)
		Indirect	72	6.82 (0.98)	0.89 (0.08)	1.21 (0.41)	0.62 (1.34)	1.14 (2.04)	3.45 (0.72)	0.99 (0.02)	2.16 (0.29)
		Control	72	6.93 (0.95)	0.87 (0.09)	1.15 (0.36)	0.67 (1.17)	1.76 (2.76)	3.39 (0.77)	0.98 (0.03)	2.11 (0.38)
		Filler	144	6.83 (0.95)	0.87 (0.10)	1.12 (0.32)	0.78 (1.52)	1.49 (2.66)	3.42 (0.77)	0.98 (0.03)	2.14 (0.35)
	Pseudoword	Direct	72	7.08 (0.87)	0.89 (0.06)	1.22 (0.42)					
		Indirect	72	7.08 (0.87)	0.88 (0.07)	1.22 (0.42)					
		Control	72	7.08 (0.87)	0.88 (0.06)	1.22 (0.42)					
		Filler	144	6.93 (1.01)	0.88 (0.05)	1.14 (0.35)					
Test	Real word	Critical prime	72	3.88 (0.63)	0.88 (0.07)		11.36 (5.21)	21.85 (7.96)	4.15 (0.94)	0.99 (0.04)	2.22 (0.34)
		Competitor pair	72	3.90 (0.65)	0.88 (0.08)		10.39 (5.02)	21.15 (7.12)	4.09 (1.11)	0.98 (0.05)	2.10 (0.42)
		Filler prime	144	3.90 (0.58)	0.90 (0.08)		10.49 (4.70)	22.27 (8.89)	4.15 (0.97)	0.99 (0.02)	2.22 (0.32)

Competitor pairs were presented as visual targets only

lexical decision accuracy from norming (LDT). The final set of exposure real words was chosen such that the four groups of stimuli—Direct, Indirect, Control, and Filler—were matched on each of these parameters. We conducted two-sample t-tests comparing each group on each parameter to ensure that they were not significantly different ($ps > .05$). The final set of auditory primes for the test task was selected in a similar way. We conducted two-sample t-tests comparing the Critical and Filler primes on every parameter except for lexical stress, which was not relevant for monosyllabic words ($ps > .05$). We also conducted paired t-tests between Critical primes and their Competitor pairs ($ps > .05$). The final set of exposure pseudowords was chosen such that the four groups of stimuli were matched on length, position of lexical stress, and mean lexical decision accuracy. The other parameters were either not available (OND and PND) or not relevant (US Zipf frequency, percent known, prevalence) for the pseudowords. The relevant parameters for each group of stimuli are shown in Table 2. Descriptive statistics for each talker’s real word VOTs by onset are provided in Table 3 (see Section 6).

An additional set of six real words and six pseudowords with filler onsets (one per onset per word type) was selected for practice. Six additional primes with filler onsets (one per onset) were also selected for practice.

Visual targets in the test task were monosyllabic real words with critical, competitor, or filler onsets. Each auditory prime had three visual targets: one was the prime itself, one was its minimal pair, and one was unrelated. The first two were determined by the prime, but the third needed to be selected separately. Each unrelated target had a filler onset and was chosen to have the same vowel as the prime but a different onset and offset. These items were primarily from the selection pool, but some were from the larger ELP pool. Two additional unrelated targets were also selected for practice.

3.1.5. Experimental lists

Combining Similarity and Variability created six between-subjects conditions: Direct-Variant, Direct-Invariant, Indirect-Variant, Indirect-Invariant, Control-Variant, and Control-Invariant. In addition, one group of participants did not receive any exposure prior to test; this will be referred to as the Test-only condition. The 288 filler items were the same across conditions and evenly divided by word type—real word and pseudoword—as well as by onset—/m/, /n/, /l/, /ɪ/, /h/, and /w/. The 144 experimental items were evenly divided by onset and word type, with onset differing by level of Similarity: Direct (critical: /p/, /t/, and /k/), Indirect (competitor: /b/, /d/, and /g/), and Control (control: /f/, /s/, and /ʃ/). As a result, there were 24 items per onset per word type.

The assignment of talkers to items differed by level of Variability. In Invariant conditions, talkers were assigned by cross-condition onset group: front (/p/, /b/, /f/, /m/, and /n/), mid (/t/, /d/, /s/, /l/, and /ɪ/), and back/other (/k/, /g/, /ʃ/, /h/, and /w/). In other words, one talker was assigned to front onsets, one to mid, and one to back/other. Within a given level of Similarity, this means that each talker produced items with one out of the three experimental onsets and two out of the six filler onsets. For Variant conditions, one third of each of the items of a given onset and word type (8) was randomly assigned to one of three sets: set 1, set 2, or set 3. Each set was assigned to one talker.

To counterbalance which talkers were assigned to which items across participants, each Variant item set

485 was paired with one Invariant onset group to create three assignment groups: group 1 (front with set 1),
486 group 2 (mid with set 2), and group 3 (back/other with set 3). In addition, we also counterbalanced the
487 talkers between the exposure and test phases. To do this, we added group 4 (test) and rotated the four
488 talkers across these four assignment groups in a Latin square design. Overall, this resulted in 24 experimental
489 lists for the exposure phase, one for each combination of exposure condition (6) and talker assignment (4).

490 Regardless of exposure condition, all participants heard the same 216 auditory primes during test. One
491 third (72) were critical items divided evenly by onset—/p/, /t/, and /k/—and two thirds (144) were filler
492 items divided evenly by onset—/m/, /n/, /l/, /x/, /h/, and /w/. Each onset was divided evenly by target
493 type—Identity, Competitor, and Unrelated. This resulted in eight items per onset per target type Three
494 experimental lists were created to counterbalance the combinations of auditory prime and visual target
495 across participants. Since the test talker was also counterbalanced across participants, this resulted in 12
496 experimental lists for the test phase, one for each combination of prime-target pair (3) and talker assignment
497 (4).

498 Exposure practice was presented in Talker 0’s voice for all participants. Test practice was presented
499 according to the participant’s talker assignment and Variability condition. Participants in the Test-only
500 condition completed the test practice in Talker 0’s voice.

501 3.1.6. Tasks

502 The experiment included a headphone check, exposure task, test task, and post-experiment questionnaire.
503 All aspects of the experiment were conducted online using Pavlovia. The headphone check, exposure task,
504 and test task were created in PsychoPy Builder (Peirce et al., 2019). The post-experiment questionnaire was
505 built using Pavlovia’s survey platform.

506 The headphone check followed the anti-phase tone test procedure from Woods et al. (2017). On each trial,
507 participants listened to three pure tones, one of which was out of phase with the other two. Participants
508 indicated which of the three tones was the quietest by selecting the appropriate button on the screen: tone 1,
509 tone 2, or tone 3. There were two trials per response for a total of six trials, which were presented in random
510 order. Participants using well-functioning headphones should have easily perceived the anti-phase tone as
511 the quietest, while those using loudspeakers should not. Participants completed the task at most two times.

512 The exposure phase featured the auditory lexical decision task from Xie and Myers (2017). On each trial,
513 participants indicated whether an auditory stimulus was a real English word or not by pressing the *d* or *k*
514 key on their keyboard. One of the two response-key relations—real-*d* or real-*k*—was assigned randomly to
515 each participant. Participants completed 12 practice trials followed by 432 main trials presented in random
516 order. Half of the practice trials (6) and half of the main trials (216) required real word responses.

517 The test phase featured a cross-modal matching task adapted from the primed cross-modal lexical decision
518 task in Xie and Myers (2017). On each trial, participants first heard a real word (auditory prime) and then
519 saw a real word written on the screen (visual target). They indicated whether the visual target matched the
520 auditory prime or not by pressing the *d* or *k* key on their keyboard. The assignment of responses to keys was
521 carried over from the exposure phase, with *match* responses mapped to the same key as *real word* responses
522 in the exposure task. Participants completed six practice trials followed by 216 main trials presented in
523 random order. Half of the practice trials (3) and one third of the main trials (72) required match responses.

524 The post-experiment questionnaire included two sets of items. The first set of items related to the
525 talker from the test task. The second set of items included questions about the participant’s own language
526 background and demographics. This second set of items was used to confirm the participant’s eligibility as
527 described in Section 3.1.2. All post-experiment questionnaire items are described in detail in the supplementary
528 materials.

529 3.1.7. Procedure

530 Eligible participants accessed the experiment through Prolific. Once a participant began the experiment,
531 they were randomly assigned to one of the experimental lists. First, they performed the headphone check. If
532 the participant achieved fewer than five correct trials out of six, they completed the task again. Participants
533 who failed the headphone check a second time were not allowed to continue; instead, they were redirected to
534 Prolific and asked to return their submission. After passing the headphone check, participants continued to

535 the exposure task and completed the practice (participants in the Test-only condition continued straight
536 to the test task). If a participant scored 50% or lower on the exposure practice, they were not allowed to
537 continue; instead, they were redirected to Prolific and asked to return their submission. After successfully
538 completing the practice session, the participant performed the exposure task. Next, the participant continued
539 to the test task and completed the practice. Regardless of their performance on the practice, the participant
540 performed the test task. Immediately following the test task, participants continued to the post-experiment
541 questionnaire. Once the questionnaire was complete, the participant was redirected to Prolific and received
542 compensation.

543 3.1.8. Analysis approach

544 Data processing and analysis were conducted with R version 4.2.2 (R Core Team, 2022). Filler items
545 were not included in any of the analyses.

546 Prior to analyzing the exposure task data, we removed responses with RTs less than 50 ms ($N = 31$;
547 0.03%). RT was calculated from the onset of the word. For the RT analyses, we filtered for correct responses.
548 We used the *robustbase* package to calculate adjusted boxplot statistics for skewed distributions (such as
549 reaction time), which were used to detect responses with reaction times outside the upper and lower fences
550 (Hubert and Vandervieren, 2008). These outlier responses were removed from further analysis ($N = 4630$;
551 4.17%). RTs were then inverse-transformed (-1000/RT) for analysis.

552 Prior to analyzing the test task data, we removed responses with RTs less than 50 ms ($N = 4$; 0.02%). RT
553 was calculated from the presentation of the visual target. RT analyses were restricted to correct responses.
554 We detected and removed outliers by target type according to the adjusted boxplot statistics ($N = 898$;
555 4.24%). Inverse RTs were used for modeling as in the exposure task analyses.

556 Mixed-effects models were fitted to trial-level data with the *lme4* package (Bates et al., 2015). Generalized
557 linear mixed-effects models with a binomial family function were used to analyze binary accuracy (1,0).
558 Linear mixed-effects models were used to analyze inverse RT. Model fitting began with the full random effects
559 structure that was relevant to the task (see below). In the case of non-convergence, singularity, or correlations
560 above 0.95, random slopes were successively removed, such that the final model for each analysis reflected the
561 maximally-supported structure (Barr et al., 2013). Type-III analysis-of-deviance tables were calculated and
562 Wald chi-square tests conducted with the *car* package (Fox and Weisberg, 2019). Estimated marginal means
563 were calculated and pairwise comparisons conducted with the *emmeans* package (Lenth, 2022). Pairwise
564 *p*-values were adjusted with the Hommel method to control the family-wise error rate (Blakesley et al., 2009).

565 Exposure task analyses modeled the effects of Variability, Similarity, Word type, and their interactions on
566 accuracy and RT. The two levels of Variability and Word type were sum contrast-coded. The three levels of
567 Similarity were Helmert contrast-coded. Scaled VOT, trial, and word frequency were included as covariates.
568 Item and participant were included as random intercepts. By-participant random slopes were included for
569 Variability, Similarity, and their interaction.

570 For the test task, we conducted two separate sets of analyses. The first set of analyses modeled the effects
571 of Exposure, Target, and their interaction on accuracy and RT. Exposure had seven levels: one for each of
572 the six exposure conditions and one for the Test-only condition. This factor was simple contrast-coded such
573 that the Test-only condition was the reference level. The three levels of Target were Helmert contrast-coded.
574 Scaled VOT, scaled trial, and the interaction between scaled prime frequency and scaled target frequency
575 were included as covariates. Participant and the interaction between auditory prime and visual target were
576 included as random intercepts.

577 The second set of analyses modeled the effects of Variability, Similarity, Target, and their interactions
578 on accuracy and RT. The coding scheme for each variable was the same as in the first test analysis. The
579 continuous predictors and random intercepts were also the same. By-participant random slopes were included
580 for Variability, Similarity, and their interaction.

581 If the first set of analyses did not reveal differences between the exposure and Test-only conditions, the
582 second set of analyses was not conducted. Only significant effects involving the exposure conditions will be
583 discussed. All results are available on GitHub (XX).

584 *3.1.9. Predictions*

585 The exposure-to-variability and similarity-based hypotheses make different predictions about the effects
586 of each type of exposure on test. The similarity-based hypothesis predicts a main effect of Similarity, with
587 test performance increasing from Control to Indirect to Direct exposure. This prediction follows directly
588 from the ideal adapter framework, such that both talker-specific and talker-independent generative models of
589 VOT-stop distributions can generalize to a test talker; what matters is how similar the exposure model is to
590 the test distribution.

591 By contrast, the exposure-to-variability hypothesis predicts a main effect of Variability, with better
592 test performance after Variant exposure than after Invariant exposure. Under this hypothesis, exposure to
593 covariation between cues and categories encourages the formation of talker-independent generative models.
594 This implies that such models have more utility for categorizing critical onsets than any talker-specific models
595 (Kleinschmidt, 2019). The exposure-to-variability hypothesis does not predict effects of Similarity on test
596 performance.

597 *3.2. Results*

598 *3.2.1. Exposure*

599 For accuracy, there was a significant three-way interaction between Variability, Similarity, and Word type
600 ($\chi^2(2, N = 3) = 15.90, p < .001$), as well as a significant two-way interaction between Variability and Word
601 type ($\chi^2(1, N = 2) = 11.39, p < .001$). For RT, we also observed a significant three-way interaction between
602 Variability, Similarity, and Word type ($\chi^2(2, N = 3) = 30.46, p < .001$). In addition, the main effect of
603 Similarity was significant ($\chi^2(2, N = 3) = 36.93, p < .001$). To follow up on these effects, we conducted two
604 sets of pairwise comparisons: one between levels of Similarity within each combination of Variability and
605 Word type and another between levels of Variability within each combination of Similarity and Word type.

606 All of the comparisons for accuracy are shown in the left column of Figure 3. Within Variant exposure,
607 pseudoword accuracy was higher for Control exposure ($M = 0.93, 95\% \text{ CI } [0.90, 0.95]$) compared to Indirect
608 exposure ($M = 0.88, 95\% \text{ CI } [0.84, 0.91]$; $z = 2.52, p = .035$). Within Control exposure, real word accuracy
609 was higher for Invariant exposure ($M = 0.86, 95\% \text{ CI } [0.83, 0.89]$) compared to Variant exposure ($M = 0.83,$
610 $95\% \text{ CI } [0.78, 0.87]$; $z = 2.59, p = .010$).

611 All of the comparisons for RT are shown in the right column of Figure 3. Comparing levels of Similarity
612 revealed overall slower RTs for Control exposure. Within Variant exposure, pseudoword RTs were slower for
613 Control exposure ($M = 1361, 95\% \text{ CI } [1310, 1417]$) compared to Direct exposure ($M = 1242, 95\% \text{ CI } [1195,$
614 $1292]$; $z = 3.25, p = .003$). In addition, real word RTs were slower for Control exposure ($M = 1254, 95\% \text{ CI }$
615 $[1211, 1301]$) compared to both Direct ($M = 1110, 95\% \text{ CI } [1074, 1148]$; $z = 3.44, p = .002$) and Indirect (M
616 $= 1135, 95\% \text{ CI } [1097, 1176]$; $z = 2.97, p = .006$) exposure. Within Invariant exposure, pseudoword RTs were
617 slower for Control exposure ($M = 1446, 95\% \text{ CI } [1389, 1508]$) compared to both Direct ($M = 1267, 95\% \text{ CI }$
618 $[1220, 1317]$; $z = 4.66, p < .001$) and Indirect ($M = 1277, 95\% \text{ CI } [1229, 1329]$; $z = 4.29, p < .001$) exposure.
619 Real word RTs were also slower for Control exposure ($M = 1254, 95\% \text{ CI } [1211, 1301]$) compared to both
620 Direct ($M = 1110, 95\% \text{ CI } [1074, 1148]$; $z = 4.93, p < .001$) and Indirect ($M = 1135, 95\% \text{ CI } [1097, 1176]$; z
621 $= 3.91, p < .001$) exposure. Comparing levels of Variability revealed differences within Control exposure,
622 such that pseudoword RTs were slower for Invariant exposure than for Variant exposure ($z = 2.33, p = .020$).

623 *3.2.2. Test: Comparison to the Test-only group*

624 For accuracy, we observed a main effect of Exposure ($\chi^2(6, N = 7) = 13.28, p = .039$); however, pairwise
625 comparisons did not reveal significant differences between the Test-only group and any of the exposure groups
626 when averaging across levels of Target ($ps > .05$). To further investigate the effect of Exposure, we conducted
627 pairwise comparisons between the Test-only group and the exposure groups within each level of Target, but
628 did not observe any significant differences ($ps > .05$). We also conducted pairwise comparisons between the
629 exposure groups; again, there were no significant differences ($ps > .05$). By-participant means and estimated
630 marginal means for each level of Exposure and Target are shown in Figure 4.

631 For RT, we observed an interaction between Exposure and Target ($\chi^2(12, N = 13) = 24.59, p = .017$);
632 however, none of the pairwise comparisons between the Test-only group and any of the exposure groups

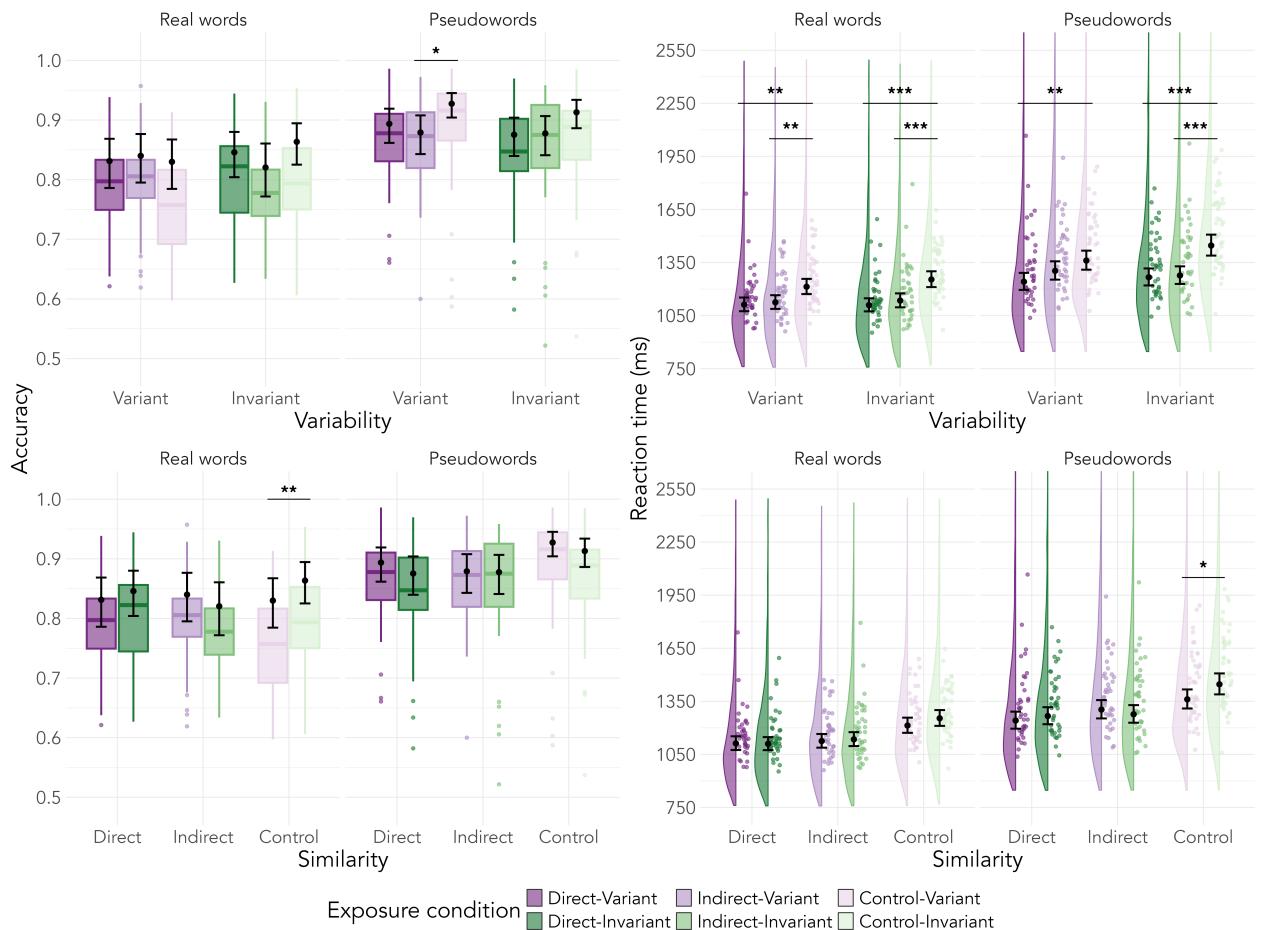


Figure 3: Experiment 1 exposure task performance. Left column presents boxplots with mean accuracy by participant. Right column presents half violin plots with reaction times for correct responses and dot plots with mean reaction times for correct responses by participant. Plots are overlaid with estimated marginal means and 95% confidence intervals. Asterisks indicate significance levels from pairwise comparisons: *** $< .001$, ** $< .01$, * $< .05$.

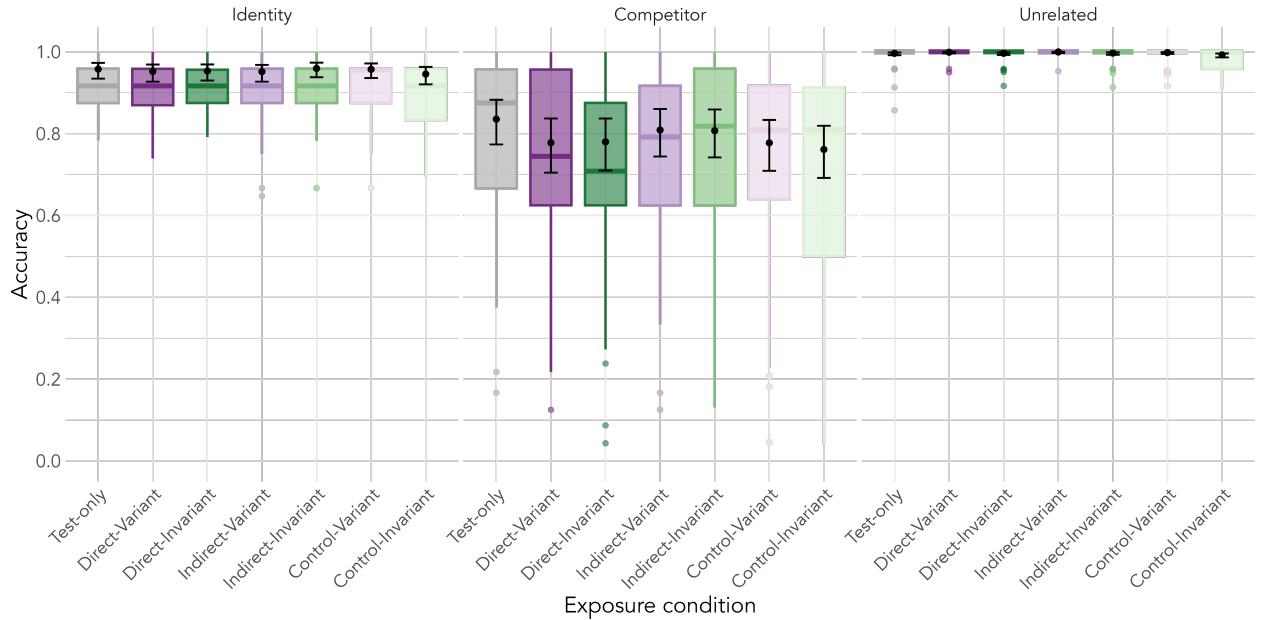


Figure 4: Experiment 1 test task accuracy. Boxplots present by-participant means for each Target type. Plots are overlaid with estimated marginal means and 95% confidence intervals.

633 within each level of Target was significant ($ps > .05$). We also conducted pairwise comparisons between each
 634 of the exposure groups, but did not observe any significant differences ($ps > .05$).

635 3.3. Discussion

636 Experiment 1 compared the effects of exposure to voiceless stops (Direct), voiced stops (Indirect), or
 637 voiceless fricatives (Control) on adaptation to Spanish-accented speech. These three levels of Similarity were
 638 crossed with two levels of Variability, such that each talker either produced exemplars of all onsets (Variant)
 639 or exemplars of one particular set of onsets (Invariant). Performance on the exposure task did not differ
 640 between the two critical levels of Similarity—Direct and Indirect—in either accuracy or RT. Between levels
 641 of Variability, performance only differed within Control exposure. The results of the exposure task show that
 642 participants were equally able to distinguish real words from pseudowords in the context of Spanish-accented
 643 speech regardless of exposure condition.

644 However, the exposure task did not benefit performance on the test task. We did not observe differences
 645 between the Test-only group, which did not receive any exposure to Spanish-accented speech prior to the
 646 test phase, and any of the exposure conditions. Under the similarity-based hypothesis, we expected Direct
 647 exposure to increase accuracy on the matching task, particularly for Competitor targets (*park-bark*). Having
 648 direct experience with the mappings between VOT and /p/, /t/, and /k/ in supporting lexical contexts was
 649 expected to decrease activation of Competitor targets and increase activation of Identity targets based on
 650 Xie and Myers (2017). Here, we observed neither.

651 A potential difference between our design and that of Xie and Myers (2017) was in the construction of
 652 the pseudowords. Specifically, we included pseudowords with the experimental onsets that participants were
 653 adapting to during exposure. For example, Direct exposure included both the real word *peanut* and the
 654 pseudoword **pachine*. While we reasoned in Section 3.1.1 that pseudowords like **pachine* should not disrupt
 655 perceptual learning of the relevant VOT-/p/ mapping, it is possible that exposure to these mappings in
 656 non-lexical contexts disrupted the benefits of exposure to these mappings in disambiguating lexical contexts.
 657 To address this issue, we removed these pseudowords in Experiment 2. Instead, participants in the Direct
 658 and Indirect groups were both exposed to pseudowords with control onsets like **fachine*. Including control
 659 onsets within the Direct and Indirect groups lessened the need for separate Control groups. To reduce the

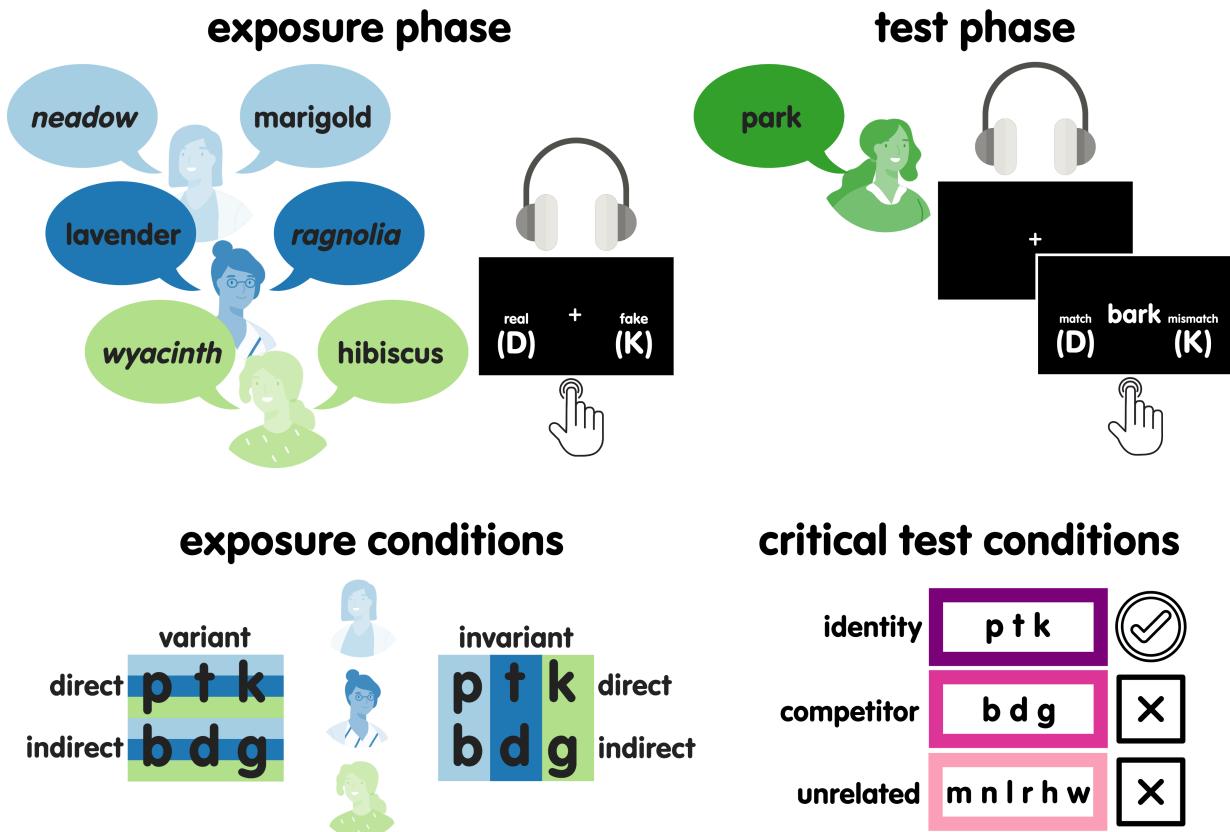


Figure 5: Experiment 2 design.

complexity of the design and home in on the effects of Direct versus Indirect exposure, the Control level of Similarity was not included in Experiment 2.

4. Experiment 2: Investigating the specificity of generalization from exposure to Spanish-accented stops

4.1. Methods

This experiment was the same as Experiment 1, with the exception of the Similarity factor.

4.1.1. Design

The Control level was removed, leaving just two levels of Similarity: Direct and Indirect. Recall that the pseudowords in the Direct and Indirect conditions had the same onsets as the experimental real words (critical and competitor, respectively). It is possible that this design disrupted adaptation to Spanish-accented VOTs by associating them with pseudowords. To improve learning, we replaced them with pseudowords with control onsets.

4.1.2. Participants

We recruited 192 participants through Prolific, none of whom had participated in Experiment 1. All aspects of recruitment and criteria for data exclusion were the same as in Experiment 1. After removing ineligible participants, 160 remained. After removing participants with poor data quality, 154 remained for

analysis (Age: $M = 31$, $SD = 6$, Min = 18, Max = 40; Sex: Female = 65, Male = 89; Race: Asian = 5, Black = 17, Multiple selected = 12, Other = 4, White = 116). Data from the participants who completed the experiment without the exposure phase in Experiment 1 were used for comparison.

4.1.3. Materials and procedure

The talkers, stimuli, tasks, and procedure were the same as those in Experiment 1. Removing the Control level of Similarity from the exposure phase left four between-subjects conditions: Direct-Variant, Direct-Invariant, Indirect-Variant, and Indirect-Invariant. The 288 filler items remained the same. The 72 experimental real words also remained the same. The 72 experimental pseudowords had control onsets regardless of Similarity. Talker assignment and counterbalancing was the same. This resulted in 16 experimental lists for the exposure phase, one for each combination of exposure condition (4) and talker assignment (4). All aspects of the test phase remained the same.

4.1.4. Analysis and predictions

The data processing, model fitting, and analysis approaches were the same as in Experiment 1; the only change was implementing sum contrasts for the two levels of Similarity. Prior to analyzing exposure task performance for real words with experimental onsets, we removed responses with RTs less than 50 ms ($N = 64$; 0.10%). We then detected and removed outliers ($N = 2917$; 4.39%). Prior to analyzing test task performance for critical prime-target pairs, we removed responses with RTs less than 50 ms ($N = 7$; 0.05%). We then detected and removed outliers ($N = 679$; 4.94%).

4.2. Results

4.2.1. Exposure

For accuracy, we observed significant interactions among Variability, Similarity, and Word type ($\chi^2(1, N = 2) = 10.78, p = .001$) and between Variability and Word type ($\chi^2(1, N = 2) = 4.64, p = .031$). To investigate these effects, we conducted pairwise comparisons between levels of Similarity within each combination of Variability and Word type and between levels of Variability within each combination of Similarity and Word type. These comparisons are shown in the left column of Figure 6. Within Variant exposure, pseudoword accuracy was higher for Direct exposure ($M = 0.91$, 95% CI [0.88, 0.93]) compared to Indirect exposure ($M = 0.94$, 95% CI [0.92, 0.95]; $z = 2.80, p = .005$).

For RT, the interaction between Variability and Word type was significant ($\chi^2(1, N = 2) = 10.66, p = .001$). Pairwise comparisons between Variant and Invariant exposure within each level of Word type were not significant ($ps > .05$). To further explore the effects of Variability and Similarity on RTs, we also conducted the same pairwise comparisons as we did for accuracy; however, there were no differences in this case ($ps > .05$). These comparisons are shown in the right column of Figure 6.

4.2.2. Test: Comparison to the Test-only group

For accuracy, we observed a significant interaction between Exposure and Target ($\chi^2(8, N = 9) = 20.20, p = .010$). Pairwise comparisons within each level of Target did not reveal significant differences between the Test-only group and any of the exposure groups ($ps > .05$). To investigate the source of the Exposure-Target interaction, we conducted pairwise comparisons between each of the exposure groups. Direct-Invariant exposure yielded significantly higher accuracy on Competitor targets ($M = 0.87$, 95% CI [0.82, 0.91]) than Direct-Variant exposure ($z = 2.99, p = .017$), Indirect-Invariant exposure ($z = 2.68, p = .030$), or Indirect-Variant exposure ($z = 2.55, p = .043$). By-participant means and estimated marginal means are shown in Figure 7.

For RT, we also observed a significant interaction between Exposure and Target ($\chi^2(8, N = 9) = 41.53, p < .001$). Pairwise comparisons within each level of Target did not reveal significant differences between the Test-only group and any of the exposure groups ($ps > .05$). Pairwise comparisons between the exposure groups revealed significantly slower RTs on Identity targets for Direct-Invariant exposure ($M = 0.95$, 95% CI [0.92, 0.97]) than for Indirect-Variant exposure ($M = 0.95$, 95% CI [0.93, 0.97]; $z = 2.79, p = .032$).

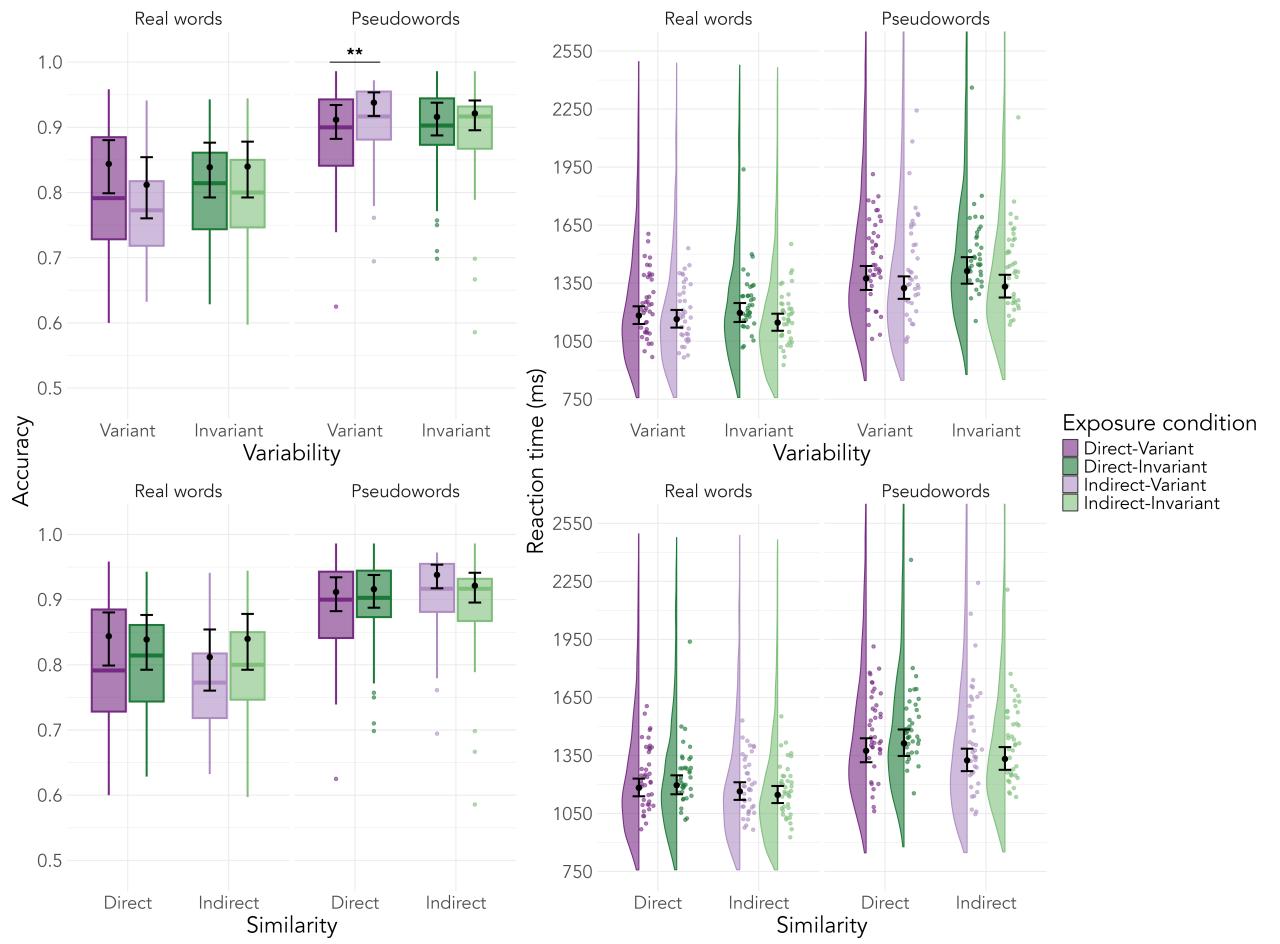


Figure 6: Experiment 2 exposure task performance. Left column presents boxplots with mean accuracy by participant. Right column presents half violin plots with reaction times for correct responses and dot plots with mean reaction times for correct responses by participant. Plots are overlaid with estimated marginal means and 95% confidence intervals. Asterisks indicate significance levels from pairwise comparisons: *** $< .001$, ** $< .01$, * $< .05$.

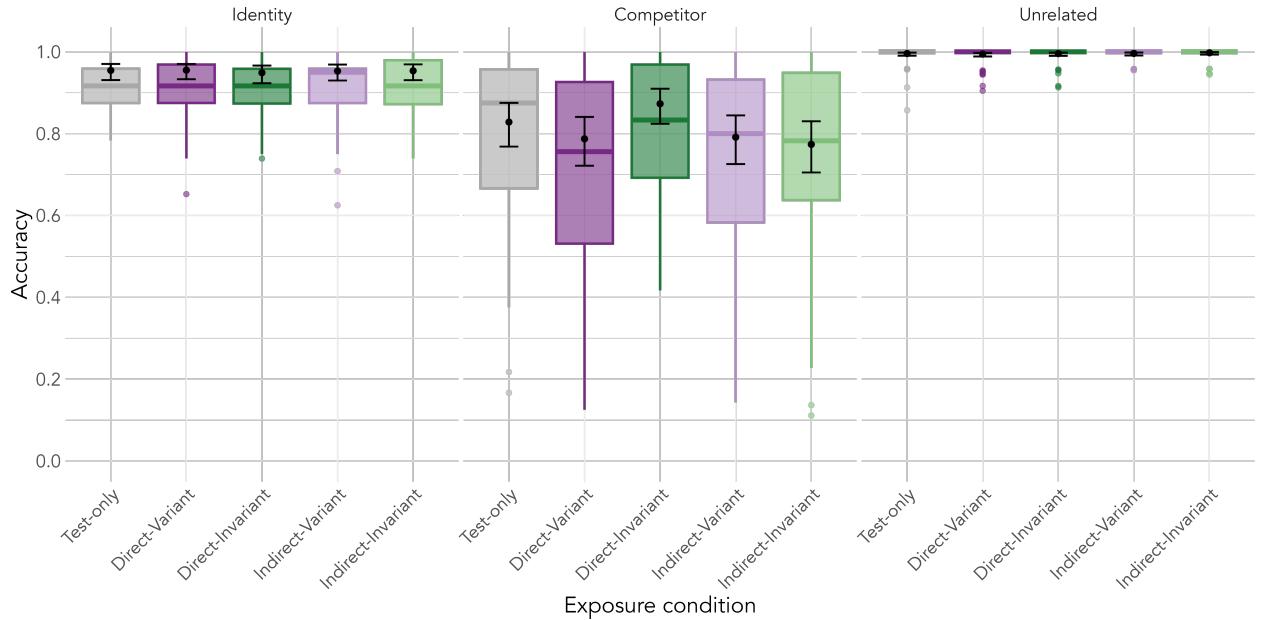


Figure 7: Experiment 2 test task accuracy. Boxplots present by-participant means for each Target type. Plots are overlaid with estimated marginal means and 95% confidence intervals.

722 4.3. Discussion

723 In Experiment 2, two levels of Similarity—Direct and Indirect—were crossed with two levels of Variability—
 724 Variant and Invariant—in order to investigate how experience with Spanish-accented stops transfers to a novel
 725 talker. During the exposure task, participants performed equally well across conditions. The only difference
 726 we observed was on pseudoword accuracy, where the Indirect-Variant group outperformed the Direct-Variant
 727 group. Overall, the results of the exposure task suggest that word recognition in Spanish-accented speech
 728 was successful regardless of Similarity or Variability.

729 As in Experiment 1, the exposure task did not benefit performance on the test task. Compared to the
 730 Test-only group, none of the exposure groups performed differently. However, we did observe differences
 731 between exposure groups. Specifically, the Direct-Invariant group was more accurate on Competitor targets
 732 than any of the other three exposure groups. This group was also slower on Identity targets than the
 733 Indirect-Variant group.

734 In order to interpret these findings in terms of generalization, we will briefly describe the structure of the
 735 exposure conditions again. Both Direct conditions exposed participants to Spanish-accented /p/, /t/, and
 736 /k/ in such disambiguating lexical contexts as *peanut*, *terminal*, and *kingdom*, respectively. However, the
 737 Direct-Invariant condition in particular exposed participants to one-to-one mappings between each critical
 738 onset and each talker, such that all /p/ onsets were produced by Talker A, all /t/ onsets by Talker B, and
 739 all /k/ onsets by Talker C. Thus, participants in the Direct-Invariant group were given the opportunity to
 740 develop talker-specific generative models for each VOT-stop mapping. Recall that a generative model refers
 741 to a listener’s mental representation of the distribution of a phonetic category (like /p/) over an acoustic
 742 cue (like VOT) under the ideal adapter framework. Since generative models are specific to pairs of cues and
 743 categories under this theory, listeners in the Invariant conditions had to organize their representations of VOT
 744 for each onset by talker. By contrast, listeners in the Variant conditions had the option to integrate across
 745 talkers to organize each VOT-stop mapping at the accent level, since all three talkers produced exemplars of
 746 all onsets. The high performance of the Direct-Invariant group on Competitor targets suggests that robust
 747 talker-specific generative models for voiceless stops may reduce the (erroneous) activation of voiced stops
 748 in response to Spanish-accented primes like *park*. The slight reduction in performance for this group on
 749 Identity targets suggests that such generative models may not increase the activation of voiceless stops

750 in turn. Together, these findings suggest that exposure-test similarity is necessary but not sufficient for
751 talker-independent perceptual adaptation.

752 Finally, we return to the lack of significant differences between the Test-only and Direct-Invariant groups.
753 The fact that participants without exposure were able to perform at a similar level to those with exposure
754 weakens the argument we put forward in the previous paragraph. If exposure facilitates generalization, but
755 generalization does not facilitate future performance, then what is the benefit of exposure? However, there is
756 a wealth of evidence that previous exposure to an L2 accent improves perception of a novel L2-accented
757 talker with the same L1 (Bent and Baese-Berk, 2021). Previous studies have generally used either sentence
758 transcription (e.g., Bradlow and Bent, 2008) or primed lexical decision (e.g., Xie and Myers, 2017) to test the
759 strength of adaptation. Here, we used a matching task, under the assumption that categorization of short lag
760 VOTs as voiceless stops should change as a function of perception of short lag VOTs as voiceless stops. For
761 example, consider the auditory prime *park* and the visual target *bark*. Participants needed to decide whether
762 the onset of the token they had heard was a /b/ or not. If they accurately perceived the onset as a /p/,
763 then they would correctly reject *bark* as a match. Thus, accuracy on the matching task was the outcome of
764 the categorization process. It is possible that this outcome-based measure was not fine-grained enough to
765 capture subtle changes in perception. To return to the Competitor target example, the lack of difference in
766 categorizing Spanish-accented *park* as *bark* may belie differences in perceiving Spanish-accented /p/. To
767 better assess changes in perception, we changed the test task for Experiment 3.

768 5. Experiment 3: Investigating fine-grained perceptual changes from exposure to Spanish-accented stops

769 5.1. Methods

770 The exposure phase was the same as Experiment 2, but the task used in the test phase was different.

772 5.1.1. Design

773 To better detect subtle changes in perception as a function of exposure to Spanish-accented speech,
774 we implemented the primed cross-modal lexical decision task from Xie and Myers (2017). This design is
775 illustrated in Figure 8. We maintained the same three types of Target: Identity, Competitor, and Unrelated.
776 However, participants performed a different task with these targets relative to Experiment 2. Specifically,
777 participants decided whether the visual target was a real English word or not. The auditory primes should
778 increase or decrease RTs on the visual targets as a function of perceptual adaptation to Spanish-accented
779 voiceless stops.

780 5.1.2. Participants

781 We recruited 195 participants through Prolific; none had participated in Experiments 1 or 2. All aspects
782 of recruitment and criteria for data exclusion were the same as in Experiments 1 and 2. After removing
783 ineligible participants, 170 remained. After removing participants with poor data quality, 158 remained for
784 analysis (Age: $M = 31$, $SD = 6$, Min = 18, Max = 40; Sex: Female = 80, Male = 78; Race: Asian = 4,
785 Black = 20, Multiple selected = 13, Other = 4, White = 117).

786 Additionally, a new group of 50 participants was also recruited to complete the experiment without
787 the exposure phase. There were 43 participants remaining after checking the eligibility criteria, and 42
788 participants remained after checking for data quality (Age: $M = 30$, $SD = 5$, Min = 21, Max = 40; Sex:
789 Female = 20, Male = 22; Race: Black = 3, Multiple selected = 5, White = 34).

790 5.1.3. Stimuli

791 The real words and pseudowords in the exposure task were the same as those in Experiment 2. The
792 auditory primes in the test task were also the same. The only difference was in the visual targets.

793 The unrelated target for each filler prime (144) was replaced with a pseudoword with a different filler
794 onset from the prime. Potential pseudowords were downloaded from the ELP's set of normed pseudowords.
795 The best possible match was selected for each prime by (orthographic) vowel and number of letters. Three
796 additional pseudoword targets were selected for practice.

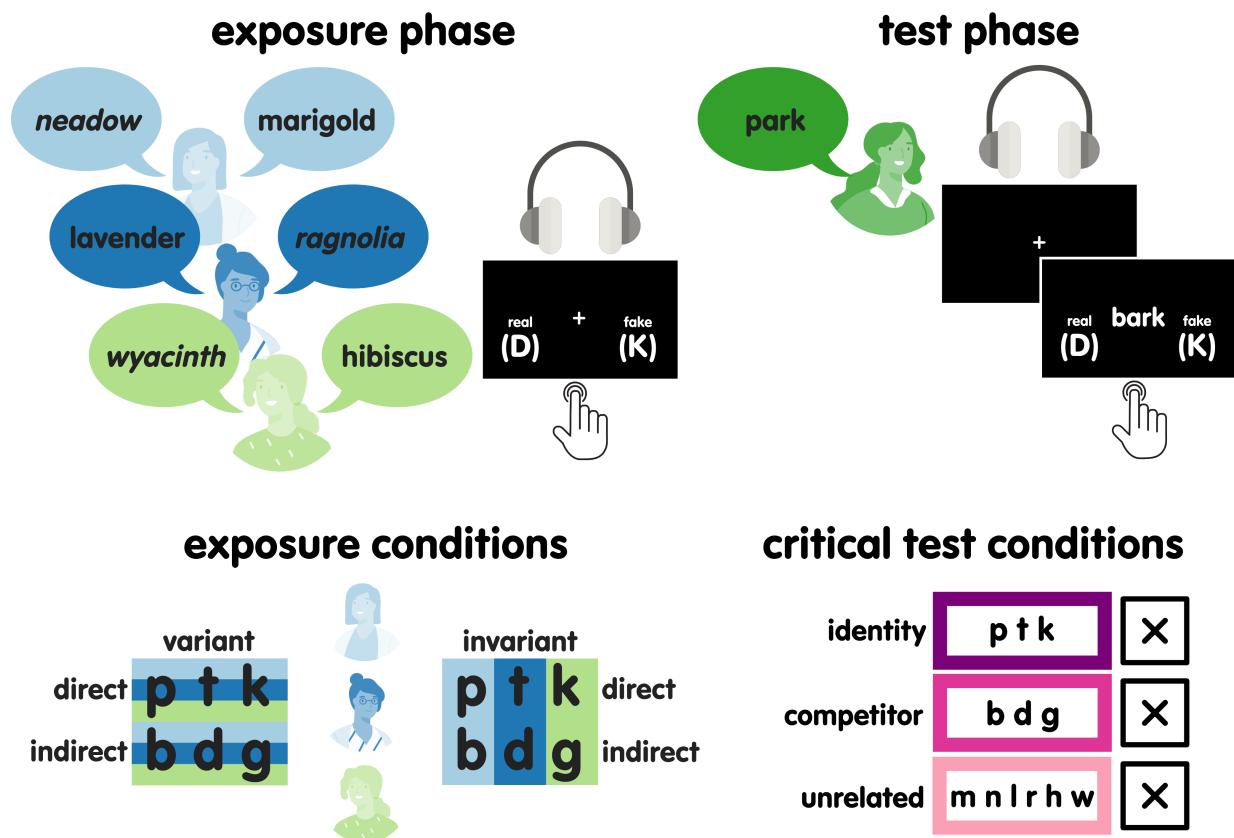


Figure 8: Experiment 3 design.

797 5.1.4. *Experimental lists*

798 The experimental lists for the exposure task were the same as in Experiment 2. Talker assignment was
799 also counterbalanced the same way.

800 For the critical primes (72), the combinations of auditory prime and visual target were counterbalanced
801 across participants in three experimental lists as in Experiment 2. For the filler primes (144), perfect
802 counterbalancing across these three lists was not possible, since three quarters of the filler items in each list
803 (108) needed to have unrelated pseudoword targets. Three sets of primes, divided evenly by onset, were
804 rotated through the assignment of Identity or Unrelated (pseudoword) target as evenly as possible. This
805 resulted in 12 experimental lists for the test phase, one for each combination of critical prime-target pair (3)
806 and talker assignment (4).

807 5.1.5. *Tasks and procedure*

808 The headphone check, exposure task, and post-experiment questionnaire were the same as Experiment 2.
809 The procedure was also the same. The only difference was in the structure of the test task.

810 The test phase featured the cross-modal primed lexical decision task from Xie and Myers (2017). On
811 each trial, participants first heard a real word (auditory prime) and then saw a real word or pseudoword
812 written on the screen (visual target). They indicated whether the visual target was a real English word or
813 not by pressing the *d* or *k* key on their keyboard. The real word response was mapped to the same key as in
814 the exposure task. Participants completed six practice trials followed by 216 main trials presented in random
815 order. Half of the practice trials (3) and half of the main trials (108) required real word responses. All of the
816 trials with critical primes (72) required real word responses.

817 5.1.6. *Analysis*

818 The data processing, model fitting, and analysis approaches were the same as in Experiment 2. Prior to
819 analyzing exposure task performance for real words with experimental onsets, we removed responses with
820 RTs less than 50 ms ($N = 22$; 0.03%). We then detected and removed outliers ($N = 2953$; 4.33%). Prior to
821 analyzing test task performance for critical prime-target pairs, we removed responses with RTs less than 50
822 ms ($N = 53$; 0.37%). We then detected and removed outliers ($N = 1076$; 7.50%).

823 5.2. *Results*

824 5.2.1. *Exposure*

825 For accuracy, there was a significant interaction between Variability and Word type ($\chi^2(1, N = 2) =$
826 17.77, $p < .001$), with higher accuracy on pseudowords in Variant groups ($M = 0.94$, 95% CI [0.92, 0.95])
827 compared to Invariant groups ($M = 0.92$, 95% CI [0.90, 0.94]; $z = 2.63$, $p = .009$). We also conducted
828 pairwise comparisons within each combination of Similarity and Word type and within each combination of
829 Variability and Word type. There were no significant effects of Variability or Similarity, respectively ($ps >$
830 .05). These comparisons are shown in the left column of Figure 9.

831 For RT, there was a significant three-way interaction among Variability, Similarity, and Word type ($\chi^2(1,$
832 $N = 2) = 9.11$, $p = .003$). The two-way interaction between Similarity and Word type was also significant
833 ($\chi^2(1, N = 2) = 4.40$, $p = .036$). To follow up on the three-way interaction, we analyzed the simple effects
834 of Variability and Similarity; however, none of the pairwise comparisons was significant ($ps > .05$). These
835 comparisons are shown in the right column of Figure 9.

836 5.2.2. *Test: Comparison to the Test-only group*

837 For accuracy, we observed a significant interaction between Exposure and Target ($\chi^2(8, N = 9) = 49.25$,
838 $p < .001$); however, pairwise comparisons within each level of Target did not reveal any significant differences
839 between the Test-only group and the exposure groups ($ps > .05$). Comparing the exposure groups to one
840 another, we found that accuracy on Identity targets was significantly higher in the Indirect-Invariant group
841 ($M = 1.00$, 95% CI [1.00, 1.00]) than in either the Indirect-Variant group ($M = 0.99$, 95% CI [0.99, 1.00];
842 $z = 2.89$, $p = .023$) or the Direct-Variant group ($M = 0.99$, 95% CI [0.99, 1.00]; $z = 2.62$, $p = .045$).
843 By-participant means and estimated marginal means are shown in the top row of Figure 10.

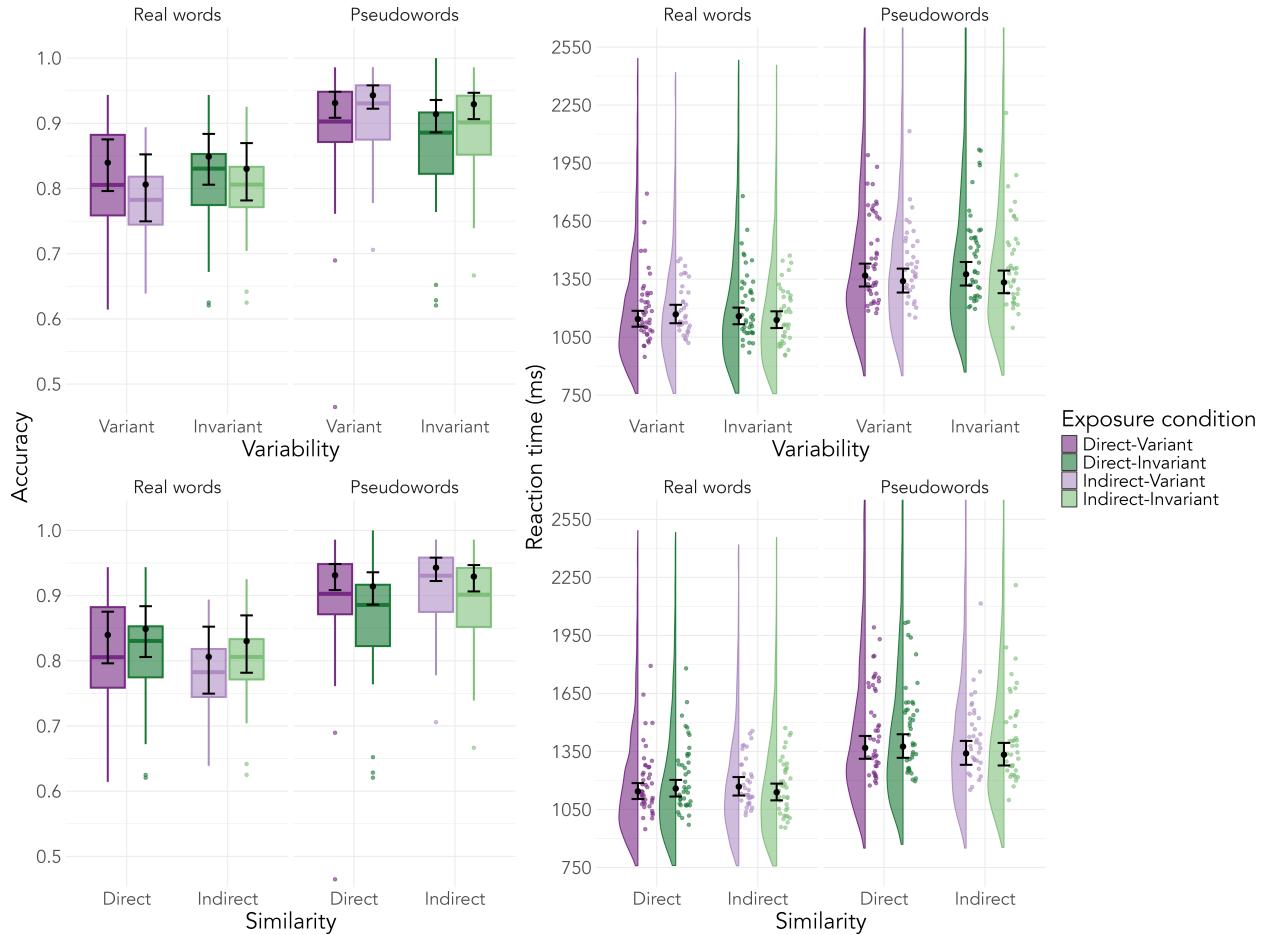


Figure 9: Experiment 3 exposure task performance. Left column presents boxplots with mean accuracy by participant. Right column presents half violin plots with reaction times for correct responses and dot plots with mean reaction times for correct responses by participant. Plots are overlaid with estimated marginal means and 95% confidence intervals.

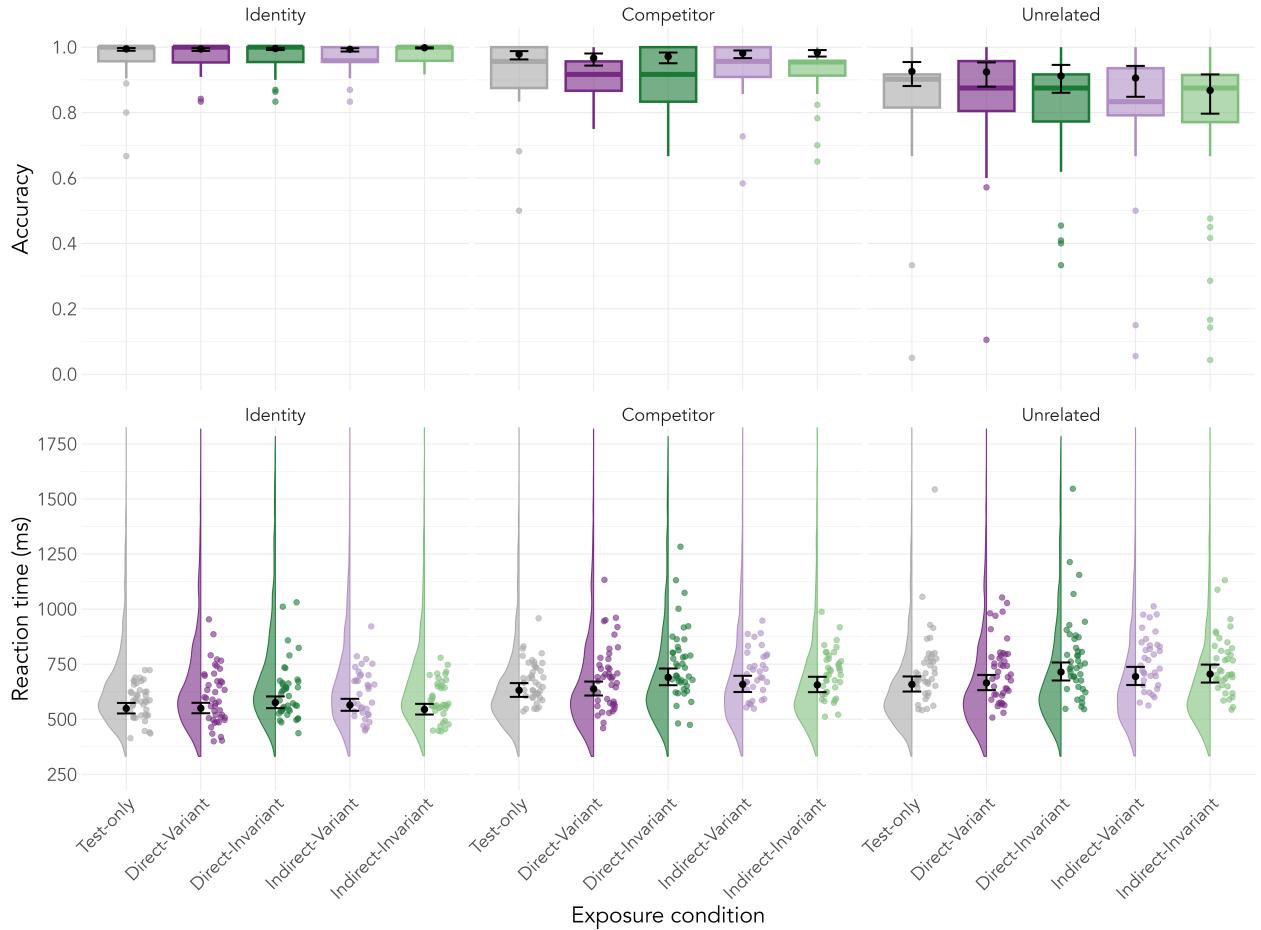


Figure 10: Experiment 3 test task performance. Top row presents boxplots with mean accuracy by participant. Bottom row presents half violin plots with reaction times for correct responses and dot plots with mean reaction times for correct responses by participant. Plots are overlaid with estimated marginal means and 95% confidence intervals.

For RT, we also observed a significant interaction between Exposure and Target ($\chi^2(8, N = 9) = 33.56, p < .001$). This interaction was driven by slower RTs on Competitor targets in the Direct-Invariant group ($M = 690, 95\% \text{ CI } [654, 730]$) than in the Test-only group ($M = 631, 95\% \text{ CI } [601, 664]; z = 2.54, p = .044$). Distributions, by-participant means, and estimated marginal means are shown in Figure 10. To investigate the differences between exposure conditions, we conducted the second set of analyses comparing the effects of Variability and Similarity on RTs.

5.2.3. Test: Comparison between Variability and Similarity

For RT, we observed significant two-way interactions between Variability and Target ($\chi^2(2, N = 3) = 13.87, p < .001$) and between Similarity and Target ($\chi^2(2, N = 3) = 9.05, p = .011$). We followed up on these interactions with separate pairwise comparisons for Variability and Similarity within each level of Target; however, the pairwise comparisons between levels of Variability and between levels of Similarity were not significant ($ps > .05$). To further probe the source of the interactions in the model, we conducted pairwise comparisons between Direct and Indirect exposure within each combination Target and Variability and between Variant and Invariant exposure within each combination of Target and Similarity. Within the Competitor level of Target and Direct level of Similarity, Invariant exposure yielded slower RTs ($M = 691, 95\% \text{ CI } [653, 732]$) than Variant exposure ($z = 2.17, p = .030$). None of the other pairwise comparisons by Target and Similarity was significant, nor were any of the pairwise comparisons by Target and Variability ($ps > .05$).

861 > .05).

862 *5.3. Discussion*

863 Experiment 3 followed up on the findings of Experiment 2 with a different test task. Recall that in
864 Experiment 2, we observed significantly higher accuracy on Competitor targets in the Direct-Invariant
865 group than in any of the other exposure groups. However, the Test-only group performed similarly to the
866 exposure groups, limiting our interpretation of the effects of Variability and Similarity. We posited that
867 Direct-Invariant training reduced lexical competition between voiced and voiceless stops, thereby increasing
868 correct rejection of Competitor targets as matches for the Spanish-accented auditory primes. The matching
869 task required participants to explicitly compare the visual target to the auditory prime. In Experiment 3, we
870 took a more implicit approach. The priming task probed the extent to which the auditory prime increased
871 or decreased activation of the visual target. In this way, we could investigate changes in the perception of
872 Spanish-accented stops.

873 When we compared the four exposure conditions to the Test-only condition, we observed slower RTs
874 for Competitor targets after Direct-Invariant exposure. We interpret this reduction in speed as a reduction
875 in lexical competition. For example, consider the auditory prime *park* and the visual target *bark*, which
876 are minimal pairs that differ only in the voicing of their onsets. The more the onset of *park* is perceived as
877 /b/, the more it will activate the target *bark*. This increase in activation will facilitate the lexical decision
878 for *bark*, resulting in faster RTs. Our results suggest that, in the absence of exposure to Spanish-accented
879 speech, Spanish-accented *park* was perceived **more** like *bark*. By contrast, with Direct-Invariant exposure,
880 Spanish-accented *park* was perceived **less** like *bark*. This reduction in the activation of voiced targets after
881 training suggests that talker-specific exposure to Spanish-accented /p/, /t/, and /k/ improved phonetic
882 categorization of short lag VOTs as voiceless. Based on this evidence for talker-independent adaptation, we
883 conducted analyses to distinguish the effects of Variability and Similarity.

884 Within the Direct conditions, Invariant exposure reduced the activation of voiced competitors more
885 than Variant exposure. This effect was illustrated by significantly slower RTs on Competitor targets for
886 the Direct-Invariant group than for the Direct-Variant group. As described in the previous paragraph, the
887 Direct-Invariant group was also slower on Competitor targets than the Test-only group. Together, these
888 results suggest that the Test-only and Direct-Variant groups exhibited similar levels of *park-bark* priming,
889 indexing increased activation of /b/ by Spanish-accented /p/. This finding is striking, considering that the
890 Direct-Variant group had the same amount of exposure to Spanish-accented /p/, /t/, and /k/ from the same
891 talkers as the Direct-Invariant group. We will interpret this effect of Variability more fully in Section 7. In
892 short, we will argue that listeners were not able to develop robust talker-independent models during Variant
893 exposure. Because this level of organization is subject to listeners' use of indexical and social information
894 (Kleinschmidt, 2019), we will investigate their perceptions of the test talkers in the next section.

895 **6. Talker analysis: Characterizing voiceless stop VOT distributions**

896 In this section, we analyze the VOT distributions of the talkers. Table 3 provides descriptive statistics of
897 the VOT distributions for both voiced and voiceless stops by task, phoneme, and talker.

898 *6.1. Talker-specific analysis of voiceless stop VOT*

899 We first investigate the extent to which each talker's voiceless stop VOTs diverge from L1 US English
900 "norms," which represent our listeners' prior beliefs about these distributions. We used the VOT dataset
901 compiled by Chodroff et al. (2019) to determine the mean VOT for each voiceless stop onset among L1
902 US English talkers (see Table 4). We then conducted one-sample t-tests to compare each of the talkers'
903 mean VOTs to the norm. These analyses included both the exposure (multisyllabic) and test (monosyllabic)
904 stimuli (48 items per onset per talker). Table 4 provides the difference in means and asterisks indicating the
905 significance level from each analysis (corrected for multiple comparisons with the Hommel method).

Table 3: Mean, standard deviation, and range of VOTs by onset by talker.

Phase	Onset	N	Talker 1			Talker 2			Talker 3			Talker 4		
			M	SD	Range	M	SD	Range	M	SD	Range	M	SD	Range
Exposure	/b/	24	-21	51	-132-48	15	27	-87-77	-44	52	-150-31	-69	53	-133-103
	/d/	24	2	42	-100-39	21	22	-75-53	8	29	-54-33	-78	53	-177-44
	/g/	24	2	55	-104-68	31	25	-75-51	-5	51	-148-47	-64	65	-229-67
	/p/	24	73	26	24-122	61	41	11-160	82	26	37-123	21	17	6-69
	/t/	24	72	20	36-112	65	22	25-109	91	28	46-157	41	15	18-74
	/k/	24	82	28	33-151	88	21	49-123	96	23	43-133	48	16	25-77
Test	/p/	24	71	28	19-126	60	35	10-118	115	39	15-201	26	24	8-106
	/t/	24	78	21	27-109	61	20	28-106	102	23	60-142	53	27	8-99
	/k/	24	94	21	70-171	94	20	50-133	119	40	65-244	65	30	21-124

Table 4: Mean difference between L1 US English and talker mean VOTs by onset. Asterisks indicate significance levels from two-sample t-tests: *** < .001, ** < .01, * < .05. English means from Chodroff and Wilson (2019).

Phone	L1 US English mean	Talker 1	Talker 2	Talker 3	Talker 4
/p/	66	6	-5	33***	-43***
/t/	78	-3	-15***	19***	-31***
/k/	70	18***	21***	37***	-14**

6.2. Between-talker analysis of voiceless stop VOT

We next investigate the extent to which each talker's mean VOTs differed from the other talkers'. This analysis provides a point of comparison to Xie and Myers' (2017) operationalization of exposure-test similarity. Here, we conducted paired two-sample t-tests on each combination of talkers by onset (48 items per onset per talker). Table 5 shows the results of these comparisons, with the value in each cell being the column talker's mean minus the row talker's mean.

6.3. Discussion

The two analyses conducted in this section show that our talkers differed from both L1 US English norms and one another in their VOT distributions. Talker 4 produced the most Spanish-like VOT distributions across places of articulation, with significantly shorter means in all comparisons. Talker 3 also differed from the L1 US English norms and the three other talkers, but in the opposite direction: her mean VOTs were significantly longer in all comparisons. Talkers 1 and 2 were the most similar to one another, with statistically equivalent mean VOTs for /p/ and /k/. Talker 1 was the most similar to the L1 US English norms, though both she and Talker 2 had significantly longer /k/ VOTs. Overall, each talker exhibited a unique pattern of VOT-stop distributions. We discuss the implications of these results for our Similarity and Variability manipulations in the next section.

7. General discussion

The current study investigated generalization of exposure to L2-accented speech. We compared two hypotheses with different explanations for how exposure supports talker-independent adaptation. The exposure-to-variability hypothesis argues that increasing covariation between acoustic-phonetic cues during exposure increases performance on a new talker (Baese-Berk et al., 2013; Bradlow and Bent, 2008). The similarity-based hypothesis argues that increasing cue-category overlap between the exposure and test talkers increases performance on the test talker (Xie and Myers, 2017). We expressed these hypotheses in terms of the ideal adapter framework to investigate their explanatory power within the same theoretical model

Table 5: Mean difference in VOT between talkers by onset. Asterisks indicate significance levels from paired two-sample t-tests corrected with Hommel method: *** < .001, ** < .01, * < .05.

Comparison	Exposure onset	Talker 1	Talker 2	Talker 3	Talker 4
Talker 1	/p/		-11	27**	-49***
	/t/		-12*	22***	-28***
	/k/		3	19**	-32***
Talker 2	/p/	11		38***	-37***
	/t/	12*		34***	-16**
	/k/	-3		16*	-35***
Talker 3	/p/	-27**	-38***		-75***
	/t/	-22***	-34***		-50***
	/k/	-19**	-16*		-51***
Talker 4	/p/	49***	37***	75***	
	/t/	28***	16**	50***	
	/k/	32***	35***	51***	

930 (Kleinschmidt and Jaeger, 2015). Under this theory, listeners develop mental models of the mappings between
931 acoustic cues (like VOT) and phonetic categories (like /p/, /t/, and /k/) according to informative and useful
932 social (here, L2 accent) and indexical (talker) groupings (Kleinschmidt, 2019). To restate the two competing
933 hypotheses using this framework, exposure to variability enhances the formation of talker-independent
934 generative models, while exposure-test similarity facilitates the selection of talker-specific generative models.
935 Using a novel experimental approach that improved the scope and precision of how variability and similarity
936 are operationalized, we conducted a series of three experiments. The combined results provide strongest
937 support for the similarity-based hypothesis.

938 The key findings come from Experiment 3. After exposure to three Spanish-accented English talkers,
939 we measured performance on a novel Spanish-accented English talker with a primed cross-modal lexical
940 decision task. The critical items in this task were monosyllabic auditory primes with voiceless stop onsets:
941 for example, *park*, *tune*, and *coal*. Importantly, these items had onset competitors with the same place and
942 manner of articulation: *bark*, *dune*, and *goal*, respectively. Participants made lexical decisions on the visual
943 targets that followed these critical primes. Performance on the different target types indexed different aspects
944 of perceptual adaptation. Specifically, priming of Identity targets indexed lexical activation of the intended
945 word, while priming of Competitor targets indexed lexical activation of perceptually similar words. For
946 example, consider the auditory prime *park*. The Identity target for this prime is *park*, while the Competitor
947 target for this prime is *bark*. If perceptual adaptation to Spanish-accented speech increased activation of
948 the lexical item “*park*” upon hearing *park*, then responses to the visual target *park* would also increase.
949 If perceptual adaptation to Spanish-accented speech decreased activation of the lexical item “*bark*” upon
950 hearing *park*, then responses to the visual target *bark* would also decrease. Ideally, exposure would both
951 increase activation of the intended word and decrease activation of competing words. However, we found that
952 Direct-Invariant exposure decreased lexical competition without increasing lexical activation. Specifically,
953 the Direct-Invariant group displayed significantly slower RTs on Competitor targets than both the Test-only
954 group and Direct-Variant group, but did *not* display significantly faster RTs on Identity targets. Thus,
955 the decrease in lexical competition we observed was not related to an increase in lexical activation for the
956 intended word. In the next section, we compare our findings to the predictions of the exposure-to-variability
957 and similarity-based hypotheses.

958 7.0.1. *Exposure to variability versus similar exposure*

959 The exposure-to-variability hypothesis is based on the findings of Baese-Berk et al. (2013). In that study,
960 exposure to multiple talkers, each with a different L2 accent, generalized to a novel talker with an unfamiliar
961 L2 accent. Exposure to multiple talkers with the same L2 accent did not generalize to an unfamiliar L2
962 accent (Bradlow and Bent, 2008). However, both types of exposure generalized to a novel talker with a
963 familiar L2 accent. The multi-accent condition included one Mandarin-accented talker out of five talkers
964 with different L1s, while the single-accent condition included five Mandarin-accented talkers. In both cases,
965 performance on a novel Mandarin-accented talker was higher than in the control condition. Together, these
966 results suggest that listeners benefit from exposure to systematic covariation across L2 accents. In other
967 words, the ways in which L2 varieties differ from L1 varieties at a high level generalizes across L2 talkers.
968 This hypothesis is also supported by recent work on cross-accent generalization (Bradlow et al., 2023), as
969 well as by the large body of work on high-variability phonetic training for L2 acquisition (for a review, see
970 Zhang et al., 2021).

971 By contrast, the similarity-based hypothesis comes from Xie and Myers (2017). In this study, two groups
972 of participants trained on multiple Mandarin-accented talkers. The experimental group was exposed to
973 disambiguating lexical contexts for word-final /d/ (e.g., *overload*), while the other was not. During test on a
974 novel Mandarin-accented talker, lexical activation for /d/-final primes was higher in the experimental group.
975 Among the Mandarin-accented talkers in the training set, one talker's productions of /d/ were similar to those
976 of the test talker on multiple acoustic measures. There was also one talker out of the five whose productions
977 of /d/ differed from the test talker's on these metrics. After single-talker exposure to the similar talker, lexical
978 activation for /t/-final competitors was lower in the experimental group. By contrast, single-talker exposure
979 to the dissimilar talker did not reduce lexical competition. Overall, single-talker exposure to the similar
980 talker and multi-talker exposure including the similar talker both facilitated generalization to a novel talker.
981 This suggests that the specific way in which an L2-accented talker produces an L2 variety determines the
982 level of generalization to other talkers. This hypothesis is supported by work on lexically-guided perceptual
983 retuning showing different patterns of generalization for different phonetic contrasts (Kraljic and Samuel,
984 2006, 2007; Reinisch and Holt, 2014).

985 Having revisited the two competing hypotheses that motivated our work, we return to the specific
986 predictions and results. The exposure-to-variability hypothesis predicted a main effect of Variability, with
987 better test performance after Variant versus Invariant exposure. Participants in both conditions were
988 also expected to outperform the Test-only group. This hypothesis does not make predictions regarding
989 activation of intended words (Identity targets) versus activation of competing words (Competitor targets).
990 However, given that the hypothesis is based on research where transcription accuracy is the measure of
991 generalization, Variant exposure was likely to emerge on Identity target performance. The similarity-based
992 hypothesis predicted a main effect of Similarity, with better test performance after Direct versus Indirect
993 (or Control) exposure. Participants in the Direct condition were also expected to outperform the Test-only
994 group. Specifically, Direct exposure was expected to both increase activation for intended words and decrease
995 activation for competitors.

996 Our results provided evidence for the similarity-based hypothesis and against the exposure-to-variability
997 hypothesis. In both Experiments 2 and 3, Direct-Invariant exposure decreased activation for competitors.
998 Specifically, we observed an increase in RTs for Competitor targets in the Direct-Invariant group relative
999 to both the Test-only group and the Direct-Variant group in Experiment 3. In Experiment 2, we observed
1000 an increase in accuracy on Competitor targets in the Direct-Invariant group relative to all three exposure
1001 groups. The performance benefits from Direct exposure align best with the predictions of the similarity-based
1002 hypothesis. In addition, the performance benefits from Invariant exposure contrast the predictions of the
1003 exposure-to-variability hypothesis. To better understand these findings, we return to how we operationalized
1004 Similarity and Variability relative to previous research.

1005 7.0.2. *Implementing similarity*

1006 Regarding similarity, we exposed participants to different sets of Spanish-accented stops through an
1007 auditory lexical decision task. In Direct conditions, participants gained experience with /p/, /t/, and /k/

1008 onsets in the context of real words like *peanut*, *terminal*, and *kingdom*, respectively. Exposure to Spanish-
1009 accented voiceless stops trained listeners to shift the VOT distributions for these sound categories from long
1010 lag (English-like) toward short lag (Spanish-like). In Indirect conditions, participants gained experience with
1011 /b/, /d/, and /g/ in the context of real words like *beehive*, *desert*, and *gallop*, respectively. Exposure to
1012 Spanish-accented voiced stops trained listeners to shift the VOT distributions for these sound categories from
1013 short lag (English-like) toward lead (Spanish-like). Shifting the VOT distribution for voiced stops leftward
1014 on the cue continuum had the potential to trigger a more general leftward shift that would move the VOT
1015 distributions for voiceless stops from long lag to short lag. This approach to investigating similarity involved
1016 exposure to different sound categories.

1017 By contrast, the approach of Xie and Myers (2017) involved exposure to different cue distributions over
1018 the same sound category. That study compared the means of three different cues to voicing in word-final
1019 position to determine which talker had the most (dis)similar cue-category distributions to the test talker;
1020 separate groups of participants were then exposed to either the similar talker or the dissimilar talker. As
1021 a result, their design featured different talkers across exposure conditions. Our design featured the same
1022 talkers across exposure conditions, allowing us to control the amount of talker-specific acoustic-phonetic
1023 variability between levels of Similarity. In turn, this allowed us to manipulate the constructs of Similarity
1024 and Variability independently. We did not, however, control the exact VOT-stop distributions that listeners
1025 heard during exposure. We ultimately found in Section 6 that Talker 4 was the only one of the four with
1026 consistently Spanish-accented voiceless stops.

1027 The fact that our talkers produced less ambiguous voiceless stops than expected may have weakened the
1028 adaptation effects we observed; however, it did not disrupt the Similarity manipulation itself. The VOT
1029 distributions for the /p/, /t/, and /k/ onsets in the Direct exposure conditions were distinct from those
1030 for the /b/, /d/, and /g/ onsets in the Indirect exposure conditions. Critically, they were also similar to
1031 those for the test items. As a result, Direct exposure was both useful and informative for adapting to the
1032 test talker in a way that Indirect exposure was not. We also note that it is relatively uncommon in the
1033 L2-accented speech recognition literature to report acoustic-phonetic measures at all (cf. Xie et al., 2021; @
1034 Alexander and Nygaard, 2019), even though they presumably comprise what we perceive as an accent itself.
1035 In the present study, these measures help us understand the structure of our listeners' generative models for
1036 Spanish-accented speech, as we describe in the next section.

1037 7.0.3. Implementing variability

1038 Regarding variability, we exposed participants to different combinations of talker and onset. In Variant
1039 conditions, participants gained experience with the VOT distributions for each onset across talkers. Since
1040 there were 24 real word stimuli per onset, this means that each of the three talkers produced eight experimental
1041 items with each onset. In Invariant conditions, participants gained experience with the VOT distributions
1042 for each onset within talkers. This means that each talker produced 24 experimental items with a single
1043 onset. The variability manipulation was also extended to the filler items, with two filler onsets assigned to
1044 each talker. Overall, listeners heard 72 real words from each talker during exposure.

1045 By contrast, in Bradlow and Bent (2008) and Baese-Berk et al. (2013), participants in the multi-talker
1046 exposure groups heard 16 sentences with 50 total key words per talker. While each set of sentences was
1047 different for each exposure talker, participants completed five repetitions per talker, increasing the total
1048 amount of exposure to each talker to 250 key words. This is nearly three times the amount of word-level
1049 exposure to each talker compared to our study. Previous research has shown that listeners need significantly
1050 more exposure to adapt to multiple talkers compared to a single talker (Luthra et al., 2021). Even in a
1051 talker-specific paradigm, test performance increases with increasing evidence for adaptation (Cummings and
1052 Theodore, 2023). Thus, the relatively limited amount of exposure our participants had to variability may
1053 have limited the strength of its benefits. However, the amount of exposure did not differ between our levels
1054 of Variability. Moreover, we observed a clear benefit for Invariant versus Variant exposure. We interpret this
1055 finding in terms of the specificity of generative models.

1056 According to Kleinschmidt (2019), generative models can be constructed according to an individual person
1057 (talker-specific) or a higher-level social category (talker-independent). The level of organization is determined
1058 by two factors. First, the model needs to be a good representation of the actual distribution. That is, the

1059 cue-category mappings that the listener has heard need to be captured by the mental representation (i.e.,
1060 informative). Second, the model needs to be a good predictor of future exemplars. In other words, the
1061 mental representation needs to help listeners categorize a given instance of a cue as a particular category (i.e.,
1062 useful). In our experiment, we gave listeners different samples of each talker’s VOT distributions. We can
1063 illustrate this idea with hypothetical Talker A and her VOT-/p/ mapping. Variant exposure give listeners a
1064 small sample of Talker A’s VOT-/p/ mapping relative to Invariant exposure, where all VOT-/p/ mappings
1065 are from Talker A. Direct exposure allows listeners to sample Talker A’s VOT distribution for /p/, while
1066 Indirect exposure only allows listeners to sample her VOT distribution for the voiced counterpart /b/. These
1067 different samples influenced how listeners constructed their representations.

1068 As explained in Section 4.3, talker-specific generative models were the only possible representations of
1069 Invariant exposure. This is because Invariant exposure only provided one talker’s cue distributions for a
1070 given category. With Variant exposure, both talker-specific and talker-independent models were possible;
1071 however, talker-independent models were optimal. This is because there were only 24 exemplars of each
1072 phonetic category divided among three talkers. Thus, talker-specific models would have been developed
1073 from a relatively sparse sample of eight exemplars. If we assume that the talkers had similar VOT-stop
1074 distributions, the talker-specific generative models developed during Invariant exposure should have been
1075 as useful as the talker-independent generative models developed during Variant exposure. However, the
1076 reduction in lexical competition we observed following Direct-Invariant exposure compared to Direct-Variant
1077 exposure suggests that this was not the case.

1078 There are two possible explanations. First, the talker-independent generative models that listeners
1079 developed from Variant exposure were not as useful as the talker-specific generative models that listeners
1080 developed from Invariant exposure. Second, listeners developed talker-specific generative models regardless
1081 of variability, which led to sub-optimal representations of Variant exposure. Both of these accounts would
1082 be consistent with exposure to different VOT distributions for each talker. The post-hoc talker analysis
1083 in Section 6 showed that this was indeed the case. We found that mean VOTs differed between all four
1084 talkers at almost every place of articulation. A third explanation would be that exposure to such distinct
1085 Spanish-accented talkers would result in a general relaxation of categorization criteria (rather than shifts
1086 in categorization boundaries); however, if this were the case, we would not have observed any differences
1087 in performance as a function of exposure group. Rather, we observed a unique benefit of Direct-Invariant
1088 exposure broadly consistent with the similarity-based hypothesis. In terms of the ideal adapter framework,
1089 these findings suggest that developing talker-specific generative models for multiple talkers may be the best
1090 approach when between-talker variation is high.

1091 7.0.4. Conclusion

1092 Overall, our results show that talker-independent adaptation to L2-accented speech is facilitated by
1093 exposure-test similarity. We also found evidence that generalization is enhanced when between-talker
1094 variability is limited, at least for the rapid adaptation effects we investigated here. This study is the first
1095 to compare variability and similarity while maintaining the same talkers across conditions. This level of
1096 control is key to understanding how exposure to between-talker variability interacts with exposure to specific
1097 acoustic-phonetic properties. We suggest that the contradictory findings in the literature are at least partially
1098 a result of comparing different talkers, who have idiosyncratic patterns of within-talker covariation (Clayards,
1099 2017; Whalen et al., 2018; Xie and Jaeger, 2020). Listeners are highly sensitive to variation in speech,
1100 particularly during perception. In future work, we hope to better understand how this sensitivity is deployed
1101 during talker- and category-independent adaptation to L2-accented speech.

1102 References

- 1103 Arthur S Abramson and Douglas H Whalen. Voice onset time (VOT) at 50: Theoretical and practical issues in measuring
1104 voicing distinctions. *Journal of phonetics*, 63:75–86, 2017.
1105 Jessica E. D. Alexander and Lynne C Nygaard. Specificity and generalization in perceptual adaptation to accented speech. *The*
1106 *Journal of the Acoustical Society of America*, 145(6):3382–3398, 2019.
1107 Melissa M Baese-Berk, Ann R Bradlow, and Beverly A Wright. Accent-independent adaptation to foreign accented speech. *The*
1108 *Journal of the Acoustical Society of America*, 133(3):EL174–EL180, 2013.

- 1109 Miriam Baigorri, Luca Campanelli, and Erika S Levy. Perception of American–English vowels by early and late Spanish–English
 1110 bilinguals. *Language and speech*, 62(4):681–700, 2019.
- 1111 David A Balota, Melvin J Yap, Keith A Hutchison, Michael J Cortese, Brett Kessler, Bjorn Loftis, James H Neely, Douglas L
 1112 Nelson, Greg B Simpson, and Rebecca Treiman. The English lexicon project. *Behavior research methods*, 39:445–459, 2007.
- 1113 Dale J Barr, Roger Levy, Christoph Scheepers, and Harry J Tily. Random effects structure for confirmatory hypothesis testing:
 1114 Keep it maximal. *Journal of memory and language*, 68(3):255–278, 2013.
- 1115 Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. Fitting linear mixed-effects models using lme4. *Journal of
 1116 Statistical Software*, 67(1):1–48, 2015. doi: 10.18637/jss.v067.i01.
- 1117 José R Benkí. Place of articulation and first formant transition pattern both affect perception of voicing in English. *Journal of
 1118 Phonetics*, 29(1):1–22, 2001.
- 1119 Tessa Bent and Melissa Baese-Berk. *Perceptual learning of accented speech*, chapter The handbook of speech perception, pages
 1120 428–464. Wiley Blackwell, 2021.
- 1121 Catherine T Best. The emergence of native-language phonological influences in infants: A perceptual assimilation model. *The
 1122 development of speech perception: The transition from speech sounds to spoken words*, 167(224):233–277, 1994.
- 1123 Catherine T Best and Michael Tyler. *Second language speech learning: The role of language experience in speech perception
 1124 and production*, chapter Nonnative and second-language speech perception: Commonalities and complementarities, pages
 1125 13–34. John Benjamins, 2007.
- 1126 Catherine T Best, Gerald W McRoberts, and Elizabeth Goodell. Discrimination of non-native consonant contrasts varying in
 1127 perceptual assimilation to the listener’s native phonological system. *The Journal of the Acoustical Society of America*, 109
 1128 (2):775–794, 2001.
- 1129 Rebecca E Bieber and Sandra Gordon-Salant. Semantic context and stimulus variability independently affect rapid adaptation
 1130 to non-native english speech in young adults. *The Journal of the Acoustical Society of America*, 151(1):242–255, 2022.
- 1131 Richard E Blakesley, Sati Mazumdar, Mary Amanda Dew, Patricia R Houck, Gong Tang, Charles F Reynolds III, and Meryl A
 1132 Butters. Comparisons of methods for multiple hypothesis testing in neuropsychological research. *Neuropsychology*, 23(2):255,
 1133 2009.
- 1134 Ann R Bradlow and Tessa Bent. Perceptual adaptation to non-native speech. *Cognition*, 106(2):707–729, 2008.
- 1135 Ann R Bradlow, Adrianna M Bassard, and Ken A Paller. Generalized perceptual adaptation to second-language speech:
 1136 Variability, similarity, and intelligibility. *The Journal of the Acoustical Society of America*, 154(3):1601–1613, 2023.
- 1137 P Broersma and D Weenink. Praat: Doing phonetics by computer, 2021. URL <http://www.praat.org/>.
- 1138 Marc Brysbaert, Paweł Mandera, Samantha F McCormick, and Emmanuel Keuleers. Word prevalence norms for 62,000 English
 1139 lemmas. *Behavior research methods*, 51:467–479, 2019.
- 1140 Taehong Cho and Peter Ladefoged. Variation and universals in VOT: Evidence from 18 languages. *Journal of phonetics*, 27(2):
 1141 207–229, 1999.
- 1142 Eleanor Chodroff and Colin Wilson. Structure in talker-specific phonetic realization: Covariation of stop consonant VOT in
 1143 American English. *Journal of Phonetics*, 61:30–47, 2017.
- 1144 Eleanor Chodroff, Alessandra Golden, and Colin Wilson. Covariation of stop voice onset time across languages: Evidence for a
 1145 universal constraint on phonetic realization. *The Journal of the Acoustical Society of America*, 145(1):EL109–EL115, 2019.
- 1146 Ja Young Choi and Tyler K Perrachione. Time and information in perceptual adaptation to speech. *Cognition*, 192:103982,
 1147 2019.
- 1148 Constance M Clarke and Merrill F Garrett. Rapid adaptation to foreign-accented English. *The Journal of the Acoustical Society
 1149 of America*, 116(6):3647–3658, 2004.
- 1150 Meghan Clayards. Individual talker and token covariation in the production of multiple cues to stop voicing. *Phonetica*, 75(1):
 1151 1–23, 2017.
- 1152 Meghan Clayards, Michael K Tanenhaus, Richard N Aslin, and Robert A Jacobs. Perception of speech reflects optimal use of
 1153 probabilistic speech cues. *Cognition*, 108(3):804–809, 2008.
- 1154 William E Cooper. Selective adaptation for acoustic cues of voicing in initial stops. *Journal of Phonetics*, 2(4):303–313, 1974.
- 1155 Shawn N Cummings and Rachel M Theodore. Hearing is believing: Lexically guided perceptual learning is graded to reflect the
 1156 quantity of evidence in speech input. *Cognition*, 235:105404, 2023.
- 1157 James E Flege. Language contact in bilingualism: Phonetic system interactions. *Laboratory phonology*, 9(353-381), 2007.
- 1158 James E Flege and O Bohn. The revised speech learning model (slm-r). *second language speech learning: theoretical and
 1159 empirical progress*, pages 3–83, 2021.
- 1160 James E Flege and Wieke Eefting. Production and perception of english stops by native spanish speakers. *Journal of phonetics*,
 1161 15(1):67–83, 1987.
- 1162 James E Flege, Carlo Schirru, and Ian R. A. MacKay. Interaction between the native and second language phonetic subsystems.
 1163 *Speech communication*, 40(4):467–491, 2003.
- 1164 John Fox and Sanford Weisberg. *An R Companion to Applied Regression*. Sage, Thousand Oaks CA, third edition, 2019. URL
 1165 <https://socialsciences.mcmaster.ca/jfox/Books/Companion/>.
- 1166 Stephen D Goldinger. Echoes of echoes? an episodic theory of lexical access. *Psychological Review*, 105(2):251, 1998.
- 1167 Mia Hubert and Ellen Vandervieren. An adjusted boxplot for skewed distributions. *Computational statistics & data analysis*,
 1168 52(12):5186–5201, 2008.
- 1169 Keith Johnson. Resonance in an exemplar-based lexicon: The emergence of social identity and phonology. *Journal of phonetics*,
 1170 34(4):485–499, 2006.
- 1171 Dave F Kleinschmidt. Structure in talker variability: How much is there and how much can it help? *Language, cognition and
 1172 neuroscience*, 34(1):43–68, 2019.
- 1173 Dave F Kleinschmidt and T Florian Jaeger. Robust speech perception: recognize the familiar, generalize to the similar, and

- adapt to the novel. *Psychological Review*, 122(2):148–203, 2015.
- Tanya Kraljic and Arthur G Samuel. Generalization in perceptual learning for speech. *Psychonomic bulletin & review*, 13(2):262–268, 2006.
- Tanya Kraljic and Arthur G Samuel. Perceptual adjustments to multiple speakers. *Journal of Memory and Language*, 56(1):1–15, 2007.
- Russell V. Lenth. *emmeans: Estimated Marginal Means, aka Least-Squares Means*, 2022. URL <https://CRAN.R-project.org/package=emmeans>. R package version 1.8.3.
- Leigh Lisker and Arthur S Abramson. A cross-language study of voicing in initial stops: Acoustical measurements. *Word*, 20(3):384–422, 1964.
- Leigh Lisker and Arthur S Abramson. The voicing dimension: Some experiments in comparative phonetics. In *Proceedings of the 6th International Congress of Phonetic Sciences*, pages 563–567, 1970.
- Sahil Luthra, Hannah Mechtenberg, and Emily B Myers. Perceptual learning of multiple talkers requires additional exposure. *Attention, Perception, & Psychophysics*, 83:2217–2228, 2021.
- Charles L Nagle and Melissa M Baese-Berk. Advancing the state of the art in l2 speech perception-production research: Revisiting theoretical assumptions and methodological practices. *Studies in Second Language Acquisition*, 44(2):580–605, 2022.
- Jonathan Peirce, Jeremy R Gray, Sol Simpson, Michael MacAskill, Richard Höchenberger, Hiroyuki Sogo, Erik Kastman, and Jonas Kristoffer Lindeløv. Psychopy2: Experiments in behavior made easy. *Behavior research methods*, 51:195–203, 2019.
- Janet B Pierrehumbert. Phonological representation: Beyond abstract versus episodic. *Annual Review of Linguistics*, 2:33–52, 2016.
- Robert F Port and Jonathan Dalby. Consonant/vowel ratio as a cue for voicing in english. *Perception & Psychophysics*, 32:141–152, 1982.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2022. URL <https://www.R-project.org/>.
- Eva Reinisch and Lori L Holt. Lexically guided phonetic retuning of foreign-accented speech and its generalization. *Journal of Experimental Psychology: Human Perception and Performance*, 40(2):539, 2014.
- Sabrina K Sidaras, Jessica E. D. Alexander, and Lynne C Nygaard. Perceptual learning of systematic variation in Spanish-accented speech. *The Journal of Acoustical Society of America*, 125(5):3306–3316, 2009.
- Michael D Tyler, Catherine T Best, Alice Faber, and Andrea G Levitt. Perceptual assimilation and discrimination of non-native vowel contrasts. *Phonetica*, 71(1):4–21, 2014.
- Christina Y Tzeng, Jessica E. D. Alexander, Sabrina K Sidaras, and Lynne C Nygaard. The role of training structure in perceptual learning of accented speech. *Journal of Experimental Psychology: Human Perception and Performance*, 42(11):1793, 2016.
- Christina Y Tzeng, Marissa L Russell, and Lynne C Nygaard. Attention modulates perceptual learning of non-native-accented speech. *Attention, Perception, & Psychophysics*, 86(1):339–353, 2024.
- M Luisa Castañeda Vicente. El VOT de las oclusivas sordas y sonoras españolas. *Estudios de fonética experimental*, pages 91–110, 1986.
- Navin Viswanathan, Annie J Olmstead, and M Pilar Aivar. The use of vowel length in making voicing judgments by native listeners of English and Spanish: Implications for rate normalization. *Language and Speech*, 63(2):436–452, 2020.
- Travis Wade, Allard Jongman, and Joan Sereno. Effects of acoustic variability in the perceptual learning of non-native-accented speech sounds. *Phonetica*, 64(2-3):122–144, 2007.
- Douglas H Whalen, Arthur S Abramson, Leigh Lisker, and Maria Mody. F0 gives voicing information even with unambiguous voice onset times. *The Journal of the Acoustical Society of America*, 93(4):2152–2159, 1993.
- Douglas H Whalen, Wei-Rong Chen, Mark K Tiede, and Hosung Nam. Variability of articulator positions and formants across nine english vowels. *Journal of phonetics*, 68:1–14, 2018.
- Lee Williams. The voicing contrast in Spanish. *Journal of phonetics*, 5(2):169–184, 1977.
- Marijt J Witteman, Andrea Weber, and James M McQueen. Foreign accent strength and listener familiarity with an accent codetermine speed of perceptual adaptation. *Attention, Perception, & Psychophysics*, 75(3):537–556, 2013.
- Kevin JP Woods, Max H Siegel, James Traer, and Josh H McDermott. Headphone screening to facilitate web-based auditory experiments. *Attention, Perception, & Psychophysics*, 79:2064–2072, 2017.
- Xin Xie and T Florian Jaeger. Comparing non-native and native speech: Are l2 productions more variable? *The Journal of the Acoustical Society of America*, 147(5):3322–3347, 2020.
- Xin Xie and Emily B Myers. Learning a talker or learning an accent: Acoustic similarity constrains generalization of foreign accent adaptation to new talkers. *Journal of Memory and Language*, 97:30–46, 2017.
- Xin Xie, Rachel M Theodore, and Emily B Myers. More than a boundary shift: Perceptual adaptation to foreign-accented speech reshapes the internal structure of phonetic categories. *Journal of Experimental Psychology: Human Perception and Performance*, 43(1):206–217, 2017.
- Xin Xie, Kodi Weatherholtz, Larisa Bainton, Emily Rowe, Zachary Burchill, Linda Liu, and T Florian Jaeger. Rapid adaptation to foreign-accented speech and its transfer to an unfamiliar talker. *The Journal of the Acoustical Society of America*, 143(4):2013–2031, 2018.
- Xin Xie, Linda Liu, and T Florian Jaeger. Cross-talker generalization in the perception of nonnative speech: A large-scale replication. *Journal of Experimental Psychology: General*, 2021.
- Xiaojuan Zhang, Bing Cheng, and Yang Zhang. The role of talker variability in nonnative phonetic learning: A systematic review and meta-analysis. *Journal of Speech, Language, and Hearing Research*, 64(12):4802–4825, 2021.