

1 Phonetic category formation is perceptually driven during the early stages of adult L2
2 development

3 Joseph V. Casillas¹

4 ¹ Rutgers University

5 Author Note

6 I would like to thank Miquel Simonet, Yukari Hirata, Joseph Kern, and three
7 anonymous reviewers for their candid comments and suggestions on an earlier version of
8 this manuscript. All errors are mine alone. This work was partially funded by the
9 Comanche Nation Higher Education Grant (N4m4n44), as well as the University of
10 Arizona Graduate and Professional Student Council Research and Project Grant
11 (RSRCH-702FY15).

12 Correspondence concerning this article should be addressed to Joseph V. Casillas,
13 Rutgers University - Department of Spanish and Portuguese, 15 Seminary place, New
14 Brunswick, New Jersey, 08904, USA . E-mail: joseph.casillas@rutgers.edu

Abstract

15

16 Research on the acquisition of L2 phonology in sequential language learners has stressed
17 the importance of language use and input as a means to accurate production and
18 perception; however, the two constructs are difficult to evaluate and control. This study
19 focuses on the role of language use during the initial stages of development of phonetic
20 categories related to stop voicing and analyzes the relationship between production and
21 perception. Native English speaking late learners of Spanish provided
22 production/perception data on a weekly basis throughout the course of a 7-week immersion
23 program in which L1 use was prohibited. The production/perception data were analyzed
24 using generalized linear mixed effects models. Generalized additive mixed models were
25 used to analyze and compare the learning trajectories of each modality. The analyses
26 revealed phonetic learning in both production and perception over the course of the
27 program. Perception gains paralleled those of native bilinguals by the conclusion of the
28 program and preceded production gains. This study is novel in that it provides
29 production/perception data in a semi-longitudinal design. Moreover, the beginning adult
30 learners are examined in a learning context in which L1 use was minimal and L2 input was
31 maximized. Taken together, the experiments suggest that L2 phonetic category formation
32 can occur abruptly, at an early stage of development, is perceptually driven, and appears
33 to be particularly fragile during the initial stages of learning.

34

Keywords: Sequential language learning, SLA, Production, Perception, Stop voicing

35

Word count: 12644

36 Phonetic category formation is perceptually driven during the early stages of adult L2
37 development

38 **Introduction**

39 A common finding in the second language (L2) speech literature is that adults who
40 learn another language typically retain a non-native accent (Caramazza et al., 1973; Flege,
41 1981, 1987a; Fowler et al., 2008; Oyama, 1976; Pallier et al., 1997; Sundara & Polka, 2008,
42 among others). The phonetic consequences of sequential language learning—acquiring an
43 L2 after, rather than at the same time as, the L1—are traditionally associated with speech
44 production. Some researchers, however, refer to L2 learners as perceiving speech with an
45 accent as well (Escudero, 2005; Strange, 1995). That is to say, L2 speech perception can
46 also differ from native listening. There is a dearth of research regarding the relationship
47 between production and perception in adult L2 learning and how the two modalities are
48 affected by L2 input and L2 use.

49 The present work is concerned with understanding how late learners manage to
50 acquire L2 sound categories and the nature of their development in reference to input and
51 L1/L2 use. Additionally, this work explores the interface between speech perception and
52 speech production in beginning adult learners. An L2 can be learned formally (i.e., in an
53 L2 classroom), informally (i.e., in a naturalistic context), or both formally and informally
54 in an immersion type context (Saville-Troike, 2005). The present work focuses on the L2
55 learning that takes place in the latter. Data were collected from learners in a domestic
56 immersion context—i.e., foreign language immersion in their country of origin, the U.S.—in
57 which they were required to minimize the use of their L1 and received large amounts of L2
58 input. Specifically, this work examines the initial stages of L2 production and perception in
59 a group of adult late learners that took part in a Spanish domestic immersion program in
60 which L1 use was prohibited.

Background

Input and use in L2 learning

While age effects have duly garnered the attention of second language acquisition (SLA) researchers, we know now that early exposure alone cannot explain L2 outcomes in all cases (See Pallier et al., 1997; Sebastián-Gallés & Soto-Faraco, 1999; Sebastián-Gallés et al., 2005; Sebastián-Gallés, 2006, among others). A possible explanation may lie in the nature of the input learners receive. For instance, in an investigation of the acquisition of English /d/-/ð/, Sundara, Polka, & Genesee (2006) determined that variable realizations of English /ð/ in the input of young French/English bilinguals may have delayed their acquisition of a “functional” /d/-/ð/ contrast (p. 382). This suggests that the difficulties in perception/production of some learners might be a reflection of the input to which they are exposed.

The role of input in the production/perception of adult learners has also been studied, albeit to a lesser extent. For instance, Flege & Liu (2001) investigated length of residency (LOR) and input in a group of 60 Chinese late learners of English. Flege & Liu (2001) found that LOR was a crucial factor regarding L2 outcomes for late learners, but only if they needed to use English regularly. It remains unclear if this was due to the quality of the input (from native speakers), the quantity, or a combination of both. Further complicating the issue, L2 input is likely associated with other factors, such as L2 feedback, motivation, and attention, to name just a few.

In comparison with input, the role of L2 use in the production/perception of adult learners has received less attention. There is, however, an illustrative example in a series of foreign accent studies conducted by Flege and colleagues. Flege et al. (1995) found that the earlier the age of learning (AOL), the more native-like the participants’ production/perception. In an unpublished follow-up to this study, Flege & MacKay revisited the learners from the original investigation 10 years later (as cited in Flege, 2012).

87 The researchers found that increased use of English was associated with more native-like
88 production when compared with the 1992 data. A decrease in English usage resulted in no
89 change in production accuracy. In a separate unpublished longitudinal study, Flege found
90 that after a period of 5 years in the U.S., foreign accent ratings of Spanish-speaking late
91 learners of English showed no improvement (as cited in Flege, 2012). However, post-hoc
92 scrutiny of the data compared the 3 “worst” (more foreign accented) learners with the 3
93 “best” (least foreign accented) learners and showed that the “best” learners reported using
94 English more. Specifically, they used English in contexts where it could be considered
95 optional (i.e., in conversations with friends), suggesting that L2 use in extensive contexts
96 may foment L2 phonological development.

97 In sum, the role of input and use in adult language learning has not been a primary
98 object of focus in the SLA literature despite the crucial status these factors are given in L1
99 acquisition. This is likely the result of both variables being overshadowed by research on
100 age effects, coupled with the fact that they are difficult to control and there is no
101 straightforward method for quantifying them. Language immersion provides an ideal point
102 of comparison for studying use and input because learner access to native speakers is high
103 and the target language is likely used often. Moreover, this learning context opens the door
104 to studying L2 phonological acquisition during the initial stages of learning. The present
105 study contributes to this literature by examining adults in a domestic immersion program
106 in which L2 use is maximized and L2 input is rampant.

107 **Production/perception interface in L2 learning**

108 The relationship between production and perception is assumed to be
109 non-controversial, as there are numerous studies demonstrating the relationship between
110 the two modalities (i.e. Flege & Eefting, 1988; Williams, 1979; Flege et al., 1999, among
111 many others). However, there are discrepancies in the literature regarding which of the two
112 modalities is the driving force in L2 learning.

113 On one hand, a long line of research supports the claim that perception of a novel
114 phonetic segment precedes its production (Williams, 1979; Borden et al., 1983; Neufeld,
115 1988; Barry, 1989; Grasseger, 1991; Flege, 1993; Rochet, 1995; Llisterri, 1995; Flege et al.,
116 1997; Leather, 1999). For instance, Flege et al. (1997) explored the production/perception
117 of /i/-/ɪ/ in Spanish-speaking learners of English and found that a subset of the
118 participants perceived the vowels similarly to a group of native controls. Crucially, only a
119 few of the native-like perceivers were able to produce the /i/-/ɪ/ contrast accurately. Flege
120 et al. (1997) took these findings as evidence that accurate perception must necessarily
121 precede accurate production. Further support for the perception-first hypothesis comes
122 from training studies. For example, some studies on perceptual training have lead to more
123 accurate production (See Rochet, 1995; Bradlow et al., 1997, 1999).

124 On the other hand, there are also investigations that cast doubt on the
125 perception-first claim (Goto, 1971; Caramazza et al., 1973; Sheldon & Strange, 1982; Mack,
126 1989; Mathews, 1997; Leather, 1997; Wang, 2002; Kissling, 2014). For instance, Sheldon &
127 Strange (1982) examined the production/perception of American English /l/-/r/ in
128 Japanese learners and found that they demonstrated more accuracy in production than in
129 discrimination. The authors contend that production accuracy may precede perceptual
130 accuracy for a given segment due to pedagogical reasons, as some of the participants
131 received instruction regarding the contrast in question based on articulatory parameters (as
132 opposed to auditory cues). Mack (1989) obtained similar results in English/French
133 bilinguals and hypothesized that there might be greater consequences associated with
134 mispronouncing a segment versus misperceiving it. Alternatively, the abundance of viable
135 acoustic cues in the speech signal might make sufficient perception possible without
136 monolingual-like mastery of any single cue.

137 One method of exploring the nature of the production/perception relationship is via
138 longitudinal data. In a semi-longitudinal analysis of Spanish /p/-/b/, Zampini (1998) and
139 Zampini & Green (2001) tested intermediate/advanced learners on 3 separate occasions

140 during the course of an academic semester. These studies found that perception changes
141 occurred after production changes, though one cannot reach a definitive conclusion from
142 these data regarding the nature of the relationship between the two modalities because the
143 participants in question were fairly advanced and, at the time of testing, already differed
144 from monolinguals in their perception of Spanish stops. Longitudinal data from absolute
145 beginners are necessary to determine which of the two modalities occurs first in L2
146 phonological development.

147 **L2 speech models and phonetic category development**

148 There are numerous L2 models used to posit hypotheses that can account for the
149 difficulties encountered by L2 learners. Relevant to the present work are the Speech
150 Learning Model (SLM, Flege, 1995) and the Linguistic Perception Model (L2LP, Escudero
151 & Boersma, 2004; Escudero, 2005, 2009). The SLM maintains that phonetic similarity or
152 perceived equivalence can predict L2 difficulties. According to the SLM, the human ability
153 to learn novel sounds is maintained throughout life. Representations of the L2 sounds
154 being acquired are stored in long-term memory and share a common phonological space
155 with the L1. The SLM proposes that bilinguals aim to maintain L1 and L2 phonetic
156 categories separate. This implies that L1 to L2, as well as L2 to L1, interactions occur,
157 with each language having some influence on the other. Formation of novel L2 categories
158 becomes more difficult—but not impossible—as the L1 sound system develops. If L1 and
159 L2 sounds are too perceptually similar, category formation is hindered because learners
160 perceive the L2 sound as being equivalent to the L1 category.

161 For its part, the L2LP also maintains that the human ability to learn novel speech
162 sounds remains active throughout life and that L2 difficulties are accounted for via
163 phonetic similarities, differences, or perceived equivalences to native contrasts. The L2LP
164 posits two distinct situations for learning L2 contrasts. A contrast can be considered novel
165 (“new scenario”) to the learner, in which case the model contends that new categories must

166 be formed. Conversely, a contrast can be familiar (“similar scenario”), and (s)he must then
167 reset the boundary for the contrast via a comparison module, the Gradual Learning
168 Algorithm (GLA).

169 Under the L2LP, a new L2 grammar is created via *full copying/full access*, and is
170 developed independently of the L1 perception grammar. Thus L2 learners have the ability
171 to become native-like in their perception of the L2 without affecting the L1 perception
172 grammar. This notion is at odds with other popular speech models, namely the SLM,
173 which proposes that L1 and L2 categories share the same phonetic space and therefore L2
174 development is hypothesized to occur simultaneously with changes in L1 categories. The
175 L2LP, on the other hand, proposes that L2 learning occurs as a result of resolving two
176 problems. The first problem deals with the fact that the perception grammar must change,
177 or adjust, in order to manage input from the L2, and is, thus, a perceptual problem. The
178 second problem, representational in nature, requires that new L2 categories be created
179 whenever L1 categories cannot be used. Similar contrasts are hypothesized to be easier to
180 learn because listeners can simply reuse L1 categories. Phonetic differences are accounted
181 for via boundary adjustments. Adjustments to the L2 perception grammar also occur
182 through the GLA. While the SLM and the L2LP clearly differ regarding how phonetic
183 category formation occurs, both models assume an underlying relationship between
184 production and perception. According to Flege, L2 phonetic segments are “[...] produced
185 only as accurately as they are perceived” (Flege, 2003, p. 25). The L2LP formalizes this
186 relationship by positing a specific perceptual grammar that works in conjunction with a
187 production grammar. Escudero (2005) maintains that the model is best utilized with
188 longitudinal data collected from beginning language learners as they progress over time.
189 Accordingly, the present work examines the production/perception of pure beginners
190 during the initial stages of L2 learning.

191 **Spanish and English stop contrasts**

192 The voicing contrasts of Spanish and English stops serve as a proxy for analyzing the
193 production/perception relationship and the roles of input and use during early L2
194 phonological development in the present study. Stops can vary based on their point of
195 articulation and voicing. Most, but not all languages, have voiceless stops, and, of these
196 languages, many have voicing contrasts. While some languages have three-way contrasts
197 for stops (e.g., Thai), most have a two-way voicing contrast. Two-way contrast languages
198 arbitrarily fall into one of two categories: true voicing or aspirating languages. English and
199 Spanish are languages of the two-contrast variety. English is an aspirating language;
200 Spanish is a true voicing language.

201 Importantly, these languages share the same two-way contrasts—phonologically
202 voiced versus voiceless phones at bilabial, coronal, and velar place—, but differ in where
203 the acoustic boundary lies between them based on voice onset time (VOT, Lisker &
204 Abramson, 1964). VOT refers to the duration of the time interval between the release of
205 the stop burst and the onset of modal voicing. English, an aspirating language, contrasts
206 /b, d, g/ with /p, t, k/ through phonetically voiceless stops with short-lag (positive) VOT
207 and phonetically aspirated stops with long-lag (positive) VOT, respectively.¹ Spanish, a
208 true voicing language, contrasts the same phonological segments through phonetically
209 prevoiced stops with lead (negative) VOT and phonetically voiceless stops with short-lag
210 (positive) VOT (Lisker & Abramson, 1964).

211 These differences may appear trivial, however, they have important implications with
212 regard to speech production/perception. L2 learners are likely to produce and perceive L2
213 stops using the acoustic boundaries of their L1 (Flege, 1987b). As productions in the L2
214 deviate from native values, they are more likely to be perceived as sounding foreign. Thus,

¹ In English phonologically voiced stops can surface partially or fully voiced in prosodically weak positions in connected speech (Davidson, 2018). This is less common in utterance initial position.

215 accurately producing the stop contrasts of a foreign language can significantly improve
216 foreign accent ratings (Sundara, Polka, & Baum, 2006). For example, if a native English
217 speaker learning Spanish produces the long-lag VOT associated with English /p/
218 (i.e. aspiration, [p^h]) when saying the Spanish word *papel* (“paper”)—which is realized as
219 short-lag [p]—, the production would sound foreign to the native Spanish speaker (Elliott,
220 1997; González-Bueno, 1997; Sundara, Polka, & Baum, 2006). The issue is also complex for
221 the L2 learner of Spanish regarding perception. Given that both languages have short-lag
222 VOT in phonologically distinct segments, the L2 learner can—and often does—mistakenly
223 “hear” an English /b/ category when the utterance, in reality, contained Spanish /p/. This
224 could be the difference between “hearing” *peso* (“I weigh”) versus *beso* (“I kiss”).

225 The task of English-speaking L2 learners of Spanish is to associate short-lag VOT
226 with Spanish voiceless /p, t, k/, and create a new category for lead VOT associated with
227 voiced /b, d, g/ in order to produce/perceive Spanish stops accurately. Recall that both
228 the SLM and the L2LP maintain that L2 learners develop new phonetic categories. Under
229 the L2LP account, learning the voicing contrasts of Spanish involves boundary resetting
230 (i.e., similar scenario) and is hypothesized to be relatively easy for the learner. The SLM,
231 on the other hand, would posit that Spanish stops might pose a challenge to the L2 learner
232 of Spanish because of the phonetic similarity between the corresponding English and
233 Spanish phones. Accordingly researchers have attempted to demonstrate phonetic category
234 formation by examining the acoustic properties of L2 speech and comparing them to those
235 of native speakers.

236 For example, Flege & Eefting (1987) examined Spanish speakers’ production of word
237 initial English stops and found that they produced the English segments with aspiration.
238 However, the acoustic analyses showed that the VOT values were longer than in Spanish
239 words, but not as long as the VOT values of a monolingual English control group. Flege &
240 Eefting (1987) concluded that the Spanish/English bilinguals did indeed develop L2
241 phonetic categories, but their productions were not “authentic” because they still differed

242 from those of the monolingual controls. They ascribed these differences to the nature of the
243 input their participants had received, i.e., Spanish-accented English. The aforementioned
244 study suggests that learners are capable of acquiring new phonetic categories, though the
245 nature of these categories may differ from those of monolingual speakers.

246 Previous research on the acquisition of Spanish stops has included traditional
247 classroom learners (Zampini, 1998; Zampini & Green, 2001; López, 2012; González López
248 & Counselman, 2013; Nagle, 2017), as well as study abroad (SA) immersion learners
249 (Stevens, 2001; Díaz-Campos, 2006; Nagle et al., 2016). For instance, a study by Stevens
250 (2001) compared the production of /p, t, k/ of a group of North American English speakers
251 participating in an immersion program in Spain to a group of traditional classroom room
252 learners. Stevens (2001) found that the SA group learned to reduce aspiration in their
253 Spanish productions, and the traditional classroom learners continued producing Spanish
254 voiceless stops with long-lag VOT. Moreover, Stevens (2001) noted a linear relationship
255 between the length of stay in the foreign country and production accuracy and suggested
256 that the extensive use of Spanish likely accounted for the SA groups' phonetic category
257 development. In a similar vein, Díaz-Campos (2006) found that SA learners showed greater
258 improvements than a traditional classroom control group in the production of /p, t, k/
259 when their speech was analyzed in a conversational style. In more formal settings, acoustic
260 analyses showed no differences between the same groups. There is also evidence of category
261 formation in traditional classroom learners. For instance, in a semi-longitudinal study of
262 Spanish bilabial stop production, Zampini (1998) showed that a group of North American
263 intermediate/advanced late learners produced Spanish voiceless stops with shorter VOT
264 than their English stops, but longer VOT than native Spanish speakers (see also Zampini
265 & Green, 2001, for other cues). The learners did not incorporate prevoicing into their
266 productions of /b/ by the end of the semester.

267 Summarizing, the literature on the adult acquisition of Spanish stops has focused on
268 the traditional classroom and immersion contexts. The findings suggest that adults can

269 form new phonetic categories for Spanish voiceless stops—especially in a learning context
270 that facilitates L2 use and provides ample native input—, but voiced stops may take longer
271 to develop. There is a gap in the SLA literature regarding the initial stages of adult L2
272 phonological acquisition, and in particular regarding Spanish stops.

273 **The present study**

274 To investigate phonetic category development during the early stages of L2 learning,
275 the present work tracked the production/perception of 10 adult native English speakers
276 who participated in a Spanish immersion program. The goals of this study were to (1)
277 explore phonetic category development related to the fine-phonetic detail of Spanish stop
278 voicing contrasts, and (2) determine whether L2 phonological development is driven by
279 production or perception. Furthermore, this research contributes to the L2 literature
280 regarding language input and use by providing semi-longitudinal data from beginning adult
281 learners in a context in which L2 use and input were maximized. Conversely, L1 use was
282 held to a minimum.

283 **General method**

284 **Questionnaires**

285 The learners completed a language background questionnaire and a weekly progress
286 questionnaire. The former excluded participants who had experience learning Spanish or
287 other languages. The questions inquired about time spent studying an L2, living in a
288 foreign country, if they had family members who spoke a language other than English, and
289 provided an overall assessment of their Spanish. The second questionnaire was a self-report
290 assessment aimed at quantifying information related to weekly language use and input,
291 along with other measures. The questionnaire asked about the participants' use of
292 Spanish/English, their self-reported speaking abilities, listening/comprehension abilities,

293 overall abilities in Spanish, and if they felt their Spanish had improved. The learners
294 completed this questionnaire every Sunday before participating in the experimental tasks.
295 Thus responses each week referred to what they had done during the previous week. The
296 questions related to their use of Spanish provided a percent estimate of time spent
297 speaking Spanish with native speakers, and non-native speakers with more, less or equal
298 proficiency. A third questionnaire, the Bilingual Language Profile (BLP, Birdsong et al.,
299 2012), was administered to a bilingual control group. The BLP provided a measure of
300 language dominance by calculating scores for 4 modules: language history, language use,
301 language proficiency, and language attitudes. The summed scores from each module
302 provide a value of language dominance that ranges from -218 to 218. A score near either
303 extreme indicates dominance in one of the two languages. Negative scores were associated
304 with dominance in English, and positive scores were associated with dominance in Spanish.
305 A score near 0 indicated balanced bilingualism.

306 **Participants**

307 The present study included 20 people who were divided into 2 groups. The first
308 group consisted of 10 adult L2 learners of Spanish whose native language was English. The
309 second was a control group comprised of 10 simultaneous Spanish-English bilinguals.²

² The use of bilingual control groups, as opposed to the more traditional practice of using monolingual speakers, is part of a recent trend in L2 phonetic research (see Sakai, 2018). Numerous investigations show even early bilinguals with ample L2 proficiency tend to differ from monolinguals in production, perception, and lexical processing (i.e., Pallier et al., 1997; Sebastián-Gallés & Soto-Faraco, 1999; Sebastián-Gallés et al., 2005; Sebastián-Gallés, 2006, among others). These differences are often ascribed to cross-linguistic interactions, which, simply put, may well be a part of bilingualism, independent of age of acquisition and L1/L2 proficiency. The present work takes the position that individuals undertake the endeavor of language learning with the goal of becoming bilingual and not to replace their native language. Therefore a more fair and useful assessment of their progress can be achieved by comparing their abilities to those of a population (bilinguals) that represents their end goal (bilingualism).

310 **Late learners.** The learners of the present study were students in a domestic
311 immersion language program at Middlebury College. The defining characteristic of the
312 program is the Language Pledge, a formal agreement the students signed by which they
313 promised to use only the target language (in this case Spanish) for 7 weeks. Failure to
314 comply with the pledge can result in expulsion. Students lived in the residence halls on the
315 campus with other students and professors, and they attended class for 4 hours in the
316 morning and participated in co-curricular activities in the afternoon. The program is
317 designed with the intention of creating an experience comparable to living abroad, though
318 it is considered intense immersion due to the pledge and the seriousness of the students and
319 the curriculum. The classes employed a communicative focus with instruction entirely in
320 Spanish.

321 Information from the background questionnaire was used to select 10 participants (4
322 males, 6 females) who reported no prior experience with any other languages. Individuals
323 who had completed a semester or more of foreign language study or had spent time living
324 in a foreign country were excluded. The late learners were 18 years old or older ($\bar{x} = 23.70$;
325 $SD = 5.27$), and considered absolute beginners. This assertion was confirmed via
326 placement testing and an interview with two faculty members of the immersion program.

327 Table 1 displays a summary of the self-report data obtained from the weekly
328 assessment questionnaire. The learner group used their L1, English, minimally, and their
329 L2, Spanish, almost exclusively. The L2 input to which the learner group was exposed was
330 provided mainly by native speakers or non-native speakers with higher levels of proficiency
331 in Spanish. Furthermore, the learners believed their listening, speaking, and overall
332 abilities in Spanish improved over the course of the program.

333 ++ INSERT TABLE 1 ABOUT HERE ++

334 **Simultaneous (native) bilinguals.** Ten simultaneous (i.e., native) bilinguals
335 participated in the present study. These native Spanish and English speakers served

336 primarily as a control group with which the production/perception of the learner group
337 was compared. Bilingual participants reported speaking both English and Spanish for as
338 long as they could remember, and that their parents were also bilingual. Moreover,
339 bilingual participants stated they used both languages on a daily basis with friends and
340 family. The mean language dominance score of the bilingual group was -2.76 (SD = 38.93),
341 suggesting balanced bilingualism according to the BLP.

342 **Overview of procedures**

343 The learners completed four distinct tasks: two related to their production in Spanish
344 and two related to their perception in Spanish. The present study reports one of the
345 production tasks and one of the perception tasks. The purpose of the tasks was to provide
346 data measuring their progress learning Spanish bilabial stops. During the initial session,
347 the learners completed the first iteration of the assessment questionnaire, a delayed
348 repetition production task, and a two-alternative forced-choice perception task (2AFC).
349 From this point forward until the final week of the program, the learners completed the
350 same tasks with the exception that the delayed repetition production task was replaced
351 with a reading production task. Experimental sessions took place every Sunday with the
352 exception of the last week, which included multiple sessions and two tasks not reported
353 here. Table 2 presents an overview of the learners' participation.

354 ++ INSERT TABLE 2 ABOUT HERE ++

355 Bilingual participants were recruited from a university in the Southwestern United
356 States. Their participation included two days of testing at a speech science laboratory with
357 a minimum of 24 hours between sessions. On the first day they completed the BLP
358 questionnaire, the 2AFC task, and the reading task. On the second day they completed
359 tasks that are not reported here. The following sections report the results of the production
360 task, proceeded by the perception task, and, finally, a comparison of the two modalities.

Longitudinal development of L2 bilabial stop production

The first task examined the ongoing development of Spanish bilabial stop production in adult L2 learners with a special focus on how the realization of stop voicing changed with increased L2 exposure.

Method

Materials.

Target phrases.

Participants repeated a series of words containing Spanish stops, /p, t, k, b, d, g/, in utterance initial position. A total of 30 nonce words—5 for each stop segment—were embedded in the carrier phrase “_____ es la palabra” (*Eng.* “_____ is the word”). The syllable structure was CV.CV with primary stress falling on the initial syllable. Stops were followed by one of the 5 Spanish vowels, /i, e, a, o, u/, and the consonant /k/. The final vowel was always the same as the first (i.e., /'a.ka/, /'e.ke/, /'i.ki/, etc.). There were 6 stops (/p, t, k, b, d, g/) \times 5 vowels (/i, e, a, o, u/) = 30 items. The present work focused on the bilabials, /p, b/. Target words were interspersed amongst 20 distractors used in another task.

Auditory stimuli.

A 29 year old native female Spanish speaker from Cádiz, Spain, provided the audio stimuli presented in the delayed repetition portion of the task. A Shure SM10A dynamic head-mounted microphone recorded the items. A Sound Devices MM-1 pre-amplifier boosted the signal and sent it to a laptop computer where it was recorded using Praat at a 44.1 kHz sample rate with 16-bit quantization (Boersma & Weenink, 2018). The recording took place in a sound attenuated booth in a phonetics laboratory at a university in the U.S. southwest.

385 **Procedure.**

386 *Recordings.*

387 The learners were recorded in a quiet classroom on site at Middlebury College. A
388 Shure SM10A dynamic head-mounted microphone captured the participants' productions.
389 A Sound Devices MM-1 pre-amplifier passed the signal to a Marantz PMD661 MKII
390 Handheld Solid State broadcast recorder. Recordings were sampled at 44.1 kHz with 16-bit
391 quantization. The same setup used to record the auditory stimuli served to record the
392 productions of the bilinguals.

393 *Acoustic analysis.*

394 Participants' productions were segmented in Praat using synchronized waveform and
395 spectrographic displays to hand-mark the onset of voicing and the burst for each stop.
396 Voicing onset was the first periodic pattern found in the waveform. The criterion for bursts
397 was the onset of broad-band sudden noise in the spectrogram. A Praat script
398 automatically extracted VOT, which was calculated as the difference (in ms) between the
399 aforementioned acoustic landmarks (i.e. onset of modal voicing and the burst). Figure 1
400 illustrates the segmenting procedures.

401 ++ INSERT FIGURE 1 ABOUT HERE ++

402 The initial experimental session took place on the second or third day of the
403 immersion program. Thereafter, data collection took place every Sunday until the end of
404 the program for a total of 8 experimental sessions. The final session took place on a
405 Wednesday, which was 3 days after the seventh session and 3 days before the program
406 ended. PsychoPy2 (Peirce, 2008) presented the stimuli randomly. The first session
407 employed a delayed shadowing technique. The stimuli were presented aurally via the audio
408 recordings of the native Spanish speaker. Participants listened and repeated the 50 items
409 embedded in the carrier phrase. Subsequent experimental sessions utilized a reading task.

410 In these cases, the stimuli appeared on a computer screen and participants read them
411 aloud.³ After saying the stimuli aloud, participants pressed a button on a keypad to
412 advance to the next item. Each learner provided the dataset with 240 bilabial stops (2
413 stops \times 5 vowels \times 3 repetitions \times 8 sessions). Thus, a total of 2,400 stop tokens were
414 collected from the learner group (240 tokens \times 10 participants = 2,400). The bilingual
415 controls only completed one session of the reading task, and provided a total of 300
416 utterance initial stops (2 stops \times 5 vowels \times 3 repetitions \times 10 participants). The task
417 took approximately 10 minutes.

418 **Statistical analyses.** Data from the production task were analyzed using a series
419 of generalized linear mixed effects models (GLMM). Specifically, there were 3 analyses.
420 The first analysis aimed to see if learners' production of Spanish bilabial stops changed by
421 end of the immersion program, and how self-reported measures of input and use affected
422 their progress. In this model change in VOT in standardized units (ΔZ_{VOT}) was the
423 criterion. To calculate ΔZ_{VOT} , raw VOT values were converted to z-scores (i.e.,
424 standardized) as a function of voicing. For each participant, the difference in standardized
425 VOT from the start of the program and the end of the program was calculated. The result
426 was a value indicating whether the VOT of learners' bilabial productions had reduced
427 (negative ΔZ_{VOT}), increased (positive ΔZ_{VOT}), or remained the same (ΔZ_{VOT} near 0)
428 after the program. Fixed effects were averaged self-report assessments of input (% of input
429 from the following sources: (1) native, (2) non-native, (3) non-native with higher
430 proficiency, (4) non-native with lower proficiency, and (5) non-native with equal
431 proficiency) and use (% time speaking Spanish, English), with by-subject random effects on
432 all continuous predictors and item repetitions. The fixed effects were standardized, thus all

³ Two separate techniques were used in the production task (delayed repetition during week 0, reading during the remaining weeks) in order to avoid intervocalic 't' being realized as [r] in a subset of the distractors that served as stimuli for another experiment. Participants did not produce intervocalic 't' as [r] in any of the experimental sessions.

433 continuous predictors had a mean of 0 and a standard deviation of 1.

434 The purpose of the second analysis was to determine the amount of exposure time
435 necessary for the learners' production of Spanish bilabial stops to change. Production data
436 for /b/ and /p/ were fit separately. VOT was the criterion and *exposure time* (day 0, day
437 7, day 21, day 28, day 35, day 42) was a fixed effect. The random effects structure included
438 by-subject and by-item intercepts with slopes for exposure time (for the subject effect) and
439 item repetitions (for the items effect). Exposure time was dummy coded with day 0 set as
440 the reference level. Thus all tests of simple effects compared VOT values after a given
441 amount of exposure (i.e., day 7, day 14, etc.) to the baseline.

442 The final analysis directly compared the learners' production of Spanish bilabial stops
443 at the end of the immersion program (day 42) with that of the bilingual control group.
444 VOT was again the criterion and the regressors were *group* (learner, native bilingual) and
445 *voicing* (voiced, voiceless). The model included by-subject and by-items intercepts with
446 random slopes for voicing and item repetitions. The fixed effects were again dummy coded
447 with native bilinguals' voiceless stops set as the reference levels.

448 For all models, visual inspection of Q-Q plots and plots of residuals against fitted
449 values were used to check for normality of the residuals. Unless noted otherwise, statistical
450 significance of main effects and higher order variables was assessed using hierarchical
451 partitioning of the variance via nested model comparisons. Marginal R^2 and conditional R^2
452 provided an indication of goodness-of-fit for each model (Nakagawa & Schielzeth, 2013).
453 Marginal R^2 specified a measure of variance explained without random effects and
454 conditional R^2 included them.

455 Results

456 The role of input and use in bilabial stop production.

457 ++ INSERT FIGURE 2 ABOUT HERE ++

458 Panels (a) and (b) of Figure 2 plot the VOT production data as a function of days of
 459 exposure, voicing and group (learner, native bilingual). The first analysis examined ΔZ_{VOT}
 460 as a function of self-reported input and use factors. The omnibus model tested the
 461 hypothesis that ΔZ_{VOT} was significantly different from 0. The intercept estimate was -0.74
 462 ± 0.10 standard errors (CI low = -0.93 ; CI high = -0.56 ; $t = -7.8$; $p < 0.001$). The negative
 463 estimate indicates a net decrease in VOT for bilabial stop production. Back-transforming
 464 to raw values showed that voiced stops had an average VOT of 13.11 ms and lowered by
 465 -38.49 ms by the end of the program. Voiceless stops were 19.2 ms higher at baseline ($\bar{x} =$
 466 32.31 ms) and decreased by approximately -10.78 ms by the end of the program.
 467 Continuous input and use predictors were included in the model using forward selection
 468 and only retained if they significantly contributed to model fit. There was only one main
 469 effect: self-reported % of English use ($\chi^2(1) = 4.98$; $p < 0.027$). Specifically, a 1-unit
 470 increase of *Z-English use* was associated with an increase in ΔZ_{VOT} of 0.23 ± 0.10
 471 standard errors (CI low = 0.03 ; CI high = 0.43 ; $t = 2.29$; $p < 0.03$). Thus learners who
 472 self-reported higher overall English use showed smaller changes in VOT (See Figure 3).
 473 The model of best fit included random effects ($R^2_m = 0.03$; $R^2_c = 0.50$).

474 ++ INSERT FIGURE 3 ABOUT HERE ++

475 **Change in bilabial stop production over time.**

476 ***Voiceless bilabial stops.***

477 The voiceless bilabial data were best fit using a maximal error term ($R^2_m = 0.07$; R^2_c
 478 $= 0.66$). There was a main effect of session ($\chi^2(6) = 27$; $p < 0.001$). VOT values had
 479 lowered by -8.86 ms ± 3.41 standard errors (CI low = -15.55 ; CI high = -2.18 ; $t = -2.6$; p
 480 < 0.02) after 21 days of exposure. The average VOT difference from the baseline value was
 481 approximately 10 ms from this point forward; however, there was appreciable variability, as
 482 seen in the standard errors of the parameter estimates for each testing session. Table 3

483 provides the complete model output and Figure 2 shows the distributions of the data at
 484 each testing session.

485 ++ INSERT TABLE 3 ABOUT HERE ++

486 ***Voiced bilabial stops.***

487 The analysis showed that the voiced bilabial stop data were also best fit using a
 488 maximal error term ($R^2_m = 0.16$; $R^2_c = 0.67$). There was a main effect of session ($\chi^2(6) =$
 489 29.95 ; $p < 0.001$). VOT was significantly lower than the week 0 initial state after 21 days
 490 of exposure (CI low = -49.89 ; CI high = -13.13 ; $t = -3.36$; $p < 0.001$), and for each session
 491 thereafter. Thus, learner VOT values for voiced bilabials decreased after three weeks in the
 492 program. Table 4 displays the percentage of prevoiced /b/ as a function of exposure time.
 493 Crucially, all learners produced prevoiced stops at least some of the time, and by the
 494 conclusion of the program approximately half of the productions included lead VOT. The
 495 average VOT of the prevoiced stops was lower ($\bar{x} = -75.3$, $SD = 35.3$), suggesting
 496 production of /b/ was inconsistent. The complete model output is displayed in Table 5 and
 497 density ridgeline plots of the distributions for each session are available in Figure 2.

498 ++ INSERT TABLE 4 ABOUT HERE ++

499 ++ INSERT TABLE 5 ABOUT HERE ++

500 **Comparison with bilinguals.** The data were best fit when including the random
 501 effects structure ($R^2_m = 0.54$; $R^2_c = 0.73$). There was no effect of *group* ($\chi^2(1) = 1.26$; p
 502 > 0.05), but there was an effect of *voicing* ($\chi^2(1) = 23.28$; $p < 0.001$), as well as an
 503 interaction between the two factors ($\chi^2(1) = 11.92$; $p < 0.002$). The learner groups'
 504 voiceless stops had a mean VOT value of 22.80 ms, approximately 4.77 ± 3.3 standard
 505 errors higher than the controls, a difference that was not statistically significant (CI low =

506 -1.7; CI high = 11.25; $t = 1.44$; $p > 0.05$). The learner groups' voiced stops, on the other
507 hand, differed from those of the control group by 53.87 ± 13.35 standard errors (CI low =
508 27.7; CI high = 80.04; $t = 4.03$; $p < 0.001$). As shown in the previous analysis, the
509 productions of /b/ that were indeed prevoiced were closer to the bilingual range. The
510 distributions for voiceless and voiced stops of the bilinguals are displayed in the final
511 ridgelines of panels (a) and (b) of Figure 2.

512 **Interim discussion**

513 The first task investigated bilabial stop production in late learners of Spanish. The
514 task was designed to determine (1) if the learner group improved its production of bilabial
515 stops after a 7-week immersion program, and the extent to which input and use factors
516 modulated this improvement, (2) how much exposure was necessary for observable phonetic
517 category development to take place, and (3) how the learners production compared to a
518 group of simultaneous bilinguals upon completion of the immersion program.

519 Upon comparing stop production from the baseline initial state and the final testing
520 session after 42 days of exposure to Spanish, the results of the first analysis suggested that
521 overall the learners reduced VOT for bilabial stops. Specifically, by the end of the program
522 they reduced aspiration for the voiceless stops and began to incorporate prevoicing into
523 their voiced stops. These changes were partially modulated by self-reported use of English,
524 such that participants who reported higher overall use of English during the program
525 showed less improvement in bilabial stop production. The second analysis found that the
526 learners' production boundaries began to shift after 21 days of exposure to Spanish. That
527 is, after the third week in the program, both /p/ and /b/ had lower VOT values, and
528 continued to decrease throughout the remainder of the program. The third analysis
529 compared the learners' stop production in the final week of the program to that of the
530 bilingual control group. Although the learners reduced VOT for both stop segments, they
531 clearly differed from the bilinguals regarding the voiced segment /b/, despite the fact that

532 this segment appeared to show a larger change (in ms) by the end of the program. An
533 analysis of the proportion of prevoiced stops showed that all learners included lead-vot
534 some of the time, and, when they did, the values were much closer to the bilingual range,
535 suggesting /b/ production was particularly unstable. The voiceless stop was produced
536 within the range of native values for VOT. In sum, the first task showed that (i) the
537 learners did improve their stop production in a 7-week immersion program, (ii) use of
538 English affected production gains, and, finally, (iii) evidence of phonetic category
539 development was observable after 21 days of exposure.

540 L2 perception of bilabial stops

541 The second task examined the ongoing development of Spanish stop perception in
542 adult L2 learners with a special focus on how stop voicing identification changed with
543 increased L2 exposure.

544 Method

545 **Materials.** In order to create a VOT continuum a twenty-three year old female
546 Spanish/English simultaneous bilingual from the Southwestern U.S. provided natural
547 productions of the bisyllabic words “bata” (Eng. *robe*) and “pata” (Eng. *paw*), each of
548 which contain stops in utterance initial position. An AKG C520 condenser microphone was
549 used to record the utterances. A Sound Devices USBPre 2 audio interface digitized the
550 signal at 44.1 kHz and 16 bit quantization. The digitized signal was sent to a laptop
551 computer and recorded using Praat (Boersma & Weenink, 2018). The best token of
552 [p^h]—one in which there were no clicks or extraneous noise—was selected for resynthesis.
553 For the stimuli with positive VOT, Praat manipulated the duration of the aspirated
554 portion of the stop via the Time-Domain Pitch-Synchronous-Overlap and Add algorithm
555 (TD-PSOLA). For the stimuli with negative VOT, periods of prevoicing were pasted into
556 the signal at zero-crossings before the release of the stop. The prevoiced portions were

557 taken from phonetically voiced stop productions of the aforementioned simultaneous
558 bilingual. The result was a VOT continuum ranging from -60 to 60 ms in 10 ms
559 increments. Finally, the stimuli were normalized for peak intensity.

560 **Procedure.** Upon completing the questionnaires, the learners participated in the
561 2AFC task. The initial experimental session took place on the second or third day of the
562 immersion program. Subsequent data collection occurred every Sunday for the remainder
563 of the program for a total of eight experimental sessions. The final experimental session
564 took place on a Wednesday, which was five days after the seventh session and three days
565 before the program ended. PsychoPy2 (Peirce, 2008) presented the stimuli described above
566 via a Macbook Pro. The program produced the audio stimuli at the same time that the
567 orthographic labels “ba” and “pa” appeared on the left-hand and right-hand sides of the
568 screen, respectively. The participants then determined whether they had heard “ba” or
569 “pa” by pressing the appropriate button on a DirectIN Rotary Controller. A red cross
570 appeared in the middle of the screen between trials, indicating a new trial was about to
571 begin. There was no set time limit for each trial; however, participants were instructed to
572 respond as quickly and as accurately as possible. The program presented one stimulus per
573 trial in ten randomized blocks (13 stimuli \times 10 blocks = 130 tokens) with the
574 inter-stimulus interval set at 500 ms. The participants finished the task in approximately 8
575 minutes. The task was administered in 8 separate sessions, once per week until the
576 conclusion of the program.

577 **Statistical analyses.** Data from the perception task were analyzed in R (R Core
578 Team, 2017) and can be separated into two principle analyses, one focused on the learners
579 perceptual behavior over time, and the other was concerned with how the learners
580 compared with bilingual controls at the offset of the immersion program. First, a series of
581 models were fit to examine the learners’ perceptual identification as exposure to Spanish
582 increased. Due to the categorical nature of the participants’ responses (i.e. “ba” or “pa”),
583 the data were analyzed using a GLMM with a binomially distributed error term and logit

584 linking function. The omnibus model was fit with *VOT* and *exposure time* as continuous
 585 fixed effects, and self-report assessments of input (% of input from the following sources:
 586 (1) native, (2) non-native, (3) non-native with higher proficiency, (4) non-native with lower
 587 proficiency, and (5) non-native with equal proficiency) and use (% time speaking Spanish,
 588 English) variables were included using forward selection. Causal priority was given to
 589 *exposure time*. All predictors were standardized such that their mean value was 0. The
 590 random effects structure included a scalar random effect for each subject with random
 591 slopes for *exposure time* and *VOT*.

592 Next, a second series of models was fit to examine how the learners' perceptual
 593 boundaries of the resynthesized continuum shifted as exposure increased. The random
 594 effects output from the aforementioned omnibus model was utilized to determine the 50%
 595 boundary crossover point (CO) for each participant at each testing session. The CO for the
 596 boundary between voiced and voiceless stops was calculated using the `cross_over` function
 597 of the package `lingStuff` (Casillas, 2018). This function calculated the perceptual
 598 boundary using the following formula:

$$CO = \frac{\beta_0}{\beta_{VOT}} \times -1 \quad (1)$$

599 where each by-subject intercept (β_0) is divided by the estimated by-subject slope for the
 600 effect *VOT* (β_{VOT}) and multiplied by -1. The CO point values were standardized and
 601 served as the dependent variable in subsequent analyses. In order to assess perceptual
 602 boundary shifts over time, the CO data were analyzed using a GLMM with *exposure time*
 603 as the dummy coded fixed effect. Day 0 was set as the reference level, thus the omnibus
 604 model provided parameter estimates of the change in CO values as time progressed with
 605 regard to the baseline perceptual boundary. The random effects structure included
 606 by-subject intercepts with a random slope for exposure time.

607 The final analysis compared the learners' perception of the resynthesized continuum

608 at the end of the program with that of the bilingual control group. This analysis included
 609 three models. The first model fit the identification response data (“ba”, “pa”) as a function
 610 of group (learners on day 47, bilingual controls) and VOT. The second model scrutinized
 611 the perceptual boundary (CO) data as a function of group. The final model examined
 612 contrast coefficient slopes (CCS) as a function of group. Contrast coefficient slopes in the
 613 logistic space were calculated for the corresponding sigmoidal curves in the probability
 614 space. The contrast coefficient slope gives a measure of “crispness” between phoneme
 615 boundaries (Morrison, 2007)⁴, and were derived for each participant using the parameter
 616 estimate for VOT from the random effects output of the omnibus model and the following
 617 equation:

$$CCS = \beta_{VOT} \times .25 \quad (2)$$

618 where the estimated by-subject slope for the effect VOT (β_{VOT}) was multiplied by .25. For
 619 the three aforementioned models *group* was deviation coded (-0.5, 0.5), thus the model
 620 parameter estimates provide an assessment of effect size. All mixed effects models were fit
 621 using the R package `lme4` (Bates et al., 2015). Main effects and higher order interactions
 622 were assessed using hierarchical partitioning of the variance via nested model comparisons.
 623 Visual inspection of Q-Q plots and plots of residuals against fitted values were used to
 624 check for normality of the residuals for linear models fit using Gaussian distributions.
 625 Marginal R^2 and conditional R^2 again provided an indication of goodness-of-fit for each
 626 model (Nakagawa & Schielzeth, 2013).

⁴ Speakers are believed to have “crisp” boundaries between native contrasts. When learning a new contrast, L2 speakers often have “fuzzier” boundaries, represented by shallower slopes. See Morrison (2007) for discussion on this topic.

627 **Results**

628 **Input, use, and perceptual categorization over time.** The GLMM yielded a
629 main effect of *VOT* ($\chi^2(1) = 24.88$; $p < 0.001$), *exposure time* ($\chi^2(1) = 4.92$; $p < 0.028$), as
630 well as a *VOT* x *exposure time* interaction ($\chi^2(1) = 4.78$; $p < 0.03$). The model containing
631 the higher order interaction was retained. There were no main effects nor interactions
632 related to the input and use predictors. The model output revealed that the log odds of
633 responding “voiceless” increased by 5.25 ± 0.4 standard errors as *VOT* increased (CI low =
634 4.46; CI high = 6.03; $z = 13.16$; $p < 0.001$). Overall, there was a change in the log odds of
635 “voiceless” responses of 0.19 ± 0.09 standard errors as a function of exposure time (CI low
636 = 0.01; CI high = 0.37; $z = 2.11$; $p < 0.05$). Moreover, the *VOT* x exposure time
637 interaction corresponded with a slope adjustment of 0.25 ± 0.11 standard errors (CI low =
638 0.03; CI high = 0.47; $z = 2.21$; $p < 0.04$), suggesting that the probability of responding
639 “voiceless” at the baseline *VOT* was higher as exposure to Spanish increased. Figure 4
640 plots predicted “voiceless” responses as modulated by *VOT* at each testing session. One
641 can observe a sigmoid function that appears to phase shift to the left over time.

642 ++ INSERT FIGURE 4 ABOUT HERE ++

643 The analysis of the boundary crossover point data revealed a main effect of session
644 ($\chi^2(7) = 18.3$; $p < 0.012$). Specifically, the crossover point was significantly different from
645 the week 0 baseline values after 14 days of exposure. At this point, the boundary had
646 decreased by -0.66 standardized units ± 0.29 standard errors (CI low = -1.23 ; CI high =
647 -0.09 ; $z = -2.29$; $p < 0.04$). The remaining sessions also showed significantly lower
648 boundary crossover points with the exception of day 35 (on the sixth experimental session).
649 The complete model output is shown in Table 6.

650 ++ INSERT TABLE 6 ABOUT HERE ++

651 **Comparison with bilinguals.** Figure 5 plots the results from the three models
 652 comparing the learners to the bilingual control group. Concretely, panel (a) displays the
 653 sigmoidal curves in the probability space along with the group CO points, and panel (b)
 654 plots the corresponding contrast coefficient slopes in the logistic space (b). The GLMM
 655 comparing the learners with the bilingual control group was best fit when including the
 656 random effects structure ($R^2_m = 0.88$, $R^2_c = 0.95$). The model yielded a main effect of
 657 *VOT* ($\chi^2(1) = 39.51$; $p < 0.001$). For both groups, a 10 ms increase in *VOT* was
 658 associated with a 0.2 ± 0.02 standard errors increase in the log odds of responding
 659 voiceless (CI low = 0.16; CI high = 0.23; $z = 11.35$; $p < 0.001$). There was no effect of
 660 *group* ($\chi^2(1) = 0.31$; $p > 0.05$), nor was there a *VOT* by *group* interaction ($\chi^2(1) = 3.04$; p
 661 > 0.05). In panel (a) of Figure 5 one can observe two nearly overlapping sigmoid functions,
 662 suggesting the two groups identified the resynthesized continuum in a similar manner.
 663 With regard to the crossover boundary data, there was no effect of *group* ($\beta = 0.24$, CI low
 664 = -0.71, CI high = 1.20, SE = 0.46, $t = 0.54$, $p = 0.60$). The vertical bars in panel (a) of
 665 Figure 5 show that the perceptual boundary for both groups nearly overlap. The model fit
 666 to the contrast coefficient slope data did not yield an effect of *group* ($\beta = 0.84$, CI low =
 667 -0.04, CI high = 1.72, SE = 0.42, $t = 2.02$, $p = 0.06$), though the effect approached
 668 significance. Panel (b) of Figure 5 plots the CCS slopes in the logistic space. The steepness
 669 of the lines suggests both groups had “crisp” boundaries, though it can be observed that
 670 the native control group has a slightly steeper slope.

671 ++ INSERT FIGURE 5 ABOUT HERE ++

672 Interim discussion

673 The perception experiment was concerned with uncovering how late, sequential
 674 language learners develop L2 perceptual strategies. Specifically, the perceptual
 675 identification task was designed with the purpose of analyzing how the /b/-/p/ stop

676 contrast was perceived as exposure to the target language increased, and how self-report
677 measures of input and use affect stop perception. The results of the task revealed that the
678 learners did indeed shift their perceptual boundaries with increased exposure to Spanish,
679 though there was no evidence that this shift was modulated by language input nor
680 language use. The analyses did suggest that after 14 days of exposure, learners began to
681 identify the resynthesized stimuli differently from how they had identified the same stimuli
682 two weeks prior. This finding supports the notion that the learners may have begun the
683 process of developing an L2-specific perceptual system. In this case, the contrast in
684 question, /b/-/p/, was one that already existed in their L1. Thus the learning that took
685 place involved the resetting of the perceptual boundary between the segments in this
686 contrast. The analyses found that by the end of the 7-week immersion
687 program—approximately 47 days—, the perceptual boundary of the learners was within
688 the range of the control group of simultaneous bilinguals. Moreover, the learners' linear
689 slope corresponding with the sigmoid functions of the boundary between the two segments
690 was also within the native bilingual range, though it did appear to be slightly less “crisp”.
691 Thus far the results of production and perception experiments support the notion that
692 phonetic category development may occur in a relatively short amount of time, at least
693 when L2 use is high and L1 use is minimized.

694 **Production/perception interface in L2 learning**

695 The final study examined the relationship between production and perception in
696 adult L2 learning. Specifically, the present analyses aimed to (1) determine if there was a
697 correlation between production gains and perceptual boundary shifts in late learners of
698 Spanish, and (2) determine if phonetic category development was perceptually driven, or if
699 production gains occurred before perception improved. To shed light on these issues, the
700 longitudinal production and perception data presented in the previous tasks were analyzed
701 together.

702 Method

703 Statistical analysis.

704 *Phoneme boundaries.*

705 Phoneme boundaries were calculated for each modality (production, perception). For
706 the production data, the boundary was the mean standardized value for all bilabials (/b/,
707 /p/) produced by a given participant for each session. That is, VOT was normalized as a
708 function of voicing and a mean was then calculated for each individual in each session.
709 Each participant provided one production boundary value per session. The perceptual
710 phoneme boundaries were the 50% crossover values analyzed in the perception task.

711 *Boundary trajectories: motivating GAMMs.*

712 The perception/production boundary data were analyzed using Generalized Additive
713 Mixed Models (GAMM, Sóskuthy, 2017; Winter & Wieling, 2016; Wood, 2006). GAMMs
714 represent an extension to the linear model framework that allow non-linear functions called
715 factor smooths to be applied to predictors. In this sense, the predictors can be classified
716 into two types: parametric terms (equivalent to fixed effects in hierarchical model
717 terminology) and smooth terms. Random smooths are conceptually similar to random
718 slopes and intercepts in the mixed-effects regression framework (Winter & Wieling, 2016).
719 Thus, they allow the by-subject trajectory shapes to vary as a function of a parametric
720 effect and are essential in avoiding anti-conservative models.

721 *Establishing production boundaries.*

722 As we have seen, the calculation of the perceptual boundary is straightforward. This
723 is not true for a production boundary given the fine-phonetic variability found in voicing
724 realizations. Previous work has utilized only voiceless stops as means to make comparisons
725 with perception. In the present analysis, production boundaries were calculated by
726 standardizing the VOT values from both stop categories and averaging them together. To

727 justify this calculation it is necessary to show that the rate of change over time was similar
728 for both segments. Recall from the production task that the learners' voiceless stops
729 decreased by approximately 10 ms by the end of the program and fell within the range of
730 the native bilinguals. For the voiced stops, VOT values decreased by approximately 40 ms
731 by the end of the program and did not fall within the native range. At first glance it
732 appears that the voiced stops showed a larger change over time. However, voiced stops
733 have a wider range of possible values with lead-VOT and short-lag VOT realizations, thus
734 mean change may mischaracterize overall production gains. For this reason VOT values for
735 each segment were standardized in order to put them on the same scale, and subsequently
736 analyzed using a GAMM to compare the change in VOT over time, that is, to analyze the
737 learning trajectory of each segment.

738 To this end, standardized VOT values were modeled as a function of the parametric
739 term *voicing* (voiced, voiceless) and a non-linear function of *exposure time*. *Voicing* was set
740 as an ordered variable with voiceless stops coded as 0. Cubic regression splines with 7 basis
741 knots were applied (1) as a reference smooth to *exposure time*, (2) as a difference smooth to
742 *exposure time* conditioned on *voicing*, and (3) as a random smooth for each participant
743 conditioned on *exposure time*. This specification uses the voiceless stop trajectory as the
744 baseline and compares it to the voiceless trajectory. Given that the VOT values were
745 standardized, we do not expect a *voicing* effect on the intercept (both levels of *voicing* have
746 an overall mean of 0). The model will, however, be informative regarding the shape of the
747 voiced and voiceless stop trajectories and how they differ from each other in terms of
748 curvature.

749 The model found no effect for the parametric voicing term, nor the corresponding
750 smoothing terms (see Table A1 for the model output), indicating that the trajectories for
751 voiced and voiceless stops did not differ from each other. Thus the standardized units were
752 averaged together to create the production boundary values which were subsequently
753 combined with the perception boundaries. In a separate GAMM this combination of

754 category boundary data was modeled as a function of modality (production, perception).

755 ***Production/perception trajectories.***

756 The model specification for the category boundary data was the same as the previous
757 voicing model and is described again here for the sake of completeness. VOT category
758 boundary values were modeled as a function of the parametric term *modality* (production,
759 perception) and a non-linear function of *exposure time*. *Modality* was set as an ordered
760 variable with perception coded as 0. Cubic regression splines with 7 basis knots were
761 applied (1) as a reference smooth to *exposure time*, (2) as a difference smooth to *exposure*
762 *time* conditioned on *modality*, and (3) as a random smooth for each participant conditioned
763 on *exposure time*. This specification uses the perception trajectory as the baseline and
764 compares it to the production trajectory. Given that the boundary values were
765 standardized, we again do not anticipate a *modality* effect on the intercept (both modalities
766 have an overall mean of 0), though the model will be informative regarding how the shape
767 of the production and perception learning trajectories differ from each other as a function
768 of *exposure time*.

769 All analyses were conducted in R using the `mgcv` package (Wood, 2004) for model
770 fitting and `itsadug` for visualization (van Rij et al., 2017). Autocorrelation was inspected
771 visually using autocorrelation function (ACF) plots of model residuals. Significance testing
772 for effects on *modality* were conducted using a combination of t-tests and approximate
773 F-tests on parametric and smooth terms, respectively, in conjunction with nested model
774 comparisons.

775 **Results**

776 Panel (a) of Figure 6 provides a scatterplot of the category boundary data. The
777 vertical axis plots the production boundaries and the horizontal axis plots the perception
778 boundaries. Lower VOT values take darker colors and higher VOT values are mapped to

779 lighter colors. Additionally, the two modalities are mapped to different parts of each
780 individual point. Production values are represented by the outer color of the points, while
781 perception values are represented by the inner color of the points. Exposure time is
782 mapped to each point based on geometric size so that the smallest points imply 0 days of
783 exposure and size increases in parallel with exposure time. One can observe that smaller,
784 lighter points are aggregated in the upper right quadrant of the plot and increase in size
785 and darkness as one moves towards the lower left quadrant. In other words, VOT values
786 for production (vertical axis) and perception (horizontal axis) appear to decrease as
787 exposure increases (point size). Furthermore, the relationship between category boundaries
788 is captured by the regression line plotted in black. One can observe that production and
789 perception boundaries decrease in tandem. The plot, while qualitative in nature, suggests
790 there is a relationship between the two modalities for these learners.

791 ++ INSERT FIGURE 6 ABOUT HERE ++

792 The modality GAMM explained 48.15% of the variance and 51.48% of the deviance.
793 Nested model comparisons suggested that the parametric and smooth terms on modality
794 significantly improved fit (DF = 3, Difference = 7.11, EDF = 9, $p < 0.003$). The perception
795 boundary varied as a function of exposure time (Reference smooth: EDF = 2.09, Ref. DF
796 = 2.54, $F = 8.41$, $p < 0.001$). The value higher than 1 for the effective degrees of freedom
797 (EDF) indicates that the trajectory was non-linear. Crucially, the production trajectory
798 differed from the perception trajectory (Difference smooth: EDF = 4.30, Ref. DF = 5.05,
799 $F = 4.68$, $p < 0.001$). The EDF value indicates that the difference between the trajectories
800 was also non-linear. Panel (b) of Figure 6 plots the modality trajectories (left side) and the
801 estimated differences between them over the time course (right side). The plots corroborate
802 the findings derived from the GAMM. Specifically, one can observe the non-linear time
803 course of both modalities. The perception boundaries shift earlier and at a consistent rate
804 before the slope flattens out around 25 days of exposure. The production boundaries, on

805 the other hand, follow a sigmoid-like trajectory, flattening out around around the same
806 time point as the production trajectory before rising and falling again around the final
807 testing session. The estimated difference plot indicates two exposure time windows of
808 significant differences: from days 2.65 - 13.83 and days 23.36 - 29.98.⁵ These are time
809 windows in which the difference in standardized VOT between voiced and voiceless stops is
810 significant at the 0.05 alpha level and are highlighted in the plot by the vertical,
811 discontinuous red lines. The full model summary is available in Table 7.

812 ++ INSERT TABLE 7 ABOUT HERE ++

813 **Interim discussion**

814 The final study presented an analysis of longitudinal data from the production and
815 perception tasks. These data were utilized to calculate category boundaries—for
816 production and perception—for each individual, for each experimental session. The
817 purpose of the present work was to determine if production and perception were related in
818 beginning L2 learners, and to determine if these learners showed perceptual learning before
819 increases in production accuracy (or *vice versa*). The analyses suggest a clear relationship
820 between production and perception in the learners' phonetic behavior. Specifically,
821 phonetic boundaries for both modalities decrease (i.e., shift towards native bilinguals'
822 Spanish boundaries) as exposure to and use of Spanish increases.

823 Of theoretical importance is that fact that the boundary shifts in speech perception
824 preceded those of speech production. The perceptual boundary shifts occurred early and
825 crossover points decreased at a steady rate over time before flattening out near the end of
826 the program. Importantly, by this point in their development, the L2 learners had

⁵ The estimated difference plot extrapolates the numeric time predictor over 100 values ranging from 1 to 42.

827 perceptual boundaries that fell within the range of the bilingual controls. The significant
828 difference between the modality trajectories supports the notion that production boundary
829 shifts occurred later in time. Thus, the data suggest that production and perception are
830 intimately related in the beginning stages of L2 phonetic category development and that
831 this development is perceptually driven.

832 **General discussion**

833 **Summary of findings**

834 The present work examined production and perception tasks, along with an analysis
835 that compared the two modalities. The production task showed that adult L2 learners of
836 Spanish reduced VOT in their production of bilabial stops at the end of the immersion
837 program. Moreover, production gains were modulated by self-reported English use.
838 Specifically, higher English use was associated with less target-like stop production.
839 Overall, the learners reduced aspiration for Spanish /p/ such that VOT was within the
840 range of the bilingual control groups' productions, though they were still slightly higher.
841 The learners also incorporated prevoicing into their production of Spanish /b/, though
842 their realization of the voiced segment was unstable. The results of the production task
843 show that with limited L1 use and high amounts of L2 input learners' pronunciation of
844 Spanish bilabial stops improved after limited exposure time (7 weeks).

845 The second task examined the learners' perception of bilabial stops in a 2AFC
846 identification task. The analyses showed that learners identified the same VOT continuum
847 differently over the course of the immersion program. Specifically, they were more likely to
848 categorize the resynthesized stimuli as being voiceless as exposure to and use of Spanish
849 increased. The increased voiceless responses led to a phonetic boundary shift to the left,
850 i.e., to a lower crossover point. The boundary shift was consistent with a more native-like
851 perception of the stimuli, that is, more towards the boundary of the bilingual group.

852 Perceptual categorization was not associated with input or use variables. The results
853 suggest that perceptual learning had taken place. Specifically, the evidence is consistent
854 with the notion that the learners were in the beginning stages of developing
855 language-specific speech perception.

856 The final analysis directly compared the longitudinal data from the production and
857 perception tasks. The purpose of the comparison was to determine if the two modalities
858 were related in the learner data, and also to uncover if perceptual learning had occurred
859 before the observed production gains (*or vice versa*). The two modalities showed a
860 decrease, or shift, in the phonetic category boundaries as exposure to and use of Spanish
861 increased. Moreover, the analyses showed that the perceptual boundary shifts occurred
862 prior to the production boundary shifts. Thus, the analysis provided evidence supporting
863 the hypotheses that (1) production and perception are indeed related in the beginning
864 stages of L2 learning and (2) that phonetic learning is perceptually driven during this early
865 period of development.

866 **L2 phonetic category development**

867 The results of the production task parallel those of Zampini (1998), who found that
868 native English intermediate/advanced late learners studying Spanish in a public university
869 reduced aspiration in voiceless stops, but still aspirated more than a control group. It may
870 be the case that the learners in question would have become more target-like over a more
871 extended period of time, or that the differences were due to the learners maintaining
872 separate categories for English and Spanish. In either case, the results from the present
873 production task extend Zampini's findings to a domestic immersion context where similar
874 changes occurred in a shorter time span (7 weeks). The present study differs from Zampini
875 (1998) in the voiced segment, /b/, which showed evidence of prevoicing by the end of the
876 program. Specifically, the learners showed a large relative change in VOT, but average
877 values still fell outside the bilingual range, suggesting there was still room for improvement.

878 However, analyzing the subset of prevoiced stops showed that all participants produced /b/
879 with lead-VOT some of the time, and, when they did, VOT values were closer to the
880 bilingual range. The advanced learners in Zampini (1998) did not incorporate prevoicing
881 into their productions of /b, d, g/.

882 The findings presented in the perception task also corroborate those of the perceptual
883 experiments conducted in Zampini (1998). Zampini (1998) showed that the late learners
884 improved their perception of Spanish bilabials over the course of a semester in a traditional
885 classroom setting. Specifically, Zampini (1998) found that perception of English and
886 Spanish stops became more target-like in conjunction with the learners' progress in
887 Spanish. The results from the perception task presented here suggested that the process of
888 perceptual learning may be sped up with increased L2 exposure and minimal L1 use.
889 Concretely, these results showed that 14 days of L2 input and high L2 usage were sufficient
890 for the adult learners of Spanish to shift their perceptual boundaries in a manner
891 consistent with more Spanish-like perception of bilabial stops.

892 The results also draw parallels with those of Stevens (2001). His pre/post test
893 analysis of Spanish stop production in a foreign immersion context found that adult
894 learners reduced aspiration in voiceless stops, and that this reduction was not observed
895 after one semester in traditional classroom learners. The present study extends these
896 findings to the domestic immersion context, at least for /p/. Stevens (2001) also found
897 that length of stay was positively correlated with production accuracy. Students that spent
898 an entire semester abroad produced more target-like stops. Stevens (2001) posited that this
899 may be due to increased use of Spanish. The present study found that, generally,
900 production became more target-like after around three weeks of exposure, and continued to
901 do so throughout the program. Due to the pre/post test nature of the experimental design
902 in Stevens (2001) it is impossible to say when exactly production gains began to occur. Be
903 that as it may, considered in conjunction with the present work, these findings support the
904 notion that maximizing L2 use and L2 input can accelerate the acquisition of L2

905 phonology. This affirmation is consistent with similar research that finds that the largest
906 phonetic gains occur during initial stages of learning (Williams, 1979; Munro et al., 2012).

907 Taken together, the results from the production and perception tasks suggest L2
908 phonetic learning can take place in an immersion context in a relatively short period of
909 time. The findings presented herein point to phonetic category development for the stop
910 voicing contrasts of Spanish, though this affirmation cannot be confirmed in the absence of
911 data from the learners' L1. Taking the most conservative view, these findings are at the
912 very least consistent with the early stages of phonetic category development.

913 The present findings also corroborate the notion that the ability to learn novel sounds
914 is maintained throughout the lifespan (Flege, 1995). With regard to the SLM, the fact that
915 the stops of Spanish are phonetically similar to those of English did not appear to impede
916 the learner group from detecting the fine-phonetic detail between them. Contrary to an
917 SLM account, the L2LP model predicts similar contrasts to be easier to learn. The case
918 presented here may be in agreement with the predictions of this model, though the notion
919 of "ease" is opaque, particularly in light of the fact that voiced bilabial production was
920 unstable. Another question to consider is whether or not adult learners are capable of
921 becoming native-like in the perception of their L2. In the context of the present study, the
922 term "native-like" is in reference to the behavior of the simultaneous (native) bilinguals.
923 The findings of the perception experiment revealed that the adult learners' boundaries
924 shifted to within the range of the native bilingual control group. It is important to note,
925 nonetheless, that both the SLM and the L2LP models conceive of the term "native-like" in
926 reference to the production/perception behavior of monolinguals. The SLM predicts
927 native-like attainment will be less likely with increased age due to cross-linguistic
928 interferences which result from the L2 categories residing in the same phonetic space as the
929 L1 categories. The L2LP model, for its part, maintains that native-like perception is indeed
930 possible due to the fact that the L1 sound system is copied and serves as the starting point
931 for L2 perception. With exposure, the boundaries between L2 categories are reset through

932 the GLA, which then operates independent of the L1 sound system. The results from the
933 perception experiment do not support any single model over the others because they do
934 not provide insight regarding the status of the L1 categories after the program.

935 In sum, it remains to be seen how the perceptual development displayed here would
936 continue over time. Does a newly acquired sound system disappear with diminished L2
937 use/input? What are the consequences for L1 sound categories? The L2LP predicts that
938 the L1 system will not be influenced by the development of the L2 system. L1-L2
939 interactions are expected to occur under an SLM account; however, if the L1 system is left
940 intact and operates independently of the newly acquired L2 system, as posited by the
941 L2LP, then the phonemic boundary for English stops should not change. Future research
942 could build upon the findings of the present work by including measures of English
943 production/perception in unambiguous English-sessions at the start and conclusion of the
944 immersion program. This would allow for a clearer understanding of cross-language
945 interference during the beginning stages of learning and over time. Furthermore, L1
946 category changes would provide evidence that the L1 and L2 perceptual systems might not
947 be separate, as suggested by the L2LP model.

948 **Production/perception interface in the acquisition of L2 phonology**

949 All major theories of speech perception posit a production/perception relationship.
950 The present findings corroborate a long line of research demonstrating a relationship
951 between speech production and speech perception in L1 and L2 acquisition (Smith, 1973;
952 Edwards, 1974; Ingram, 1977; Menyuk, 1977; Flege & Eefting, 1988; Williams, 1979; Flege
953 et al., 1999). Moreover, the results presented here contribute to previous research regarding
954 the production/perception interface by demonstrating the relationship between the two
955 modalities in beginning L2 learners with longitudinal data.

956 Another question addressed in production/perception analyses dealt with the causal

957 relationship between production and perception in L2 learning. Many researchers believe
958 that perception precedes production in language acquisition, though the literature has
959 shown that this assertion is not without controversy. The analyses presented here showed
960 that the perceptual boundary shifts occurred prior to the changes in production. This is
961 taken as evidence supporting the claim that changes in perception precede accurate
962 production in the beginning stages of L2 acquisition in adults. This finding is in line with
963 many studies suggesting that perception drives production in L1 and L2 acquisition
964 (Williams, 1979; Borden et al., 1983; Neufeld, 1988; Barry, 1989; Grasseger, 1991; Flege,
965 1993; Rochet, 1995; Llisterri, 1995; Flege et al., 1997; Leather, 1999). Furthermore,
966 perceptually driven L2 learning coincides with the contention that phonetic segments are
967 “[...] produced only as accurately as they are perceived” (Flege, 2003, p. 25). It appears,
968 at least in these data, that perceptual readjustments lead to, or at least precede, changes in
969 speech production.

970 Regarding evidence suggesting production precedes perception in L2 learning, it
971 remains possible that methodological concerns (see Escudero, 2006) can account for those
972 cases. Alternatively, these findings may simply be the result of a time confound that arises
973 based on the current state of the learners grammar at the moment in which the data are
974 collected. Another possibility is that subtle gains in one modality incite gains in the other
975 in a time-lagged relationship. This further underscores the importance of studying
976 phonological acquisition during the initial stages of learning. Future research could address
977 these matters by taking temporal resolution into account. For instance, one could
978 implement a longitudinal design that collects data in more experimental sessions over a
979 shorter period of time. It is also important to note that the data examined in the present
980 work dealt with the development of bilabial stops. Further support for a perception-first
981 developmental path would be provided by demonstrating that perceptual boundary shifts
982 precede production accuracy in other segments as well. The stop segments under analysis
983 likely represent an “easy” boundary adjustment, given that the phonological contrast is

984 already part of the participants' native phonology. Future investigations ought to extend
985 the aforementioned findings to other L2 segments, such as laterals, nasals, and vowels. This
986 would further our understanding of how the production/perception of different phonetic
987 segments develop over time. Of particular interest is the longitudinal development of
988 spirantization of voiced stops and the acquisition of the rhotic trill, /r/. The former is
989 contrastive in American English (as opposed to allophonic in Spanish), and the latter is not
990 part of the American English phonemic inventory. Thus, both segments provide clear
991 instances where novel category formation would be necessary for production/perception,
992 and neither segment is likely to be explained via boundary resetting. In essence, the
993 acquisition of these segments could provide an interesting point of comparison to stop
994 voicing, and the amount of exposure necessary for category formation to occur could shed
995 light on the relative difficulty of boundary resetting versus learning a new allophonic
996 distribution versus learning a new sound altogether.

997 Finally, the production/perception analysis showed that the trajectories for each
998 modality differed during two points in the time course (see Figure 6b, right panel), once
999 during the initial point of exposure, and later again between days 23 and 30. This second
1000 window merits further consideration. Specifically, at this time, the perceptual boundary
1001 trajectory was beginning to flatten out as the production boundary continued to decrease.
1002 Interestingly, the trajectories approached one another again. A priori one would not expect
1003 the trajectories to stop being significantly different. One would expect the production
1004 boundary to continue to decline as the learners incorporated more prevoicing into their
1005 production. This pattern is not observed. Conversely, the production trajectory begins to
1006 rise around day 30, reaches a peak around day 35, and then continues to fall. This
1007 variability may be explained by the participants extra-curricular activities during this time
1008 frame of the immersion program. Specifically, the Language Schools were celebrating their
1009 centennial anniversary and the language pledge was suspended during one evening for a
1010 banquet and party. In other words, the participants spent the night prior to the

1011 penultimate testing session (day 35) hearing and possibly speaking English. A posteriori
1012 the participants were asked if they had spoken English during the centennial celebration.
1013 Most reported that they had and that it had been a cathartic experience. This suggests
1014 that the variability in the phonetic behavior of the learners may be explained by L1
1015 language use. This hypothesis is supported by the finding that self-reported use of English
1016 was correlated with overall change in VOT in the production data. Importantly, English
1017 use did not modulate perceptual categorization. This poses the possibility that speech
1018 production is more affected by native language activation than speech perception. Taken
1019 together, the results suggest that phonetic category development during the early stages of
1020 learning appears to be particularly fragile and susceptible to cross-linguistic interference.
1021 That said, the trajectory analysis suggests both voiced and voiceless stops would have
1022 continued to improve, though it is likely that production gains would eventually slow,
1023 possibly to maintain distinct categories for each language. A likely scenario is that Spanish
1024 stop production would become more stable with increased exposure to and use of the L2.
1025 Future research should directly compare the domestic immersion learning context with
1026 study abroad and the traditional classroom with longitudinal designs. This would
1027 ultimately help tease apart the effects of L1/L2 use and the relative importance of different
1028 types of target language input, which, in the present work, were operationalized to include
1029 various factors, such as L2 feedback, motivation, and attention.

1030

Conclusion

1031 The present investigation analyzed early second language learning in adults. The
1032 studies undertaken for this work regarding the ongoing development of the fine-phonetic
1033 detail of Spanish stop voicing add to our understanding of the acquisition of L2 phonology
1034 in adult learners. The longitudinal data suggest that L2 phonetic category formation can
1035 occur abruptly at an early stage of development, and is perceptually driven. Moreover,
1036 early, developing L2 sound representations are fragile, and especially susceptible to

¹⁰³⁷ cross-linguistic interference during the initial stages of learning.

References

1038

- 1039 Barry, W. J. (1989). Perception and production of English vowels by German learners:
1040 instrumental-phonetic support in language teaching. *Phonetica*, 46(4), 155–168.
- 1041 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models
1042 using lme4. *Journal of Statistical Software*, 67(1), 1–48. doi: 10.18637/jss.v067.i01
- 1043 Birdsong, D., Gertken, L. M., & Amengual, M. (2012). *Bilingual language profile: an*
1044 *easy-to-use instrument to assess bilingualism*. Retrieved from
1045 <https://sites.la.utexas.edu/bilingual/>
- 1046 Boersma, P., & Weenink, D. (2018). *Praat: doing phonetics by computer [computer*
1047 *program]*. Retrieved from <http://www.praat.org/>
- 1048 Borden, G., Gerber, A., & Milsark, G. (1983). Production and perception of the /r/-/l/
1049 contrast in Korean adults learning English. *Language Learning*, 33(4), 499–526.
- 1050 Bradlow, A. R., Akahane-Yamada, R., Pisoni, D. B., & Tohkura, Y. (1999). Training
1051 Japanese listeners to identify English /r/ and /l/: Long-term retention of learning in
1052 perception and production. *Perception & psychophysics*, 61(5), 977–985.
- 1053 Bradlow, A. R., Pisoni, D. B., Akahane-Yamada, R., & Tohkura, Y. (1997). Training
1054 Japanese listeners to identify English /r/ and /l/: IV. Some effects of perceptual learning
1055 on speech production. *Journal of the Acoustical Society of America*, 101(4), 2299–2310.
- 1056 Caramazza, A., Yeni-Komshian, G., Zurif, E. B., & Carbone, E. (1973). The acquisition of
1057 a new phonological contrast: The case of stop consonants in French-English bilinguals.
1058 *The Journal of the Acoustical Society of America*, 54(2), 421–428.
- 1059 Casillas, J. V. (2018). lingstuff: Tools and gems for linguistics related research [Computer
1060 software manual]. Retrieved from <https://github.com/jvcasill/lingStuff> (R
1061 package version 0.1.1)

- 1062 Davidson, L. (2018). Phonation and laryngeal specification in American English voiceless
1063 obstruents. *Journal of the International Phonetic Association*, 48(3), 331–356.
- 1064 Díaz-Campos, M. (2006). The effect of style in second language phonology: An analysis of
1065 segmental acquisition in study abroad and regular-classroom students. In *Selected*
1066 *proceedings of the 7th conference on the acquisition of spanish and portuguese as first and*
1067 *second languages* (pp. 26–39).
- 1068 Edwards, M. L. (1974). Perception and production in child phonology: The testing of four
1069 hypotheses. *Journal of Child language*, 1(2), 205–219.
- 1070 Elliott, A. R. (1997). On the teaching and acquisition of pronunciation within a
1071 communicative approach. *Hispania*, 95–108.
- 1072 Escudero, P. (2005). *Linguistic perception and second language acquisition* (Unpublished
1073 doctoral dissertation). Utrecht University, Utrecht, Holland.
- 1074 Escudero, P. (2006). Second-language phonology: the role of perception. In
1075 M. C. Pennington (Ed.), *Phonology in context* (pp. 109–134).
- 1076 Escudero, P. (2009). Linguistic perception of ‘similar’ L2 sounds. In P. Boersma &
1077 S. Hamann (Eds.), *Phonology in perception* (pp. 151–190). Berlin: Mouton de Gruyter.
- 1078 Escudero, P., & Boersma, P. (2004). Bridging the gap between L2 speech perception
1079 research and phonological theory. *Studies in Second Language Acquisition*, 26(4),
1080 551–585.
- 1081 Flege, J. E. (1981). The phonological basis of foreign accent: A hypothesis. *Tesol*
1082 *Quarterly*, 15(4), 443–455.
- 1083 Flege, J. E. (1987a). A critical period for learning to pronounce foreign languages? *Applied*
1084 *Linguistics*, 8, 162–177.

- 1085 Flege, J. E. (1987b). The production of “new” and “similar” phones in a foreign language:
1086 Evidence for the effect of equivalence classification. *Journal of phonetics*, 15(1), 47–65.
- 1087 Flege, J. E. (1993). Production and perception of a novel, second-language phonetic
1088 contrast. *The Journal of the Acoustical Society of America*, 93(3), 1589–1608.
- 1089 Flege, J. E. (1995). Second language speech learning: Theory, findings, and problems. In
1090 W. Strange (Ed.), *Speech perception and linguistic experience issues in cross-language*
1091 *research* (pp. 229–273). Timonium, MD: York Press.
- 1092 Flege, J. E. (2003). Assessing constraints on second-language segmental production and
1093 perception. In *Phonetics and phonology in language comprehension and production*
1094 *differences and similarities* (pp. 319–355). Berlin: De Gruyter Mouton.
- 1095 Flege, J. E. (2012). The role of input in second language (L2) speech learning. Lodz,
1096 Poland.
- 1097 Flege, J. E., Bohn, O.-S., & Jang, S. (1997). Effects of experience on non-native speakers’
1098 production and perception of English vowels. *Journal of Phonetics*, 25(4), 437–470.
- 1099 Flege, J. E., & Eefting, W. (1987). Production and perception of English stops by native
1100 Spanish speakers. *Journal of Phonetics*, 15, 67–83.
- 1101 Flege, J. E., & Eefting, W. (1988). Imitation of a VOT Continuum by Native Speakers of
1102 English and Spanish - Evidence for Phonetic Category Formation. *Journal of the*
1103 *Acoustical Society of America*, 83(2), 729–740.
- 1104 Flege, J. E., & Liu, S. (2001). The effect of experience on adults’ acquisition of a second
1105 language. *Studies in Second Language Acquisition*, 23(4), 527–552.
- 1106 Flege, J. E., MacKay, I. R. A., & Meador, D. (1999). Native Italian speakers’ perception
1107 and production of English vowels. *The Journal of the Acoustical Society of America*,
1108 106(5), 2973–2987.

- 1109 Flege, J. E., Munro, M. J., & MacKay, I. R. A. (1995). Factors affecting strength of
1110 perceived foreign accent in a second language. *The Journal of the Acoustical Society of*
1111 *America*, 97(5 Pt 1), 3125–3134.
- 1112 Fowler, C. A., Sramko, V., Ostry, D. J., Rowland, S. A., & Hallé, P. (2008). Cross
1113 language phonetic influences on the speech of French–English bilinguals. *Journal of*
1114 *Phonetics*, 36(4), 649–663.
- 1115 González-Bueno, M. (1997). Voice onset time in the perception of foreign accent by native
1116 listeners of Spanish. *International Review of Applied Linguistics in Language Teaching*,
1117 35(4), 251–262.
- 1118 González López, V., & Counselman, D. (2013). L2 acquisition and category formation of
1119 Spanish voiceless stops by monolingual English novice learners. In J. Cabrelli Amaro,
1120 G. Lord, A. de Prada Pérez, & J. E. Aaron (Eds.), (pp. 118–127). Somerville, MA:
1121 Cascadilla Proceedings Project.
- 1122 Goto, H. (1971). Auditory perception by normal Japanese adults of the sounds “L” and
1123 “R”. *Neuropsychologia*, 9(3), 317–323.
- 1124 Grasseger, H. (1991). Perception and production of Italian plosives by Italian learners. In
1125 *Actes du xii congrès international des sciences phonétiques* (Vol. 5, pp. 290–293).
- 1126 Ingram, D. (1977). *Phonological disability in children* (Vol. 2). New York: Elsevier.
- 1127 Kissling, E. M. (2014). Phonetics instruction improves learners’ perception of l2 sounds.
1128 *Language Teaching Research*, 1362168814541735.
- 1129 Leather, J. (1997). Interrelation of perceptual and productive learning in the initial
1130 acquisition of second-language tone. In A. James & J. Leather (Eds.), *Second-language*
1131 *speech: structure and process* (pp. 75–102). Berlin, New York: DE GRUYTER
1132 MOUTON.

- 1133 Leather, J. (1999). Second-language research: An introduction. *Language Learning*,
1134 49(s1), 1–56.
- 1135 Lisker, L., & Abramson, A. S. (1964). A Cross-language Study of Voicing in Initial Stops:
1136 Acoustical Measurements. *Word*, 20.3, 384–422.
- 1137 Llisterri, J. (1995). Relationships between speech production and speech perception in a
1138 second language. *Proceedings of the 13th International Congress of Phonetic Sciences*, 4,
1139 92–99.
- 1140 López, V. G. (2012). Spanish and english word-initial voiceless stop production in
1141 code-switched vs. monolingual structures. *Second Language Research*, 28(2), 243–263.
- 1142 Mack, M. (1989). Consonant and vowel perception and production: Early English-French
1143 bilinguals and English monolinguals. *Perception & Psychophysics*, 46(2), 187–200.
- 1144 Mathews, J. (1997). The influence of pronunciation training on the perception of
1145 second-language contrasts. In J. Leather & A. James (Eds.), *New sounds 97: Proceedings*
1146 *of the third international symposium on the acquisition of second language speech*
1147 (p. 223–229). Klagenfurt: University of Klagenfurt.
- 1148 Menyuk, P. (1977). *Language and maturation*. Cambridge, MA: MIT Press.
- 1149 Morrison, G. S. (2007). Logistic regression modelling for first- and second-language
1150 perception data. In M. J. Sole, P. Prieto, & J. Mascaró (Eds.), *Segmental and prosodic*
1151 *issues in romance phonology* (pp. 219–236). John Benjamins Publishing Company.
- 1152 Munro, M. J., Derwing, T. M., & Saito, K. (2012). English l2 vowel acquisition over seven
1153 years. *PRONUNCIATION AND ASSESSMENT*, 112.
- 1154 Nagle, C. L. (2017). A Longitudinal Study of Voice Onset Time Development in L2
1155 Spanish Stops. , 54(5), 13-23.

- 1156 Nagle, C. L., Moorman, C., Sanz, C., et al. (2016). Disentangling research on study abroad
1157 and pronunciation: Methodological and programmatic considerations. In *Handbook of*
1158 *research on study abroad programs and outbound mobility* (pp. 673–695). IGI Global.
- 1159 Nakagawa, S., & Schielzeth, H. (2013). A general and simple method for obtaining R^2 from
1160 generalized linear mixed-effects models. *Methods In Ecology And Evolution*, *4*, 133–142.
- 1161 Neufeld, G. G. (1988). Phonological asymmetry in second-language learning and
1162 performance. *Language Learning*, *38*(4), 531–559.
- 1163 Oyama, S. (1976). A sensitive period for the acquisition of a nonnative phonological
1164 system. *Journal of Psycholinguistic Research*, *5*, 261–283.
- 1165 Pallier, C., Bosch, L., & Sebastián-Gallés, N. (1997). A limit on behavioral plasticity in
1166 speech perception. *Cognition*, *64*(3), B9–17.
- 1167 Peirce, J. W. (2008). Generating stimuli for neuroscience using PsychoPy. *Frontiers in*
1168 *Neuroinformatics*, *2*(10), 1–8.
- 1169 R Core Team. (2017). R: A language and environment for statistical computing [Computer
1170 software manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- 1171 Rochet, B. L. (1995). Perception and production of second-language speech sounds by
1172 adults. *Speech perception and linguistic experience: Issues in cross-language research*,
1173 379–410.
- 1174 Sakai, M. (2018). Moving towards a bilingual baseline in second language phonetic
1175 research. *Journal of Second Language Pronunciation*, *4*(1), 11–45.
- 1176 Saville-Troike, M. (2005). *Introducing second language acquisition*. Cambridge University
1177 Press.
- 1178 Sebastián-Gallés, N. (2006). Native-language sensitivities: evolution in the first year of life.
1179 *Trends in cognitive sciences*, *10*(6), 239–241.

- 1180 Sebastián-Gallés, N., Echeverría, S., & Bosch, L. (2005). The influence of initial exposure
1181 on lexical representation: comparing early and simultaneous bilinguals. *Journal of*
1182 *Memory and Language*, 52(2), 240–255.
- 1183 Sebastián-Gallés, N., & Soto-Faraco, S. (1999). Online processing of native and non-native
1184 phonemic contrasts in early bilinguals. *Cognition*, 72(2), 111–123.
- 1185 Sheldon, A., & Strange, W. (1982). The acquisition of /r/ and /l/ by Japanese learners of
1186 English: Evidence that speech production can precede speech perception. *Applied*
1187 *Psycholinguistics*, 3(03), 243–261.
- 1188 Smith, N. V. (1973). *The acquisition of phonology: A case study*. Cambridge University
1189 Press.
- 1190 Sóskuthy, M. (2017). Generalised additive mixed models for dynamic analysis in
1191 linguistics: a practical introduction. *arXiv preprint arXiv:1703.05339*.
- 1192 Stevens, J. (2001). Study abroad learners' acquisition of the Spanish voiceless stops.
1193 *MIFLC Review*, 10, 137–151.
- 1194 Strange, W. (1995). Cross-language studies of speech perception: A historical review. In
1195 W. Strange (Ed.), *Speech perception and linguistic experience: Issues in cross language*
1196 *research* (pp. 3–45). Baltimore, MD: York Press.
- 1197 Sundara, M., & Polka, L. (2008). Discrimination of coronal stops by bilingual adults: the
1198 timing and nature of language interaction. *Cognition*, 106(1), 234–258.
- 1199 Sundara, M., Polka, L., & Baum, S. (2006). Production of coronal stops by simultaneous
1200 bilingual adults. *Bilingualism: Language and Cognition*, 9(01), 97–114.
- 1201 Sundara, M., Polka, L., & Genesee, F. (2006). Language-experience facilitates
1202 discrimination of /d/-/ð/ in monolingual and bilingual acquisition of english. *Cognition*,
1203 100(2), 369–388.

- 1204 van Rij, J., Wieling, M., Baayen, R. H., & van Rijn, H. (2017). *itsadug: Interpreting time*
1205 *series and autocorrelated data using gamms*. (R package version 2.3)
- 1206 Wang, X. (2002). Training Mandarin and Cantonese Speakers to Identify English Vowel
1207 Contrasts: Long-term Retention and Effects on Production (Unpublished doctoral
1208 dissertation). British Columbia.
- 1209 Williams, L. (1979). The modification of speech perception and production in
1210 second-language learning. *Perception & Psychophysics*, 26(2), 95–104.
- 1211 Winter, B., & Wieling, M. (2016). How to analyze linguistic change using mixed models,
1212 growth curve analysis and generalized additive modeling. *Journal of Language*
1213 *Evolution*(1), 7-18.
- 1214 Wood, S. (2004). Stable and efficient multiple smoothing parameter estimation for
1215 generalized additive models. *Journal of the American Statistical Association*, 99(467),
1216 673-686.
- 1217 Wood, S. (2006). *Generalized additive models: an introduction with r* (2nd ed.). Boca
1218 Raton: CRC Press.
- 1219 Zampini, M. L. (1998). The Relationship between the Production and Perception of L2
1220 Spanish Stops. *Texas Papers in Foreign Language Education*, 3(3), 85–100.
- 1221 Zampini, M. L., & Green, K. P. (2001). The voicing contrast in English and Spanish: The
1222 relationship between perception and production. In *One mind, two languages bilingual*
1223 *language processing* (pp. 23–48). Malden, Mass; Oxford: Blackwell.

ID	Age	Use		Input type					Estimated ability		
		Sp.	En.	NI	NNI	NNI+	NNI-	NNI=	Speaking	Listening	Overall
101	21	98.3	3.3	45.0	55.0	43.3	56.7	50.0	51	74	60
102	26	93.3	6.7	53.3	70.0	68.3	18.3	65.0	69	69	70
103	28	100.0	16.7	51.7	48.3	78.3	10.0	11.7	44	64	50
104	18	95.0	0.0	48.3	53.3	68.3	30.0	45.0	57	71	59
105	20	70.0	31.7	35.0	70.0	61.7	30.0	66.7	44	57	49
106	34	86.7	13.3	33.3	61.7	58.3	15.0	21.7	50	64	60
107	22	80.0	15.0	33.3	63.3	63.3	13.3	43.3	34	50	40
108	29	93.3	6.7	26.7	71.7	98.3	1.7	0.0	14	19	14
109	19	75.0	23.3	20.0	53.3	43.3	28.3	40.0	37	64	41
110	20	91.7	18.3	58.3	91.7	71.7	46.7	83.3	49	64	53
<i>Avg.</i>	23.7	88.3	13.5	40.5	63.8	65.5	25.0	42.7	45	59	49

Table 1

Averaged self-report values (from weekly assessment questionnaire) of Spanish/English use, estimated native input (NI), non-native input (NNI), non-native input from speakers with a higher level (NNI+), non-native input from speakers with a lower level (NNI-), non-native input from speakers with the same level (NNI=), and estimated speaking/listening/overall ability (in Spanish). All measures represent percentages.

Session	Questionnaires		Perception		Production	
	Demographic	Assessment	2AFC (a)	2AFC (b)	Repetition	Picture Naming
Week 0	✓	✓	✓		✓	
Week 1		✓	✓		✓	
Week 2		✓	✓		✓	
Week 3		✓	✓		✓	
Week 4		✓	✓		✓	
Week 5		✓	✓		✓	
Week 6		✓	✓		✓	
Week 7		✓	✓	✓	✓	✓

Table 2

Timetable of experimental sessions for the learner group. 2AFC (b) and picture naming are not reported.

Term	β	SE	CI low	CI high	Statistic	Pr(> t)	
Intercept	32.06	5.43	21.42	42.71	5.90	0.001	*
Day 7	2.59	3.41	-4.09	9.27	0.76	0.460	
Day 14	-3.50	3.41	-10.18	3.18	-1.03	0.310	
Day 21	-8.86	3.41	-15.55	-2.18	-2.60	0.013	*
Day 28	-12.21	3.41	-18.89	-5.53	-3.58	0.002	*
Day 35	-9.60	3.41	-16.28	-2.92	-2.82	0.008	*
Day 42	-10.75	3.41	-17.43	-4.06	-3.15	0.004	*

Table 3

Model output for VOT of voiceless stops as a function of exposure time. The intercept represents VOT on day 0. Parameter estimates from subsequent sessions represent the change in VOT with regard to the intercept.

ID	Day 1	Day 7	Day 14	Day 21	Day 28	Day 35	Day 42	Avg.
101	0.00	26.67	60.00	100.00	86.67	100.00	100.00	67.62
102	0.00	0.00	0.00	0.00	6.67	26.67	40.00	10.48
103	0.00	0.00	0.00	0.00	6.67	20.00	33.33	8.57
104	0.00	0.00	0.00	0.00	6.67	13.33	33.33	7.62
105	0.00	0.00	0.00	6.67	0.00	20.00	33.33	8.57
106	6.67	0.00	13.33	46.67	93.33	66.67	53.33	40.00
107	0.00	0.00	0.00	0.00	13.30	26.67	33.33	10.47
108	5.88	6.67	0.00	100.00	87.50	93.75	100.00	56.26
109	0.00	0.00	0.00	26.67	20.00	6.67	86.67	20.00
110	0.00	0.00	0.00	20.00	33.33	33.33	26.67	16.19
Avg.	1.25	3.33	7.33	30.00	35.41	40.71	54.00	24.58

Table 4

Proportion of prevoiced /b/ realizations as a function of exposure time.

Term	β	SE	CI low	CI high	Statistic	Pr(> t)
Intercept	13.54	9.14	-4.37	31.45	1.48	0.149
Day 7	-2.44	9.38	-20.81	15.94	-0.26	0.797
Day 14	-5.24	9.37	-23.62	13.13	-0.56	0.579
Day 21	-31.51	9.38	-49.89	-13.13	-3.36	0.002 *
Day 28	-34.30	9.38	-52.69	-15.92	-3.66	0.002 *
Day 35	-30.11	9.37	-48.48	-11.73	-3.21	0.003 *
Day 42	-38.49	9.37	-56.86	-20.12	-4.11	0.001 *

Table 5

Model output for VOT of voiced stops as a function of exposure time. The intercept represents VOT on day 0. Parameter estimates from subsequent sessions represent the change in VOT with regard to the intercept.

Term	β	CI low	CI high	SE	DF	Statistic	Pr(> t)
Intercept	0.57	-0.01	1.15	0.30	27.81	1.94	0.064
Day 7	-0.14	-0.71	0.43	0.29	70.00	-0.49	0.627
Day 14	-0.66	-1.24	-0.10	0.29	70.00	-2.29	0.026 *
Day 21	-0.64	-1.21	-0.07	0.29	70.00	-2.21	0.040 *
Day 28	-0.80	-1.37	-0.23	0.29	70.00	-2.74	0.009 *
Day 35	-0.47	-1.04	0.10	0.29	70.00	-1.62	0.110
Day 42	-0.80	-1.37	-0.23	0.29	70.00	-2.74	0.009 *
Day 47	-1.07	-1.64	-0.50	0.29	70.00	-3.68	0.001 *

Table 6

Model output for perceptual boundary crossover point as a function of time. The intercept represents perceptual boundaries on day 0. Parameter estimates from subsequent sessions quantify boundary shifts with regard to the intercept.

Model output

Parametric coefficients:

	Estimate	Std. Error	t value	p-value
Intercept	0	0.08	0	> 0.05
Intercept Δ Production	0	0.098	0	> 0.05

Approximate significance of Smooth terms:

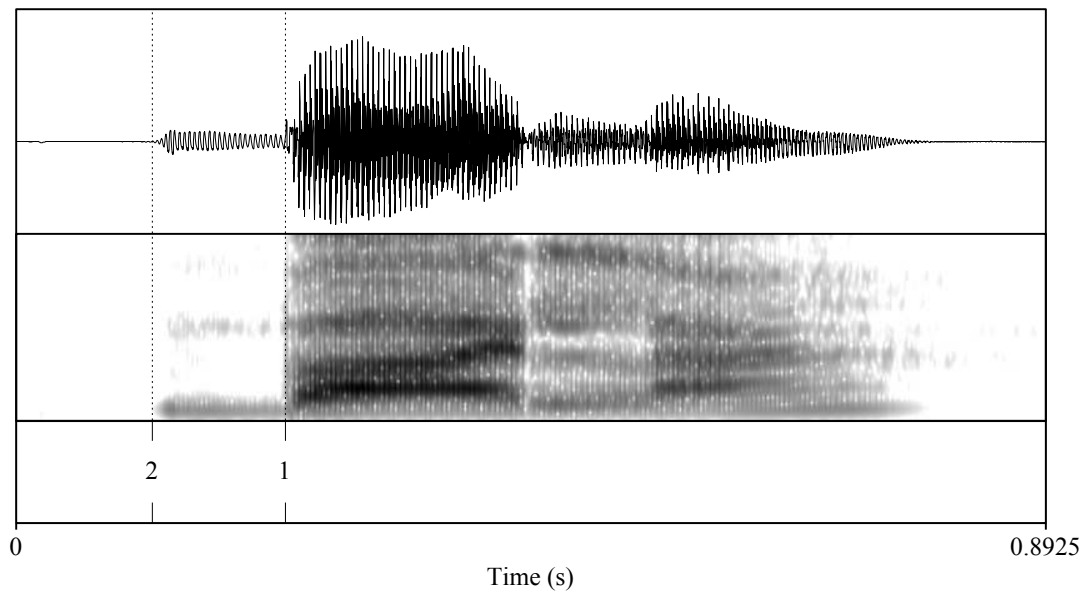
	EDF	Ref. DF	F	p-value
Reference smooth: exposure time	2.087	2.539	8.408	< 0.001
Difference smooth: production	4.298	5.048	4.684	< 0.001
Random smooth: exposure time x participant	6.016	68	0.14	> 0.05

$R^2 = 0.48$; Deviance explained: 51.48%

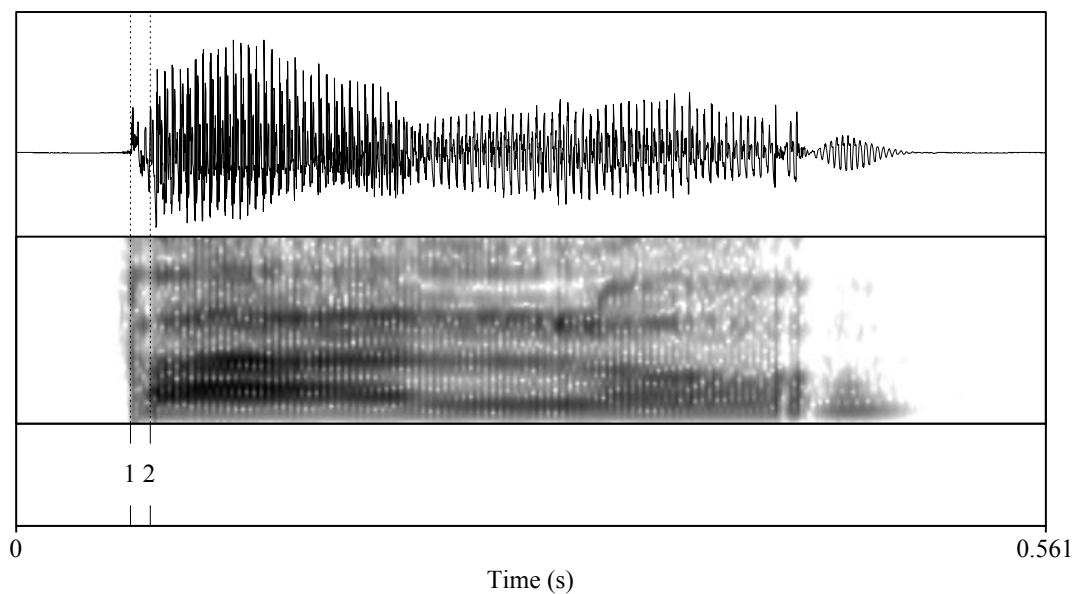
n = 140

Table 7

Summary of production/perception GAMM model output. The full model fit VOT category boundaries as a function of modality.



(a) Spanish *bala* [ˈba.la] (Eng. ‘bullet’).



(b) Spanish *palo* [ˈpa.lo] (Eng. ‘stick’).

Figure 1. Segmenting procedures for Spanish stops. Panels (a) and (b) illustrate pre-voiced and short lag VOT, respectively. The release of the stop is labeled (1) and the onset of modal voicing is labeled (2). VOT was calculated as the time (in ms) from (1) to (2).

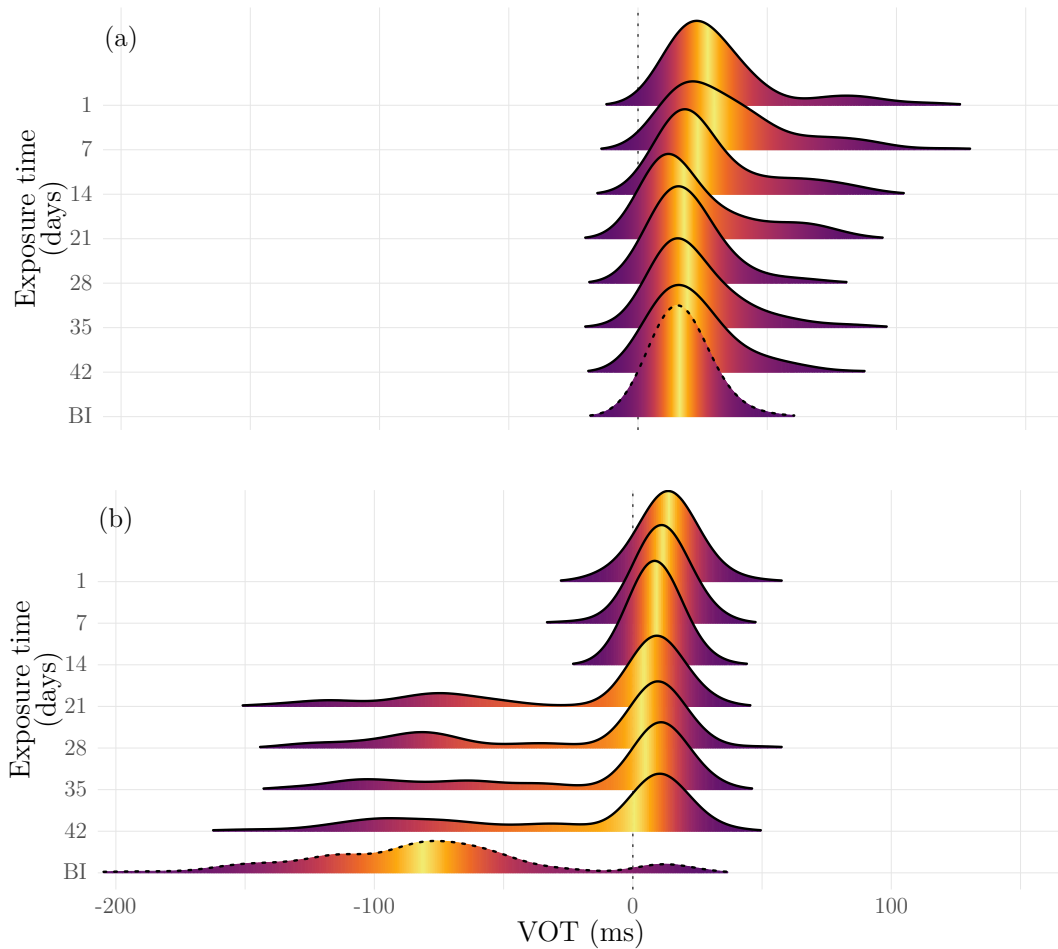


Figure 2. Ridgeline density plots of voiceless (panel ‘a’) and voiced (panel ‘b’) stop VOT (ms) as a function of exposure time. The final ridgeline of each panel, outlined with a discontinuous line, represents the VOT data from the bilingual control group. Lighter colors indicate proximity to the geometric center of gravity of the distribution for each experimental session.

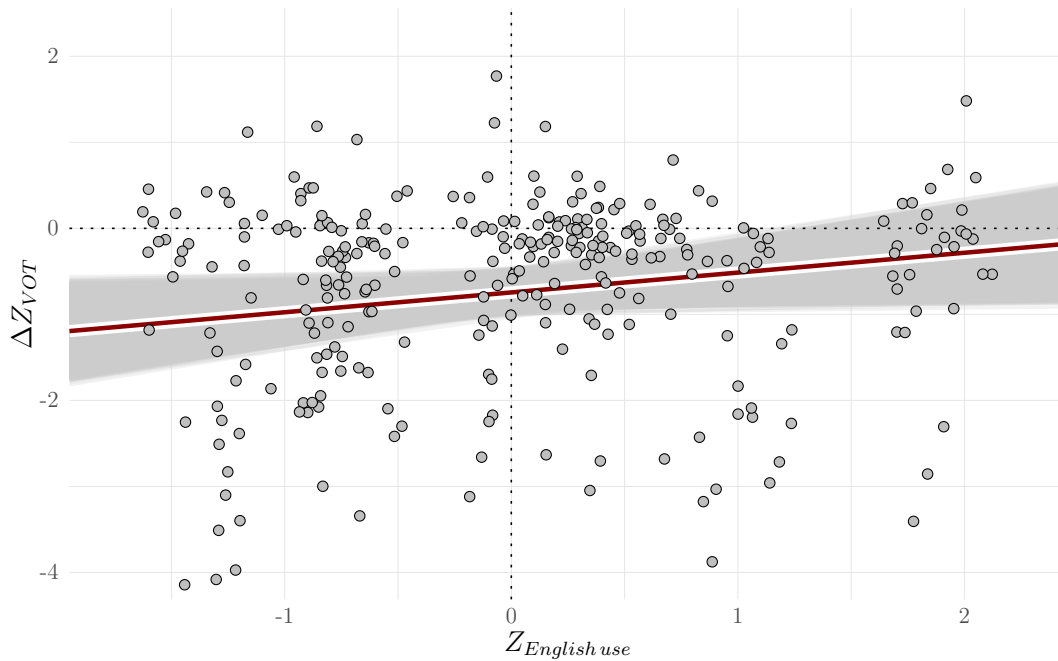


Figure 3. Scatterplot and line of best fit $\pm 95\%$ CI for ΔZ_{VOT} as a function of Z -English use. ΔZ_{VOT} represents the change in standardized VOT for bilabial stops at the end of the immersion program and Z -English use represents the average of self-reported time spent speaking English in standardized units. The x-axis is jittered horizontally to show overlapping points.

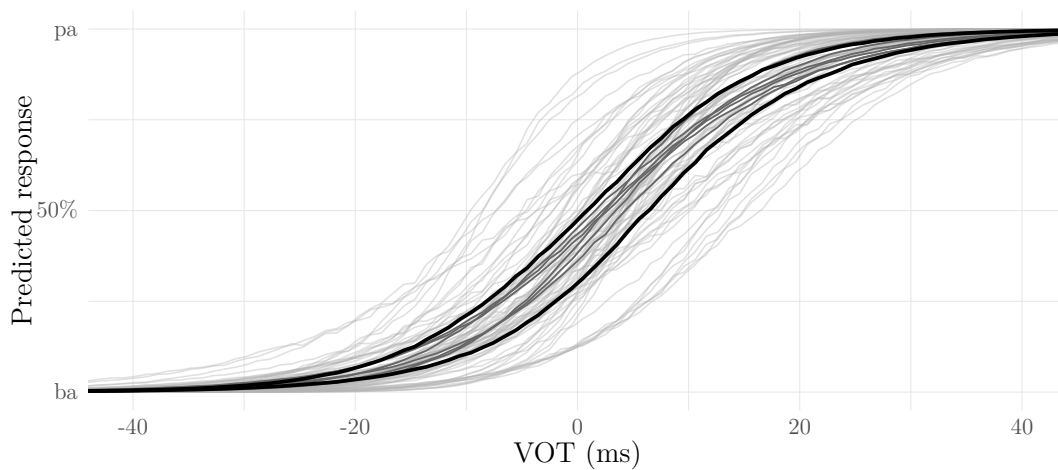


Figure 4. Voiceless responses as a function of VOT. Each line represents an experimental session. The black lines represent the first and last sessions of the program.

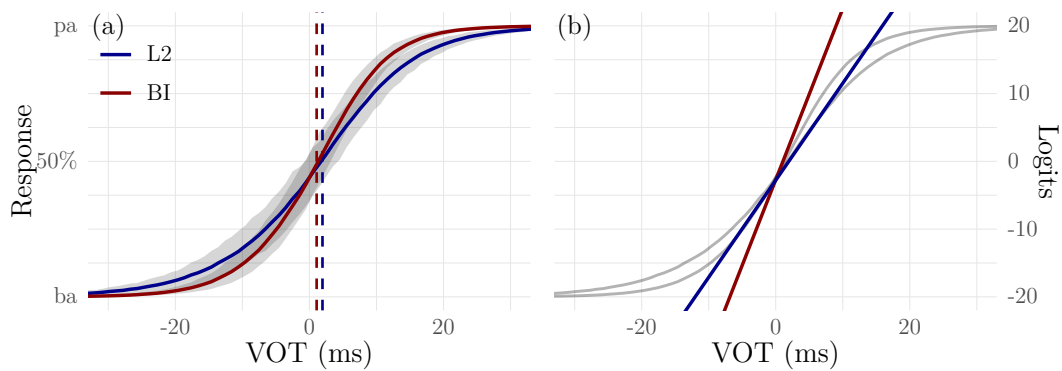


Figure 5. Sigmoidal curves in the probability space (a) and the corresponding contrast coefficient slopes in the logistic space (b) for the learner and bilingual groups. In panel (a), the vertical bars indicate the boundary crossover point for each group. In both panels, the red lines represent the bilingual controls, and the blue lines indicate the learner group.

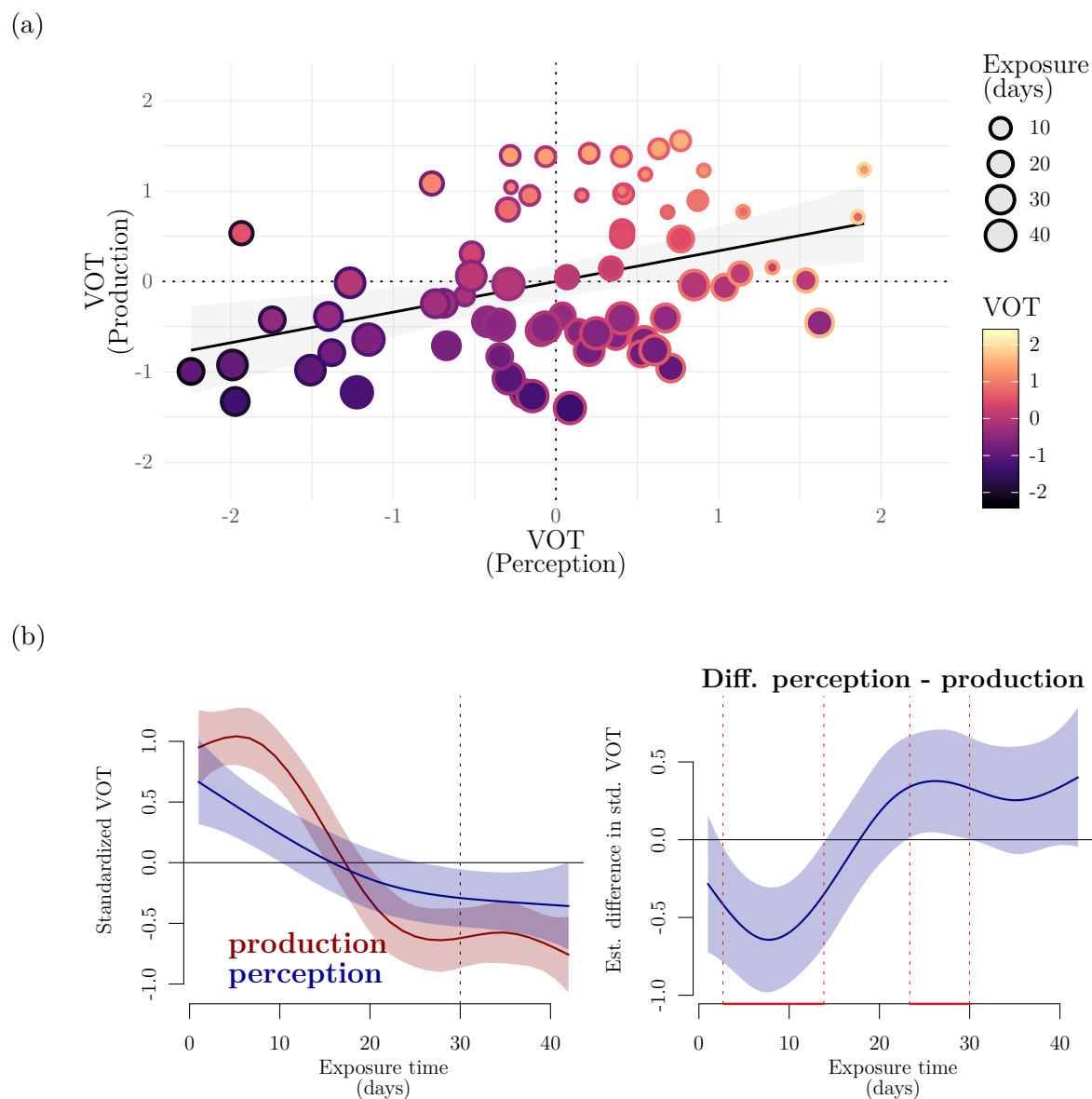


Figure 6. Scatterplot of production and perception boundaries (panel a), non-linear trajectories of production and perception boundaries (panel b, left-side), and estimated difference (voiceless - voiced) in standardized VOT as a function of exposure time (panel b, right-side). In panel (b) the black, discontinuous vertical bar highlights the 30-day mark of exposure time and the red, discontinuous vertical bars represent time windows of significant differences.

Appendix

Model output

Parametric coefficients:

	Estimate	Std. Error	t value	p-value
Intercept	0.004	0.075	0.048	> 0.05
Intercept Δ Voiced	0	0.102	0	> 0.05

Approximate significance of Smooth terms:

	EDF	Ref. DF	F	p-value
Reference smooth: exposure time	1.882	1.95	39.891	< 0.01
Difference smooth: voiced	1	1	2.715	> 0.05
Random smooth: exposure time x participant	3.57	28	0.209	> 0.05

$R^2 = 0.57$; Deviance explained: 59.7%

n = 140

Table A1

Summary of production GAMM model output. The full model fit VOT as a function of voicing.