

**Voice onset time and vowel formant measures in online testing and laboratory-based testing with(out) surgical face masks**

Antje Stoehr,<sup>1,a</sup> Christoforos Souganidis,<sup>1,2</sup> Trisha B. Thomas,<sup>1,3</sup> Jessi Jacobsen,<sup>1</sup> and Clara D. Martin<sup>1,4</sup>

<sup>1</sup> *Basque Center on Cognition, Brain and Language, Donostia–San Sebastián, 20009, Spain*

<sup>2</sup> *University of the Basque Country, Vitoria–Gasteiz, 01006, Spain*

<sup>3</sup> *University of the Basque Country, Leioa, 48940, Spain*

<sup>4</sup> *Ikerbasque Basque Foundation for Science, Bilbao, 48009, Spain*

Since the COVID-19 pandemic started, conducting experiments online is increasingly common, and face masks are often used in everyday life. It remains unclear whether phonetic detail in speech production is captured adequately when speech is recorded in internet-based experiments or in experiments conducted with face masks. We tested 55 Spanish–Basque–English trilinguals in picture naming tasks in three conditions: online, laboratory-based with surgical face masks, and laboratory-based without face masks (control). We measured plosive voice onset time (VOT) in each language, the formants and duration of English vowels /i:/ and /ɪ/, and the Spanish/Basque vowel space. Across conditions, there were differences between English and Spanish/Basque VOT and in formants and duration between English /i:/–/ɪ/; between conditions, small differences emerged. Relative to the control condition, the Spanish/Basque vowel space was larger in online testing and smaller in the face mask condition. We conclude that testing online or with face masks is suitable for investigating phonetic detail in within-participant designs although the precise measurements may differ from those in traditional laboratory-based research.

<sup>a</sup> Email: a.stoehr@bcbl.eu

1 **I. INTRODUCTION**

2 The COVID-19 pandemic has posed challenges to conducting laboratory-based psycholinguistic  
3 research and caused research studies to move online. In 2020, many psycholinguistic research  
4 laboratories around the world remained closed and internet-based studies were the only option for  
5 collecting data. When laboratories reopened, participants were often obliged to wear face masks  
6 during experiments. As of the time this article was written, many countries have lifted official mask  
7 mandates, and people infected with COVID-19 are generally no longer required to isolate  
8 themselves. However, it is still strongly advised to wear face masks in several countries if you have  
9 symptoms consistent with COVID-19 or have been in contact with infected people to reduce  
10 airborne disease transmission. Given the proven efficacy of face masks in reducing the spread of  
11 respiratory diseases, it is likely that their use may be mandated again in the future to ensure public  
12 health and safety. Both online testing and the use of face masks in on-site research challenge  
13 phonetically oriented research, as the phonetic properties of speech may be altered. In this paper, we  
14 present a systematic comparison of a set of phonetic properties in Spanish–Basque–English  
15 trilinguals’ speech production elicited in three conditions: online, in the laboratory with surgical face  
16 masks, and the laboratory without face masks. We examine the phonetic detail through measures of  
17 voice onset time (VOT) and vowel formants, the most widely used measures in bi-/multilingual  
18 phonetic research (Cabrelli Amaro and Wrembel, 2016; Hansen Edwards and Zampini, 2008).

19 **A. Face masks in speech production studies**

20 To date, only two published studies have investigated the direct influence of face masks on the  
21 phonetic properties of vowel production (Bond et al., 1989; Georgiou, 2022a) and none have  
22 investigated the consequences on plosive production. Long before the COVID-19 pandemic, Bond  
23 et al. (1989) found that wearing oxygen masks reduced the (American English) vowel space,

24 especially along the first formant (F1) dimension. This finding may be explained by the physical  
25 barrier face masks constitute. Face masks may restrict jaw and lip movements, which may limit the  
26 F1 range of vowels and affect the articulation of labial consonants, respectively (Saedi et al., 2015).  
27 Thus, a restricted jaw movement can explain the Bond et al. finding of a reduced vowel space along  
28 the F1 dimension. However, it is important to note that the oxygen masks used by Bond et al. were  
29 quite distinctive from types of masks that are commonly used today to prevent disease transmission.  
30 Therefore, their findings might not be generalizable to other (more flexible) types of masks. More  
31 recently, Georgiou (2022a) studied the effect of the now commonly used surgical and cotton face  
32 masks on phonetic detail in the production of the Cypriot Greek vowels /i e a o u/. Unlike Bond et  
33 al., Georgiou (2022a) did not find that F1 was detectably altered by wearing either face mask. This  
34 could mean that the more flexible face masks used in Georgiou’s study did not restrict jaw  
35 movement as did the more static oxygen masks in Bond et al.

36 Though Georgiou (2022a) did not find evidence of altered F1, he did find that wearing face masks  
37 affected the production of vowels along the second formant (F2) dimension. However, not all  
38 vowels were affected similarly: surgical face masks were associated with increased F2 in /e/ and /u/  
39 and decreased F2 in /a/, but cotton face masks were associated with decreased F2 in /e/ and /a/.  
40 This could be the result of face masks filtering certain frequencies, and allowing others to pass. A  
41 systematic comparison of the acoustic attenuation caused by various types of surgical, respirator,  
42 cloth, transparent, and shield face masks showed that face masks generally attenuate frequencies  
43 above 1 kHz (2 kHz for surgical face masks) with the strongest attenuation above 4 kHz (Corey et  
44 al., 2020). This attenuation may affect F1 as well as the higher frequency F2. Across all face masks,  
45 surgical face masks provided the best acoustic performance, which may explain why Georgiou  
46 (2022a) did not observe effects on F1. However, in Georgiou (2022a), alterations in F2 were not  
47 limited to high frequency as reported by Corey et al. (2020). Instead, also the relatively low F2 of

48 /u/ (surgical face masks) and intermediate F2 of /a/ (surgical and cotton face masks) were affected,  
49 which means that the filtering properties reported in Corey et al. (2020) cannot fully explain the  
50 Georgiou (2022a) results.

51 In contrast to the sparse research into the effect of face masks on speech production, the  
52 consequences of face masks on speech perception have received considerable attention. This  
53 perception research supports the idea that wearing face masks during speech production impacts  
54 phonetic speech properties, which may lead to either impaired speech intelligibility (Atcherson et al.,  
55 2017; Corey et al., 2020; Goldin et al., 2020; Magee et al., 2020) or enhanced speech intelligibility, at  
56 least under certain conditions (Cohn et al., 2021 for enhanced intelligibility with face masks during  
57 clear speech production; Pycha et al., 2022 under noise (but see Toscano and Toscano, 2021); Zellou  
58 et al., 2023 for enhanced intelligibility with face masks for coarticulatory vowel nasalization in lax  
59 vowels). Impaired speech perception can be explained by the filtering and attenuating properties of  
60 face masks (Bond et al., 1989; Corey et al., 2020; Georgiou, 2022) and, if applicable, the lack of  
61 visual information (Lalonde and Werner, 2019). Improved intelligibility of face-masked speech can  
62 be explained by speakers' compensation for the face masks' restricting properties. They may be  
63 speaking louder (Asadi et al., 2020) producing Lombard speech (Bond et al., 1989)—which is  
64 characterized by speaking more loudly, with higher fundamental frequency and longer vowel  
65 durations—and by articulating more clearly than when not wearing face masks, at least under certain  
66 conditions (Cohn et al., 2021; Pycha et al., 2022; Zellou et al., 2023). In this sense, an altered acoustic  
67 signal in face-masked speech may not only result from the physical properties of face masks  
68 themselves, but also from the psychological state of the speaker who may change their speaking style  
69 in order to adapt to the face mask.

70 In sum, wearing face masks likely affects the phonetic properties of speech due to the physical  
71 characteristics of the masks themselves and/or their psychological effect on the speaker. This may

72 lead to unrepresentative phonetic measurements during data analysis, which can cause researchers to  
73 reach faulty conclusions about the speech system. Moreover, we are not aware of any studies that  
74 have systematically investigated the effect of face masks on the production of various speech sounds  
75 in a multilingual’s languages. The present study investigates the effect of commonly used surgical  
76 face masks on trilinguals’ vowel and plosive production.

### 77 **B. Online testing in speech production studies**

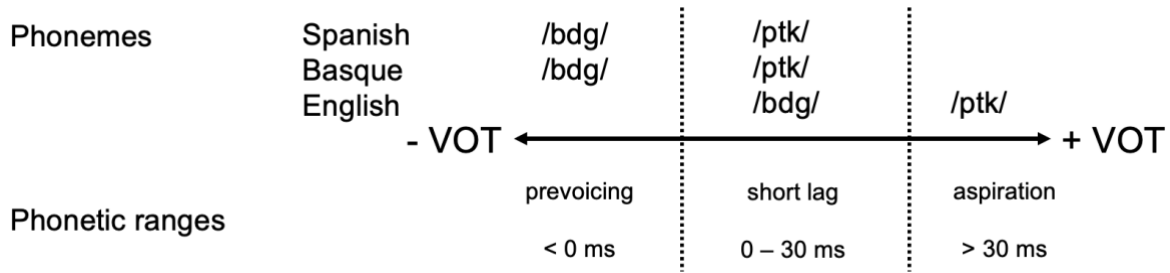
78 Online testing is an attractive option for speech production research. First, speech data can be safely  
79 acquired in an environment where no face masks are required. Second, independent of safety  
80 concerns, larger populations can be accessed, leading to better-powered research studies. The latter  
81 would be especially beneficial for bi-/multilingualism research, which is generally restricted by the  
82 local availability of participants, often resulting in underpowered studies (Brysbaert, 2021). However,  
83 testing participants online means that speech is recorded using diverse devices in diverse  
84 environments, which may cause variability, especially between participants. There is growing  
85 evidence that online studies relying on reaction time and speech onset time measures can reliably  
86 detect well known psycholinguistic effects, such as the word frequency effect (e.g., Anwyl-Irvine et  
87 al., 2020; de Leeuw and Motz, 2016; Fairs and Strijkers, 2021; Hilbig, 2016; Vogt et al., 2022).  
88 To our knowledge, only a small number of studies investigated the suitability of online audio  
89 recordings for phonetic research. Four relatively small-scale studies with English-speaking  
90 participants suggest that the use of different remote recording devices affects phonetic analyses of  
91 vowel production (Bulgin et al., 2010; Calder et al., 2022; Freeman and De Decker, 2021; Zhang et  
92 al., 2021). In all studies, participants recorded themselves with several devices or applications  
93 simultaneously, meaning any differences in vowel measurements could be attributed to the recording  
94 device or application rather than the characteristics of the specific production. Zhang et al. had  
95 participants produce isolated vowels and record themselves at home with three devices in parallel:

96 (1) the built-in microphones on their laptops running the Zoom cloud meeting application (32 kHz  
97 .m4a files), (2) the built-in microphones of their mobile phones using the Awesome Voice Recorder  
98 application (256-bps, 44.1 kHz .wav files) and the Recorder application (256-bps, 44.1 kHz .ogg  
99 files), and (3) a high-quality Zoom H6 Handy Recorder (24-bit, 44.1 kHz .wav files). Vowel formants  
100 were extracted using Praat (version 6.1.08; Boersma and Weenink, 2019) and VoiceSauce (Shue,  
101 2010). In laptop recordings using the Zoom cloud meeting application, F1, F2, and F3 were lower  
102 than when recorded with the high-quality H6 recorder (F3 was only lower when extracted with  
103 Praat). In mobile phone recordings, only F2 was lower than in recordings made with the high-quality  
104 H6 recorder when extracted with VoiceSauce but not when extracted with Praat. The participants in  
105 the Freeman and De Decker (2021) study (one female and one male American English speaker)  
106 recorded themselves in the laboratory using five devices while reading a word list: a high-quality H4n  
107 Pro field recorder (16-bit, 44.1 kHz .wav files) and four identical iPads using different recording  
108 applications: (1) the offline Voice Memos application (32-bit, 44.1 kHz .m4a files), (2) the Zoom  
109 cloud meeting application (32-bit, 32 kHz .m4a files), (3) the Microsoft Skype application (32-bit, 16  
110 kHz .mp4 files), or (4) the Microsoft Teams application (32-bit, 16 kHz .mp4 files). The Zoom  
111 recordings were also saved by two receivers using either a high-quality, medium-quality, or low-  
112 quality internet connection. Critical measures were F1 and F2 to assess vowel space shape and vowel  
113 overlap patterns, as well as spectral tilt to assess vowel nasalization. The measurements varied by  
114 recording device/application and transmission most strongly for the female speaker and in  
115 frequencies between 750 and 1500 Hz. However, these differences were generally small, and a larger  
116 sample size would be needed to draw generalizable conclusions. Just recently, Calder et al. (2022)  
117 tested 18 English-speaking participants with various language backgrounds who recorded  
118 conversations and isolated words elicited in a reading task in their homes using the Zoom cloud  
119 meeting application (32 kHz .m4a files) and also portable audio recorders (15/18 Olympus, 1/18

120 TASCAMDR-100MKII, 1/18 VoicetracerDVT 2050, 1/18 Philips VTR8060; all 16-bit, 44.1 kHz  
121 .wav files) and a SLINT omnidirectional condenser lavalier lapel microphone. Critical measures were  
122 vowel F1 and F2. Vowel formants were analyzed either raw, normalized using the Lobanov method  
123 or the Watt and Fabricius modified method. In Zoom recordings, raw F1 was lower and raw F2 was  
124 higher compared to recordings done with the portable recorder. Using Lobanov normalization, no  
125 effect of the recording condition was detected. Using the Watt and Fabricius modified  
126 normalization, both F1 and F2 were higher when recorded with Zoom.  
127 To summarize, online speech recordings may fundamentally improve access to multilingual  
128 communities, but it is possible that phonetic measurements taken from speech recorded online  
129 differ from recordings made under controlled conditions in the laboratory. These differences may  
130 arise from differences in sampling rate, internet connection, and from differences in the recording  
131 environment. The present study tests whether online studies are suited to investigating phonetic  
132 detail in trilinguals' speech production through measures of VOT in plosives and formants in  
133 vowels.

### 134 C. VOT

135 VOT is the time interval between a plosive's burst release and voicing onset and the most important  
136 acoustic differentiator of phonologically voiced from phonologically voiceless plosives (Lisker and  
137 Abramson, 1964). The VOT continuum (Figure 1) can generally be divided into three phonetic  
138 ranges: prevoicing (negative VOT), short lag (short positive VOT, usually <30 ms) and aspiration  
139 (long positive VOT >30 ms, usually around 70 ms).



140

141 FIG. 1. VOT in Spanish, Basque, and English.

142 Spanish, (Standard)<sup>1</sup> Basque, and English have a voicing contrast between voiceless /p t k/ and  
 143 voiced /b d g/. English, however, differs from Spanish and Basque in the phonetic implementation  
 144 of the voicing contrast (e.g., Lisker and Abramson, 1964; Souganidis et al., 2022). English /p t k/ are  
 145 produced with aspirated VOT, but both Spanish and Basque /p t k/ are produced with short lag  
 146 VOT. Bi-/multilingual speakers are generally known to produce language-specific VOT for voiceless  
 147 plosives if their languages differ in the phonetic implementation of voicing (among many others:  
 148 Flege, 1987, 1991; Stoehr et al., 2017). Importantly, although Spanish–Basque–English trilinguals  
 149 often do not produce monolingual-like aspiration in English, they generally produce longer VOT in  
 150 English than in Spanish and Basque (Stoehr et al., 2023). English /b d g/ fall within the short lag  
 151 range, but Spanish and Basque /b d g/ are produced with prevoicing, that is, negative VOT. Native  
 152 speakers of true-voicing languages like Spanish or Basque often carry over prevoicing to voiced  
 153 plosives in their aspirating nonnative language, but they may differ in the proportion of voiced  
 154 plosives produced with prevoicing (for Dutch–German bilinguals: Stoehr et al., 2017; for German–  
 155 Italian–English trilinguals: Geiss et al., 2022) or they may produce distinct prevoicing durations in  
 156 their true-voicing and aspirating languages (for Portuguese–English bilinguals: Osborne and  
 157 Simonet, 2021).

158 In this study, we test whether the expected crosslinguistic VOT production differences between  
 159 English and Spanish/Basque are detectable when speech production is elicited online or while



160 participants wear surgical face masks in the laboratory. We hypothesize that Spanish–Basque–  
161 English trilinguals produce language-specific VOT in each condition. If this hypothesis is true, we  
162 predict that Spanish–Basque–English trilinguals produce voiceless plosives with longer VOT in  
163 English than in Spanish and Basque in each condition and a smaller proportion of voiced plosives  
164 with prevoicing in English than in Spanish and Basque in all conditions. However, it is possible that  
165 face masks constitute a physical barrier that shortens the duration of aspiration in English voiceless  
166 plosives, which may reduce the VOT production difference between English and Spanish/Basque  
167 when participants wear surgical face masks. We further hypothesize that prevoicing in voiced  
168 plosives cannot always be measured in online testing because it is a subtle acoustic signal that may  
169 not be captured by all recording devices and in uncontrolled environments, which may lead to a  
170 lower proportion of prevoiced plosives in the online condition.

#### 171 **D. Vowel formants**

172 Formants refer to resonant frequencies of the vocal tract and are one of the primary acoustic cues  
173 for distinguishing vowels (Peterson and Barney, 1952). The first formant (F1) corresponds to vowel  
174 height and is correlated with tongue height and jaw position, such that vowels produced with a  
175 higher tongue and a more closed jaw position have smaller F1. The second formant (F2)  
176 corresponds to vowel backness and is correlated with the length of the vocal tract, such that vowels  
177 produced further back in the mouth have smaller F2. Formants are usually defined by automatic  
178 tracking algorithms like the one used by Praat (Boersma and Weenink, 2021). Formants are  
179 traditionally measured in Hertz (Hz) or on the psychoacoustical Bark scale (Zwicker, 1961).

180 The Spanish and Basque vowel inventories comprise the same five vowels /i e a o u/ (Hualde, 1991;  
181 Ladefoged and Johnson, 2010). The vowels /i a u/ form the vowel space, which is delimited by the  
182 distance between /i/–/a/, /a/–/u/, and /u/–/i/. The vowel space size is an important measure in  
183 various disciplines, including speech development (Flipsen and Lee, 2012; Pettinato et al., 2016),

184 speech directed to infants (Rattanasone et al., 2013) and foreigners (for a review: Piazza et al., 2022),  
185 clinical linguistics (Sapir et al., 2010; Skodda et al., 2012), and sociolinguistics (Fox and Jacewicz,  
186 2017; Pierrehumbert et al., 2004). Among measures for determining vowel space size, a particularly  
187 promising measure is the Formant Centralization Ratio (FCR), which maximizes sensitivity to vowel  
188 centralization and reduces inter-speaker variability (Sapir et al., 2010).

189 In the present study, we use FCR measures to test whether the size of the Spanish and Basque  
190 vowel space differs when participants' speech is recorded online or while they wear surgical face  
191 masks in the laboratory. As such, this research question focuses on a general influence of the testing  
192 condition on speech production in the (near) native languages and does not address multilinguals'  
193 speech production per se. As Spanish and Basque have the same vowel inventory, we consider these  
194 two languages together. We hypothesize that the size of the Spanish/Basque vowel space in  
195 Spanish–Basque–English trilinguals differs by condition. If this hypothesis is true, we predict that  
196 surgical face masks restrict the jaw to some extent, thereby resulting in a smaller vowel space when  
197 participants wear face masks compared to when they do not (Bond et al., 1989). Since the online  
198 recordings are made with various recording devices, we predict differences in vowel space size in  
199 speech elicited online versus in the laboratory without face masks, but it is unclear whether the  
200 different recording devices result in a smaller or larger vowel space.

201 The English vowel inventory is considerably larger than the Spanish and Basque vowel inventories,  
202 but the number and type of vowels differ by variety and dialect. The production of certain English  
203 vowel contrasts, such as the contrast between tense /i:/ (e.g., in “sheep”) and lax /ɪ/ (e.g., in “ship”)   
204 are reportedly difficult for native speakers of Spanish and other languages lacking this contrast (e.g.,  
205 Cebrian, 2007; Cebrian et al., 2021; Georgiou, 2022b). The vowels /i:/ and /ɪ/ differ in three  
206 dimensions: vowel height (/i:/ is higher/has smaller F1), vowel backness (/i:/ is more frontal/has  
207 larger F2), and duration (/i:/ is longer). Importantly, although speakers of languages lacking the

208 /i:/–/ɪ/ contrast differ from native English speakers’ production of /ɪ/ (e.g., Cebrian et al., 2021),  
209 they generally distinguish /i:/ and /ɪ/ in production to some extent, either producing distinct  
210 formants (Georgiou, 2022b) and/or duration (Cebrian, 2007; Cebrian et al., 2021; Georgiou, 2022b).  
211 In this study, we test whether the /i:/–/ɪ/ contrast is measurable in Spanish–Basque–English  
212 trilinguals’ speech elicited online or while they wear surgical face masks in the laboratory. We  
213 hypothesize that Spanish–Basque–English trilinguals produce English /i:/ and /ɪ/ distinctly in each  
214 condition. If this hypothesis is true, we predict that Spanish–Basque–English trilinguals’ production  
215 of English /i:/ and /ɪ/ differs in at least one of the following three measures in each condition: F1,  
216 F2, or duration. However, we expect the exact formant values to differ by condition. If face masks  
217 reduce the vowel space size (Bond et al., 1989) as predicted above, it may reduce the distance  
218 between /i:/ and /ɪ/ in the F1–F2 space. If this hypothesis is true, we predict that the formant  
219 differences between /i:/ and /ɪ/ will be smaller in the face mask condition than in the control  
220 condition.  
221 In online testing, F1 is expected to be smaller than in the laboratory without face masks (Calder et  
222 al., 2022; Zhang et al., 2021) and F2 may be smaller (Zhang et al., 2021) or larger (Calder et al.,  
223 2022); in laboratory-based testing with surgical face masks, F2 is expected to be larger than in  
224 laboratory-based testing without face masks (Georgiou, 2022a).

## 225 II. METHODS

### 226 A. Participants

227 Fifty-five Spanish–Basque–English trilinguals participated (41 women,  $M_{age} = 25.15$  years,  $SD_{age} =$   
228 5.90 years, range 18–39 years; see Section *Statistical analyses* for sample size determination).  
229 Participants reported their ages of acquisition for each language to research assistants trained to  
230 obtain this information (Table I). Forty-eight participants acquired Spanish from birth and Basque

231 during childhood. Five participants acquired Spanish and Basque within their first year of life, and  
 232 two acquired Basque from birth and Spanish at age 1 or 2 years, respectively. All participants learned  
 233 English as a foreign language through formal instruction in school and reported no active  
 234 knowledge and use of other languages. Participants lived in the vicinity of Donostia–San Sebastián  
 235 in the Basque Autonomous Community in Spain at the time of testing.  
 236 Participants were recruited from the Basque Center on Cognition, Brain and Language’s (BCBL)  
 237 participant pool. As part of the BCBL’s participant pool registration process, participants complete  
 238 the Basque–English–Spanish Test (BEST; de Bruin et al., 2017), which measures three proficiency  
 239 components, namely vocabulary knowledge through picture naming, word recognition through  
 240 lexical decisions in line with the original LexTALE (Lemhöfer and Broersma, 2012), and general  
 241 language proficiency through semi-structured interviews guided by a multilingual linguist and scored  
 242 on a Likert-like scale from 1 (“lowest level”) to 5 (“native-like level”). Across measures, the recruited  
 243 participants had ceiling proficiency in Spanish, intermediate to high proficiency in Basque and  
 244 intermediate proficiency in English. Their self-reported exposure to Spanish was highest, followed  
 245 by Basque and then by English (Table 1).

246

247 TABLE I. Participant characteristics.

	Spanish			Basque			English		
	M	SD	Range	M	SD	Range	M	SD	Range
AoA <sup>a</sup> (years)	0.05	0.30	0–2	2.83	1.80	0–9	5.75	1.90	2–10
Vocabulary (0–65)	64.84	0.54	62–65	52.24	9.56	24–65	44.38	11.12	11–63

Word recognition (% correct)	93.75	6.20	73–100	86.11	7.98	49–97	67.16	8.47	47.5–88.75
Interview (1–5)	5	0	5–5	4.04	0.74	3–5	3.13	0.55	2–4
Self-reported exposure (%)	63.64	14.45	30–90	25.82	13.57	0–60	10.38	7.13	0–30

248 <sup>a</sup> Age of acquisition.

249 **E. Apparatus and general procedure**

250 Participants completed four sessions. They first completed an online familiarization session in which  
 251 they saw all pictures paired with their written and auditory forms to enhance naming congruence  
 252 during the test phases. Next, participants completed Spanish, Basque, and English picture naming  
 253 tasks (PNT; blocked by language) in three conditions, each administered in different sessions: online,  
 254 on-site in the laboratory with surgical face masks, and on-site in the laboratory without face masks  
 255 (hereafter, control condition). The order of sessions was counterbalanced. Within each condition,  
 256 Spanish and Basque blocks were counterbalanced and the English block was always administered  
 257 last. We chose this order to reflect our participants’ use of Spanish and Basque (but not English) in  
 258 their day-to-day interactions. We argue that the influence of Spanish and Basque on English would  
 259 persist even if the English block were presented first and furthermore presenting the English block  
 260 first would unnaturally influence the next languages.

261 The online familiarization phase and the online condition were programmed in jsPsych (de Leeuw,  
 262 2015) using the open-source study management system JATOS (Lange et al., 2015). Fifty  
 263 participants reported completing the online condition on a laptop and five on a desktop computer.  
 264 Forty-three reported using the microphone integrated into their laptop (as instructed) and 12  
 265 reported using an external microphone. In the on-site conditions, participants were tested

266 individually in sound-attenuating chambers at the BCBL's satellite laboratory at the University of the  
267 Basque Country in Donostia–San Sebastián. The PNT was run on a laptop computer (HP EliteBook  
268 Folio 1040 G3) using OpenSesame software (version 3.2.7 Kafkaesque Koffka; Mathôt et al., 2012).  
269 Voice recordings were made with a Samson C01U PRO professional USB condenser microphone  
270 (Samson Technologies, Hicksville, NY) (set to 100%). In the face mask condition, participants were  
271 provided with a standard surgical face mask type IIR with bacterial filtration efficiency  $\geq 98\%$ ,  
272 which they wore fully covering the mouth and nose. In all test conditions, participants' responses  
273 were saved as .wav files with a 44.1 kHz sampling rate. At the beginning of each session, participants  
274 gave informed consent. At the beginning of the online familiarization session, they performed a  
275 microphone check; at the end of the online condition, they completed a questionnaire. After the  
276 control condition, they also completed a reading aloud task for a different project. Participants were  
277 compensated with €36 paid via bank transfer or PayPal and three stamps on their fidelity card (ten  
278 stamps merit an additional gift). The BCBL's Ethics Committee approved the study.

## 279 **F. Materials**

280 The Spanish and Basque PNT included 43 words each and the English PNT included 44 words (see  
281 supplementary material<sup>2</sup>). The Spanish and Basque word lists comprised 24 items for VOT analysis  
282 and 30 items for vowel analysis (11 items were used for both analyses). The English word list  
283 comprised 24 items for VOT analysis and another 20 items for vowel analysis. All words were  
284 repeated once, resulting in a total of 86 Spanish and Basque productions each and 88 English  
285 productions per condition.

286 In each language, the 24 VOT items were composed of four items per plosive (/b/, /d/, /g/, /p/,  
287 /t/, /k/). These items had a plosive–vowel onset, were one or two syllables long, and had first-  
288 syllable stress. Across languages, items were matched for the number of phonemes and syllables, and  
289 for the vowel following the plosive. As the English vowel inventory differs from the Spanish and

290 Basque vowel inventories and since nonnative speakers are often influenced by orthography when  
291 producing nonnative vowels (for a review: Hayes-Harb and Barrios, 2021), we mostly matched the  
292 words on the orthographic vowel (e.g., Basque “porru” /poru/ leak, Spanish “pollo” /poʎo/  
293 chicken and English “pocket” /pʌkɪt/ were considered matched on the vowel).

294 The 30 Spanish and Basque vowel items were composed of ten items per corner vowel (/i/, /a/,  
295 /u/). These corner vowels appeared in the stressed position of the word. Since we did not expect  
296 the vowel space size to differ between Spanish and Basque, we measured the vowel space size for  
297 both languages combined; therefore, we did not match the vowel stimuli on any variables across  
298 languages. The 20 English vowel items consisted of 10 (near) minimal pairs between /i:/ and /ɪ/,  
299 such as “sheep” and “ship”.

300 Throughout, items with the lowest cognate rate possible between Spanish and English  
301 ( $M_{\text{SpanishItems}} = 0.17$ ;  $M_{\text{EnglishItems}} = 0.15$ ) and between Basque and English ( $M_{\text{BasqueItems}} = 0.10$ ;  $M_{\text{EnglishItems}}$   
302  $= 0.16$ ) were selected (0 = no orthographic overlap; 1 = full orthographic overlap; supplementary  
303 material<sup>2</sup>). Items were represented by color drawings selected from the MultiPic database  
304 (Duñabeitia et al., 2018) when they were available (Spanish 32/43; Basque 26/43; English 24/44),  
305 and the remaining pictures were selected from open content online sources.

## 306 **G. Procedure**

307 Each trial began with a fixation cross in the center of the screen for 500 ms. Afterwards, a picture  
308 appeared for 4000 ms, during which the recorder was active. There were two naming cycles in each  
309 language. Within each, the pictures were presented in random order. Participants were offered  
310 breaks in between the different language blocks. The PNT took around 25 min.

311 **H. Analyses**

312 ***A Acoustic analyses***

313 Phonetic measurements were taken in Praat Software (version 6.1.08; Boersma and Weenink, 2021).  
314 VOT of voiced plosives was measured as the (negative) interval in milliseconds between the onset of  
315 prevoicing and the release of the burst; VOT of voiceless plosives was measured as the (positive)  
316 interval in milliseconds between the release of the burst and the onset of the following vowel. VOT  
317 measurements were determined through visual inspection of the waveform and the spectrogram  
318 viewed at 0–5000 Hz. For vowel formant analysis, the critical vowels were labeled by hand. Vowel  
319 onset and offset were defined as the first and last reliable glottal pulses with visible formants in the  
320 spectrogram viewed at 0–10 000 Hz. Afterwards, F1 and F2 were extracted at vowel midpoint using  
321 ceilings of 5500 Hz and 5000 Hz for female and male participants, respectively. Vowel duration was  
322 extracted in parallel.

323 ***B Statistical analyses***

324 Statistical data analyses were conducted in RStudio (version IDE 2022.02.2+485; RStudio Team,  
325 2022) run on R (version 4.2.0; R Core Team, 2022) and using the lme4 package (version 1.1-29;  
326 Bates et al., 2015). We obtained  $p$  values for  $t$  statistics through the lmerTest package (version 3.1-3;  
327 Kuznetsova et al., 2017). We used the performance package (version 0.9.0; Lüdtke et al., 2021) to  
328 check model assumptions. In linear mixed-effects models, data points with standardized residuals  
329 more than 2.5 standard deviations from 0 were removed using the LMERCvenienceFunctions  
330 package (version 3.0; Tremblay and Ransijn, 2020). Significant interactions were investigated with  
331 Bonferroni-corrected pairwise comparisons using the emmeans package (version 1.7.4-1; Lenth,  
332 2022). Vowel formants were converted to barks with the barktools package (version 0.2.0; Stanley,  
333 2022), which uses Traunmüller (1990) formula. Data were visualized using the ggplot2 package



334 (version 3.3.6; Wickham, 2016). The result tables for all analyses are provided in the supplementary  
335 material<sup>2</sup>.

336 A power calculation using the pwr package (version 1.3-0; Champely et al., 2020) showed that a  
337 sample of 55 participants was needed to reach 80% power for a medium effect size in the linear  
338 regression investigating changes in the Spanish/Basque vowel space by condition. A databased  
339 power calculation using the MixedPower package (version 0.1.0; Kumle et al., 2021) and data for the  
340 first 15 participants showed that a sample size of 55 provided high power for the dependent  
341 variables of the remaining research questions (i.e., 90% power for the fixed effect *Language* in the  
342 linear mixed-effects model on voiceless plosives; 99% power for the fixed effect *Language* in the  
343 logistic regression on voiced plosives; 84% power for the fixed effect *Vowel* in the linear mixed-  
344 effects model on English vowels; all based on 1000 simulations). Adding more participants did not  
345 increase power for the fixed effect *condition*, the interactions between *Language* and *condition* or *Vowel*  
346 and *condition*, which is why 55 was selected as the sample size.

### 347 III. RESULTS

#### 348 I. Plosive production

349 There were 23 760 possible productions (55 participants  $\times$  3 languages  $\times$  3 conditions  $\times$  48  
350 productions). Recording problems for one participant in the Spanish and English online conditions  
351 resulted in the loss of 96 trials (0.4% of the data). Another 1471 trials (453 trials in the control  
352 condition, 479 trials in the face mask condition, and 539 trials in the online condition; 6.22% of the  
353 data) were excluded from the analyses because of an incorrect response (wrong word or no  
354 response) or because VOT could not be measured reliably, for example due to background noise  
355 masking the onset of the burst and/or voicing or due to coarticulation (e.g., “um basket”), which  
356 does not allow for determining the voicing onset.

357 **C** *Voiceless plosives*

358 The analysis tested whether participants produced /p/, /t/, /k/ with longer VOT in English than  
359 in Spanish and Basque in all conditions. Because VOT reportedly differs by plosive (/p/ < /t/ < /k/;  
360 e.g., Volaitis and Miller, 1992; Stoehr et al., 2023 for a similar population), the factor *Plosive* was  
361 included in the model. The linear mixed-effects model had *VOT in ms* as the dependent variable  
362 with fixed effects for *Language*, *condition*, and *Plosive*, as well as an interaction between *Language* and  
363 *condition*. The model included random intercepts for *Participant* and *Item*, as well as by-*Participant*  
364 random slopes for *Language* and *condition*. No by-*Item* random slope for *condition* was included, as it  
365 caused a singularity warning [lmer formula:  $VOT \sim Language * condition + Plosive$   
366  $+ (1 + Language + condition | Participant) + (1 | Item)$ ]. 334 outliers (2.90% of the data) were removed  
367 (see Section *Statistical analyses*). We used Helmert contrast coding for the three-level categorical  
368 variable *Language* to create two contrasts of interest. The first contrast (hereafter, *Language\_ESB*)  
369 compared the difference between English and the mean of Spanish and Basque. The second  
370 contrast (hereafter, *Language\_SB*) compared the difference between the means of Spanish and  
371 Basque. This coding scheme provided the maximal power to test for a difference between English  
372 versus Spanish and Basque (Schad et al., 2020). We used deviation coding for the three-level variable  
373 *Plosive* to create two contrasts of interest. The first contrast (hereafter, *Plosive\_pt*) compared /p/  
374 [0.5] to /t/ [-0.5], and the second contrast (hereafter, *Plosive\_tk*) compared /t/ [-0.5] to /k/ [0.5].  
375 These two contrasts allowed us to capture the predicted VOT increase from /p/ to /t/ and from  
376 /t/ to /k/ (Volaitis and Miller, 1992). The same coding scheme was used for *condition* to compare  
377 the face mask condition [0.5] to the control condition [-0.5] (*condition\_Mask*) and the online  
378 condition [0.5] to the control condition [-0.5] (*condition\_Online*).

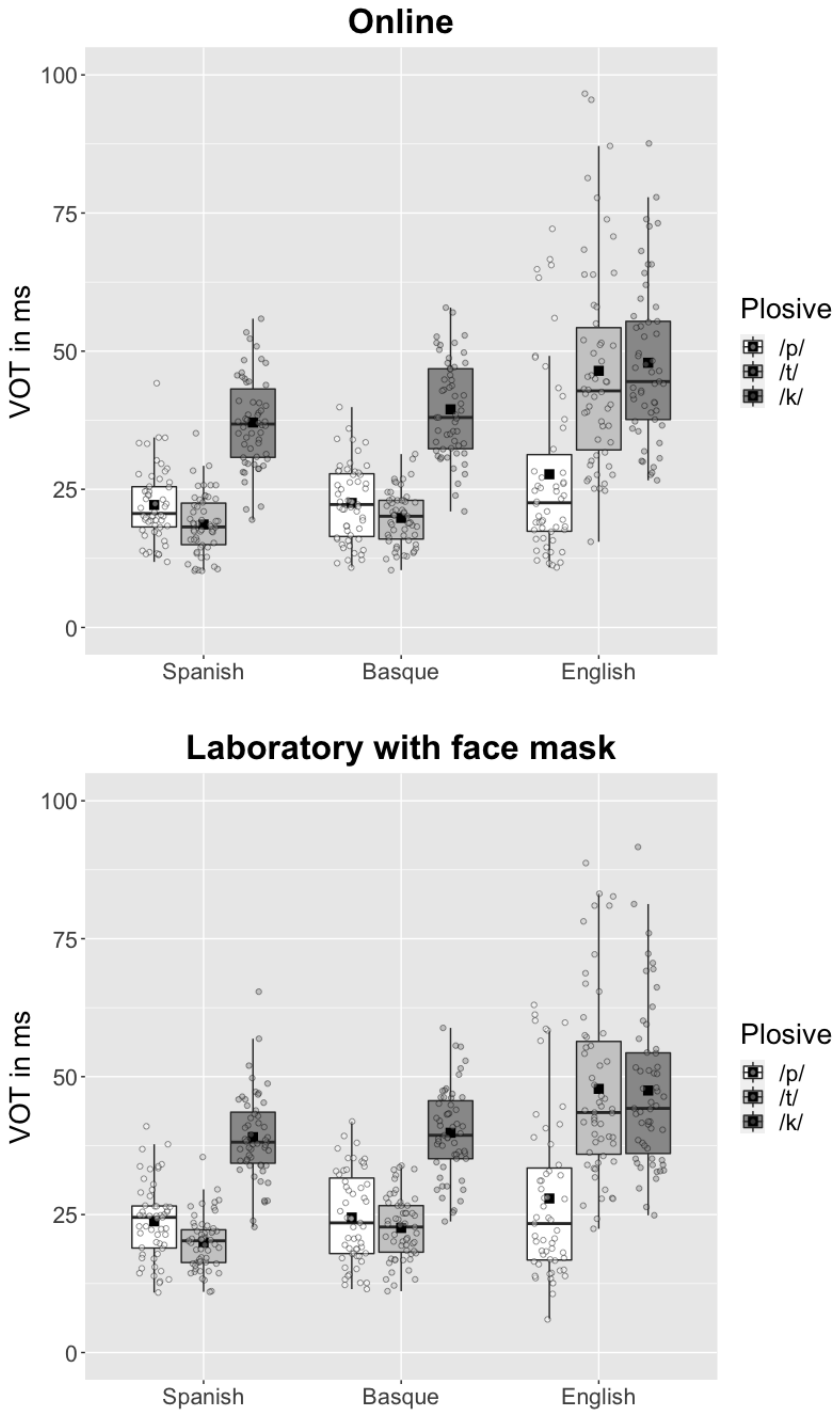
379 Participants produced shorter VOT for /p/ than /t/ ( $\beta = -14.609$ ; standard error,  $SE = 3.347$ ;  $t$   
380  $= -4.365$ ;  $p < 0.001$ ) and shorter VOT for /t/ than /k/ ( $\beta = 20.175$ ;  $SE = 3.347$ ;  $t = 6.028$ ;  $p < 0.001$ ).

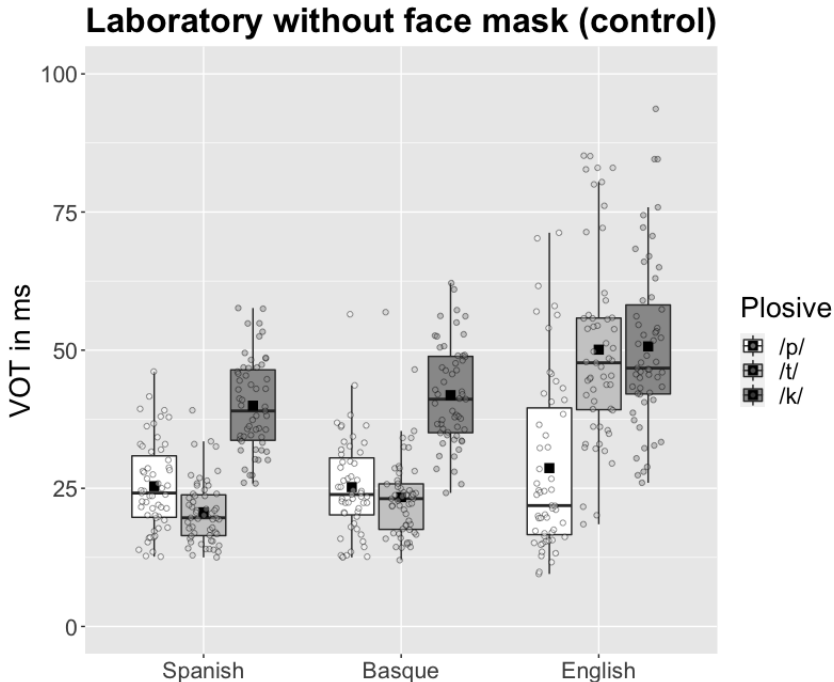
381 A significant effect of *Language\_ESB* showed that participants produced longer VOT in English than  
382 in Spanish and Basque ( $\beta = 13.419$ ;  $SE = 3.198$ ;  $t = 4.196$ ;  $p < 0.001$ ). No significant effect of  
383 *Language\_SB* was observed, suggesting that Spanish and Basque VOT were not detectably different  
384 ( $\beta = -1.476$ ;  $SE = 2.903$ ;  $t = -0.508$ ;  $p = 0.615$ ). A significant effect of *condition\_Online* showed that  
385 VOT recorded online was shorter than control ( $\beta = -2.483$ ;  $SE = 0.743$ ;  $t = -3.341$ ;  $p = 0.002$ ).  
386 There was no detectable effect of *condition\_Mask* ( $\beta = -0.173$ ;  $SE = 0.507$ ;  $t = -0.341$ ;  $p = 0.735$ ).  
387 The model detected significant interactions between *Language\_ESB* and *condition\_Mask* ( $\beta = -1.532$ ;  
388  $SE = 0.584$ ;  $t = -2.625$ ;  $p = 0.009$ ) and *Language\_ESB* and *condition\_Online* ( $\beta = 1.157$ ;  $SE = 0.589$ ;  $t$   
389  $= 1.966$ ;  $p = 0.049$ ). No other significant interactions were observed. The results are visualized in  
390 Figure 2.

391 Pairwise comparisons by *condition* confirmed that English VOT was longer than Spanish and  
392 Basque VOT in all conditions (online: English vs. Spanish:  $\beta = 14.67$ ;  $SE = 3.54$ ;  $t = 4.141$ ;  $p < 0.001$ ;  
393 English vs. Basque:  $\beta = 13.33$ ;  $SE = 3.52$ ;  $t = 3.791$ ;  $p = 0.001$ ; laboratory-based with face mask:  
394 English vs. Spanish:  $\beta = 13.43$ ;  $SE = 3.54$ ;  $t = 3.792$ ;  $p = 0.001$ ; English vs. Basque:  $\beta = 11.88$ ;  $SE =$   
395  $3.52$ ;  $t = 3.379$ ;  $p = 0.004$ ; control: English vs. Spanish:  $\beta = 14.38$ ;  $SE = 3.54$ ;  $t = 4.060$ ;  $p < 0.001$ ;  
396 English vs. Basque:  $\beta = 12.84$ ;  $SE = 3.52$ ;  $t = 3.652$ ;  $p = 0.002$ ). Pairwise comparisons by *Language*  
397 showed that English VOT was affected differently than Spanish and Basque VOT by condition.  
398 Compared to control, we found a) English VOT was shorter in both the laboratory-based condition  
399 with face mask ( $\beta = -2.050$ ;  $SE = 0.549$ ;  $t = -3.734$ ;  $p < 0.001$ ) and the online condition ( $\beta = 2.308$ ;  
400  $SE = 0.725$ ;  $t = 3.184$ ;  $p = 0.006$ ), b) Spanish and Basque VOT were shorter only in the online  
401 condition (Spanish:  $\beta = 2.600$ ;  $SE = 0.714$ ;  $t = 3.642$ ;  $p = 0.001$ ; Basque:  $\beta = 2.799$ ;  $SE = 0.713$ ;  $t =$   
402  $3.927$ ;  $p < 0.001$ ).

403 The results demonstrate that experiments conducted online and in the laboratory with surgical  
404 face masks are suitable for detecting VOT differences between aspirating and true-voicing languages

405 in voiceless plosives. However, VOT duration is shorter when recorded online across languages  
406 compared to control. Furthermore, English aspiration is shorter when participants wear surgical face  
407 masks compared to control.





408 FIG. 2. VOT by language and plosive in the online (top), laboratory-based with face mask (middle),  
409 and laboratory-based without face mask (control; bottom) conditions. Each dot shows an individual  
410 participant; the black square shows the mean; the horizontal line shows the median.

#### 411 ***D Voiced plosives***

412 The analysis tested whether participants produced a larger proportion of /b/, /d/, /g/ with  
413 positive VOT in English than in Spanish and Basque in all conditions<sup>3</sup>. The logistic mixed-effects  
414 model had the proportion of devoiced productions (positive VOT coded as 1; negative VOT coded  
415 as 0) as a dependent variable. The three-level variable *Plosive* was coded as two contrasts of interest.  
416 The first contrast (hereafter, Plosive\_bd) compared /b/ [0.5] to /d/ [-0.5]; and the second  
417 (hereafter, Plosive\_dg) compared /d/ [-0.5] to /g/ [0.5]). The remainder of the model and the  
418 coding schemes were the same as in the model on voiceless plosives reported above [glmer formula:

---

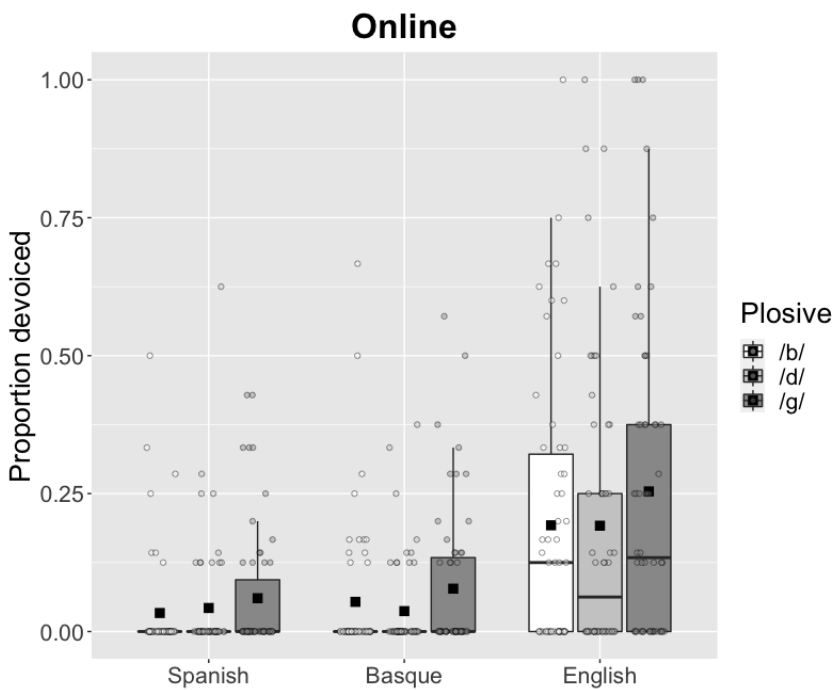
<sup>3</sup> We also ran a linear mixed-effects model testing for prevoicing duration differences across languages and conditions. This model did not detect any significant differences in prevoicing duration (supplementary material<sup>4</sup>).

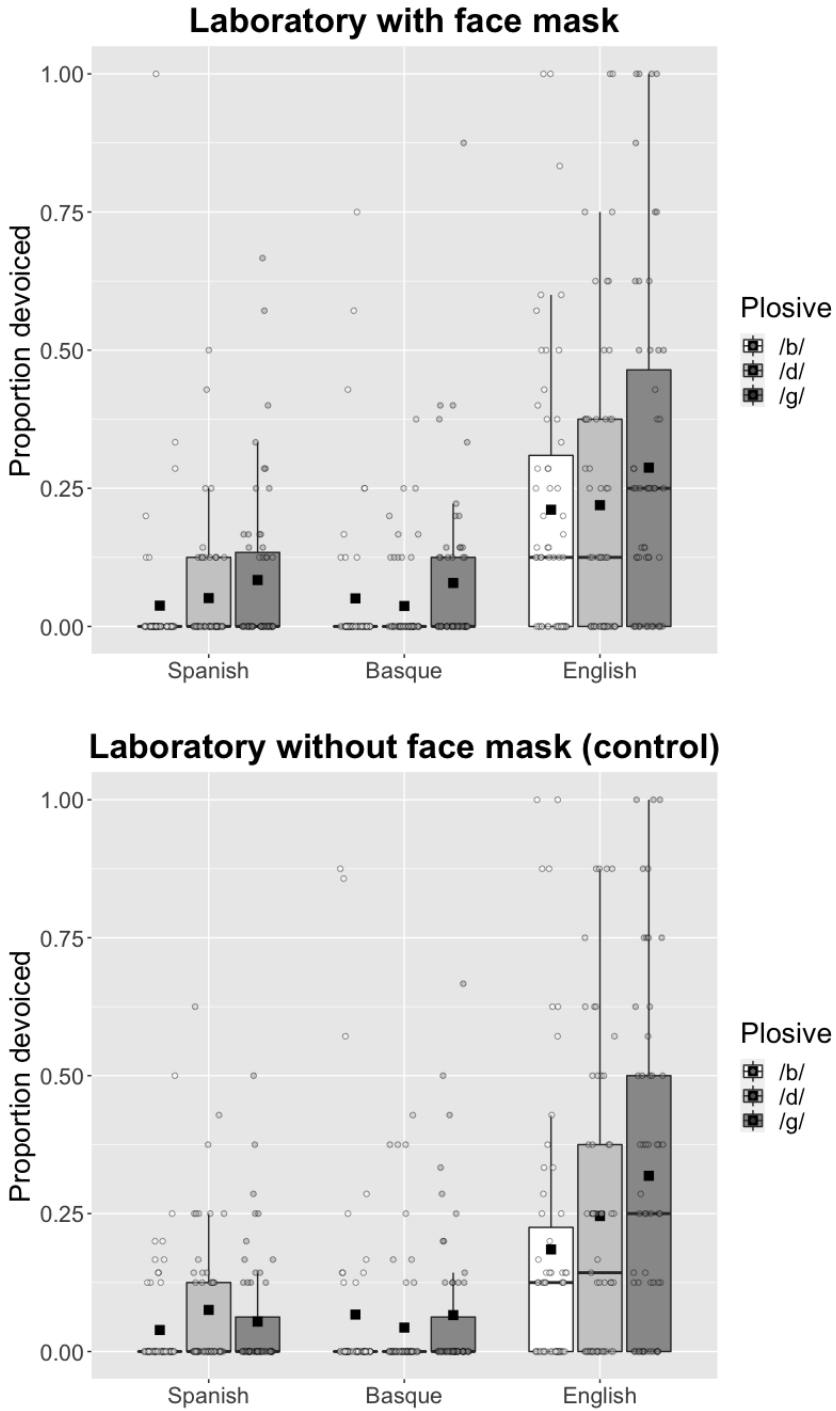
419 Proportion

420 Devoiced~Language\*condition+Plosive+(1+Language+condition | Participant)+(1 | Item)].

421 Participants produced a larger proportion of /b/, /d/, /g/ with positive VOT in English than in  
422 Spanish and Basque ( $\beta = 2.113$ ;  $SE = 0.274$ ;  $\chi = 7.703$ ;  $p < 0.001$ ). In addition, participants devoiced  
423 /g/ more frequently than /d/ ( $\beta = 0.694$ ;  $SE = 0.253$ ;  $\chi = 2.742$ ;  $p = 0.006$ ). No other significant  
424 main effects or interactions were observed. Results are visualized in Figure 3.

425 These results suggest that experiments conducted online and in the laboratory with surgical face  
426 masks are suitable for detecting differences in the proportion of voiced plosives produced with  
427 positive VOT between aspirating and true-voicing languages.





428 FIG. 3. Proportion of devoiced productions in the online (top), laboratory-based with face mask  
429 (middle), and laboratory-based without face mask (control; bottom) conditions. Each dot shows an  
430 individual participant; the black square shows the mean; the horizontal line shows the median.

431 **J. Vowel production**

432 In Spanish and Basque, there were 19 800 possible productions (55 participants × 2 languages × 3  
433 conditions × 60 productions). In English, there were 6600 possible productions (55 participants × 2  
434 vowels × 3 conditions × 20 productions). Recording problems for one participant in the Spanish  
435 and English online conditions resulted in loss of all 60 Spanish and all 40 English trials (1.52% of  
436 the data). Across conditions, 79 Spanish (0.80% of the data), 168 Basque (1.70% of the data), and  
437 138 English trials (2.10% of the data) were excluded from the analyses because a participant failed to  
438 respond within the time limit, produced a wrong word or there was interference from background  
439 noise.

440 To exclude any formant tracking errors, we removed productions with z-scored formant values  
441 larger than 3 or smaller than -3 (Kirkham and McCarthy, 2021). This resulted in removal of 338  
442 Spanish (3.46% of the data) and 361 Basque (3.71% of the data) trials in the vowel space analysis  
443 and 76 English trials in the F1 analysis (1.18% of the data) and 108 English trials in the F2 analysis  
444 (1.68% of the data). For all languages combined, F1 tracking errors amounted to 1.16% of the data  
445 in the online condition, 1.39% of the data in the control condition, and 1.59% of the data in the face  
446 mask condition. F2 formant tracking errors amounted to 1.80% of the data in the control condition,  
447 2.16% of the data in the face mask condition, and 2.36% of the data in the online condition.

448 ***E Production of the English /i:/-/ɪ/ contrast across conditions***

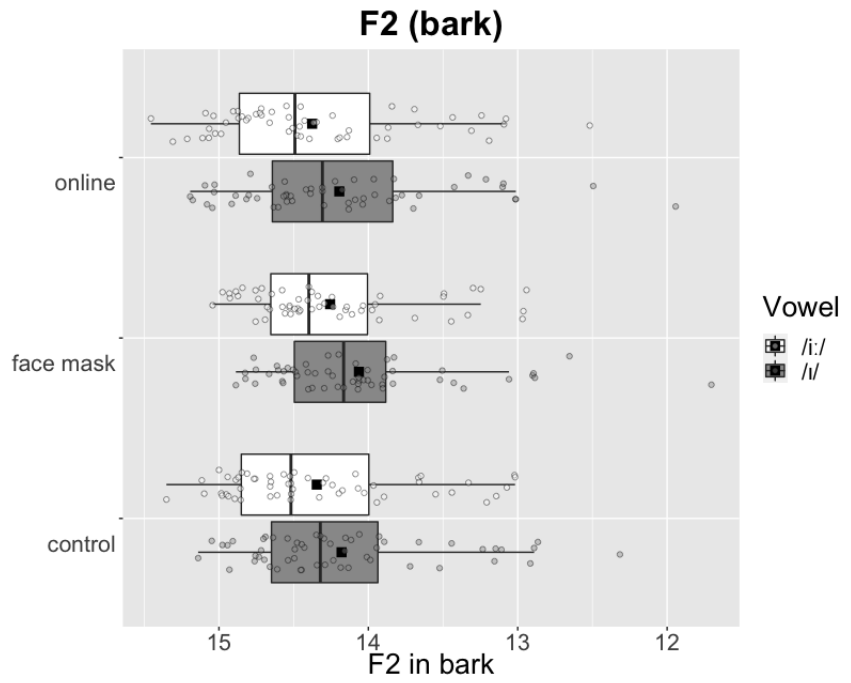
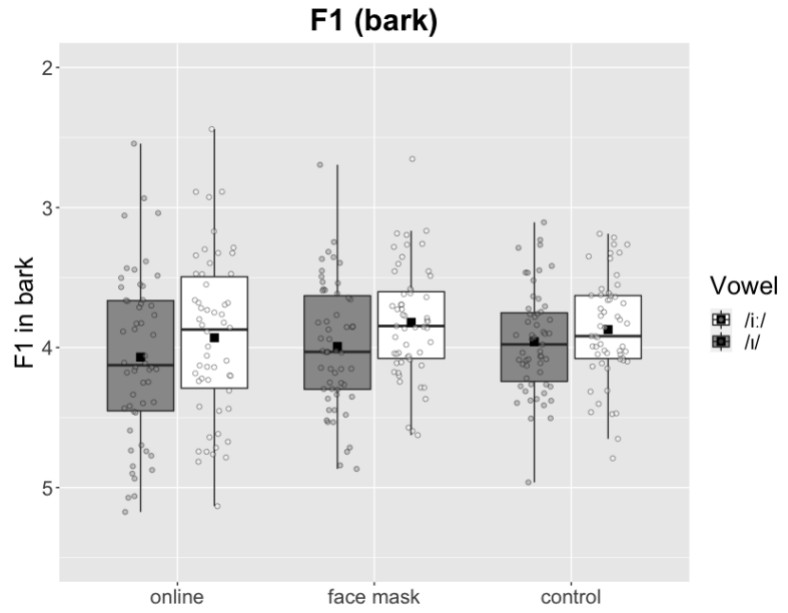
449 The analyses tested whether participants produced the English vowels /i:/ and /ɪ/ distinctly in all  
450 conditions. We fitted three linear mixed-effects models with *F1 (bark)*, *F2 (bark)*, and *duration (ms)* as  
451 dependent variables. All models had fixed effects for *Vowel* and *condition* as well as an interaction  
452 term between *Vowel* and *condition*. The models included random intercepts for *Participant* and *Item*, as  
453 well as by-*Participant* random slopes for *Vowel* and *condition* [lmer formula:

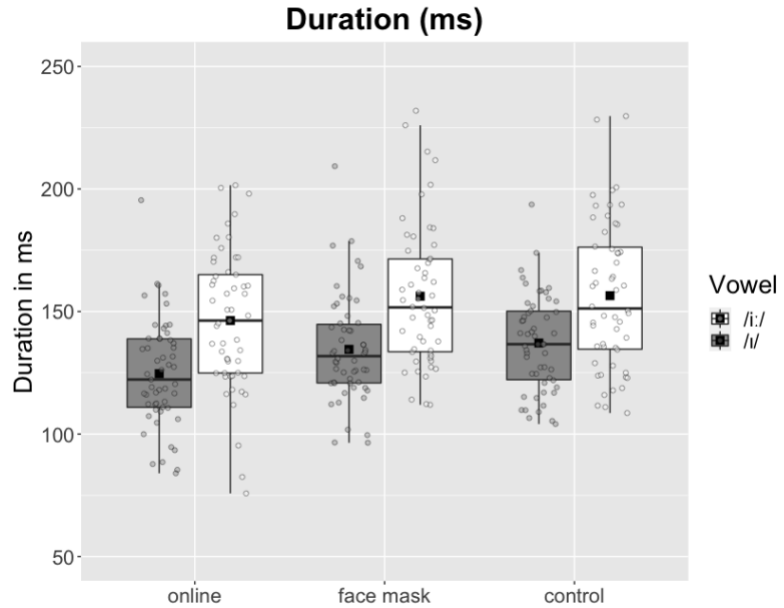
454  $F1/F2/Duration \sim Vowel * condition + (1 + Vowel + condition | Participant) + (1 | Item)$ ]. We used



455 deviation coding for the binary variable *Vowel* (/I/ coded as -0.5; /i:/ coded as 0.5). As in the VOT  
456 models, we used deviation coding for *condition* to compare the laboratory-based condition with face  
457 mask [0.5] to control [-0.5] (*condition\_Mask*) and the online condition [0.5] to control [-0.5]  
458 (*condition\_Online*). 195 outliers (3.07% of the data) were removed in the F1 analysis, 177 (2.80% of  
459 the data) in the F2 analysis, and 114 (1.78% of the data) in the duration analysis (see Section  
460 *Statistical analyses*).

461 Across measures, participants produced /i:/ and /I/ distinctly (Figure 4). Significant results are  
462 presented below. The F1 model detected that participants produced /i:/ with smaller F1 than /I/ ( $\beta$   
463 = -0.133;  $SE = 0.054$ ;  $t = -2.473$ ;  $p = 0.018$ ). There was a significant interaction between *Vowel* and  
464 *condition\_Mask* ( $\beta = -0.084$ ;  $SE = 0.028$ ;  $t = -2.993$ ;  $p = 0.003$ ). Pairwise comparisons by *condition*  
465 showed that the F1 difference between /i:/ and /I/ was significant when participants wore a  
466 surgical face mask in the laboratory and when they were tested online but not in the control  
467 condition (laboratory-based with face mask:  $\beta = 0.175$ ;  $SE = 0.055$ ;  $t = 3.150$ ;  $p = 0.003$ ; online:  $\beta =$   
468  $0.140$ ;  $SE = 0.056$ ;  $t = 2.510$ ;  $p = 0.016$ ; control:  $\beta = 0.084$ ;  $SE = 0.055$ ;  $t = 1.518$ ;  $p = 0.137$ ). The  
469 F2 model detected that participants produced /i:/ with larger F2 than /I/ ( $\beta = 1.797e-01$ ;  $SE =$   
470  $8.278e-02$ ;  $t = 2.171$ ;  $p = 0.038$ ). Moreover, participants produced smaller F2 in the laboratory-based  
471 condition with face mask than in control ( $\beta = -1.501e-01$ ;  $SE = 3.912e-02$ ;  $t = -3.836$ ;  $p < 0.001$ ).  
472 The duration model detected that /i:/ productions were longer than /I/ productions ( $\beta = 20.481$ ;  
473  $SE = 4.675$ ;  $t = 4.381$ ;  $p < 0.001$ ). In addition, vowel duration was longer in the laboratory-based  
474 condition with face mask compared to control ( $\beta = 6.065$ ;  $SE = 2.541$ ;  $t = 2.387$ ;  $p = 0.021$ ) and  
475 shorter in the online condition compared to control ( $\beta = -14.953$ ;  $SE = 3.401$ ;  $t = -4.397$ ;  $p < 0.001$ ).  
476 Overall, experiments conducted online and in the laboratory with surgical face masks are suitable  
477 for detecting small formant and duration differences in Spanish–Basque–English trilinguals’  
478 production of the /i:/–/I/ contrast.





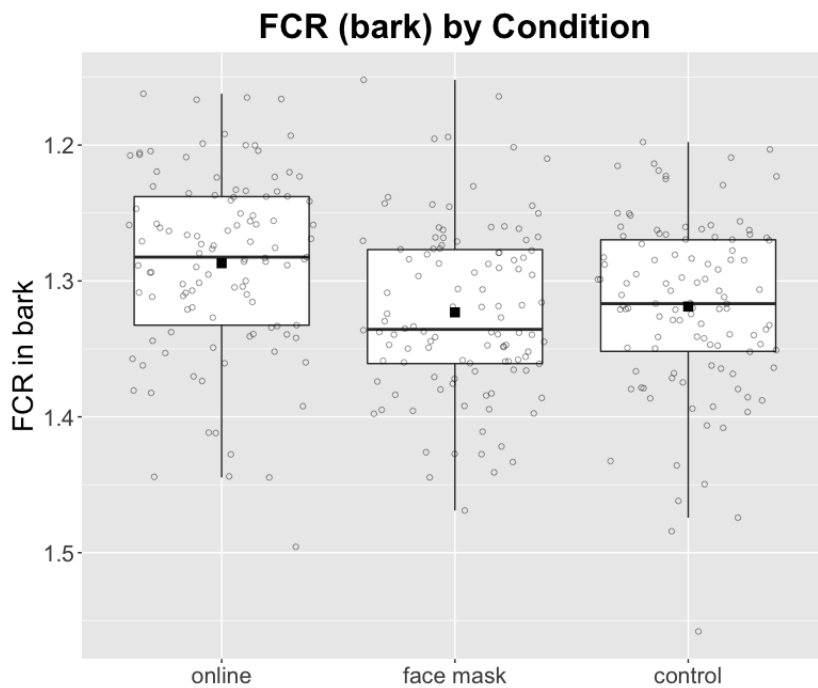
480 FIG. 4. F1 (top), F2 (middle), and duration (bottom) of English /i:/ and /ɪ/ by condition. Each dot  
481 shows an individual participant; the black square shows the mean; the horizontal line shows the  
482 median.

### 483 *F The Spanish/Basque vowel space across conditions*

484 This analysis tested whether the size of the Spanish/Basque vowel space differed by condition. We  
485 calculated the vowel space for each participant in each language and condition as the FCR (Sapir et  
486 al., 2010) in bark, expressed as  $(F2/u/+F2/a/+F1/i/+F1/u)/(F2/i/+F1/a/)$ . A larger FCR  
487 means that the vowel space is smaller (more centralized), and a smaller FCR means that the vowel  
488 space is larger (less centralized). We fitted a linear regression model with the *FCR in bark* as the  
489 dependent variable. The model had *condition* as fixed effect. *Condition* was deviation coded to  
490 compare the laboratory-based condition with face mask [0.5] to control [-0.5] (*condition\_Mask*) and  
491 the online condition [0.5] to control [-0.5] (*condition\_Online*; lm formula:  $FCR \sim condition$ ). As  
492 noted in the Section *Vowel formants*, we did not predict the Spanish and Basque vowel spaces—which  
493 are composed of the same vowels—to be differently affected by the testing conditions. Therefore,

494 we averaged the results across Spanish and Basque and did not include *Language* as a factor in the  
495 model. The vowel space was smaller in the laboratory-based condition with surgical face mask ( $\beta =$   
496  $0.027$ ;  $SE = 0.010$ ;  $t = 2.683$ ;  $p = 0.008$ ) and larger in the online condition ( $\beta = -0.045$ ;  $SE = 0.010$ ;  $t$   
497  $= -4.506$ ;  $p < 0.001$ ) compared to control (Figure 5).

498 These results show that compared to control, experiments conducted online are associated with a  
499 larger vowel space, and experiments conducted while participants wear surgical face masks are  
500 associated with a smaller vowel space.



501  
502 FIG. 5. FCR by condition as a measure of the Spanish/Basque vowel space size. Each dot shows an  
503 individual participant; the black square shows the mean; the horizontal line shows the median. The  
504 smaller the FCR the larger the vowel space.

#### 505 IV. DISCUSSION

506 The present study investigated whether recording participants' speech while they wear surgical face  
507 masks in the laboratory and recording their speech online using jsPsych (de Leeuw, 2015) and

508 JATOS (Lange et al., 2015) are reliable options when investigating phonetic detail in speech  
509 production. To that end, we compared these two methods to speech production elicited on-site in  
510 the laboratory without face masks. We focused on phonetic detail through measures of VOT in  
511 voiceless and voiced plosives, and vowel formants in isolated words produced by Spanish–Basque–  
512 English trilingual adults.

### 513 **K. Plosive production across conditions**

514 Production differences between English versus Spanish and Basque were present in all conditions:  
515 participants produced voiceless plosives (/p/, /t/, /k/) with longer VOT in English than in Spanish  
516 and Basque and produced voiced plosives (/b/, /d/, /g/) more frequently with positive VOT in  
517 English than in Spanish and Basque, thus confirming our predictions. As such, testing participants in  
518 the laboratory when surgical face masks are required or testing them online are suitable options for  
519 investigating crosslinguistic differences in plosive production.

520 However, the exact VOT duration of voiceless plosives differed by condition. In the present  
521 study, wearing surgical face masks reduced VOT duration of English—but not Spanish and  
522 Basque—voiceless plosives by 2 ms on average compared to control. This finding shows that  
523 surgical face masks specifically affect the duration of aspiration (in English voiceless plosives) but  
524 not the duration of voiceless plosives in general, as Spanish and Basque short lag voiceless plosives  
525 were not affected by participants wearing face masks. This is likely because surgical face masks are  
526 positioned close to the lips and act like a physical barrier to the aspiration air stream passing through  
527 the lips. Short lag VOT as common in Spanish and Basque is likely too short to be affected by this  
528 physical barrier. Therefore, the finding that surgical face masks reduce aspiration duration in English  
529 but not short lag VOT in Spanish and Basque appears to be reflective of the phonetic characteristics  
530 of English as an aspirating language and Spanish/Basque as true-voicing languages rather than being  
531 the result of differences in language proficiency between languages. Importantly, this shortening of

532 English aspirated VOT in the face mask condition did not affect the crosslinguistic VOT difference  
533 between English and Spanish/Basque, which remained significant.

534 In online testing, across languages and relative to the control condition (in the laboratory  
535 without a face mask) VOT in voiceless plosives was on average 3 ms shorter. Participants were likely  
536 more relaxed during the online session conducted in their homes, which may have led to more  
537 natural—and thus more representative—VOT production. The formal laboratory environment may  
538 have imposed more pronunciation effort, thus leading to longer VOT production than in the online  
539 condition. Previous research partly supports this assumption (Robb et al., 2005): native English  
540 speakers produced longer syllable durations in speech recorded in the laboratory compared to  
541 speech recorded outside the laboratory; however, VOT duration for these speakers did not  
542 statistically differ by environment. It is possible that the formal laboratory setting affected syllable  
543 and VOT durations but that the relatively small sample size of 20 was not sufficient to detect small  
544 VOT differences by environment (3 ms in the present study) in Robb et al.

545 We did not find evidence for our prediction that the proportion of voiced plosives produced  
546 with prevoicing would differ between the online and control conditions. We predicted that since  
547 prevoicing is a subtle acoustic signal, the uncontrolled environment and recording devices in the  
548 online condition would not capture the presence of prevoicing as reliably as the professional  
549 recorder and environment in the control condition. As we did not detect any differences between  
550 the online and control conditions in the proportion of voiced plosives produced with prevoicing, the  
551 present results are encouraging for online testing, showing that this uncontrolled environment is  
552 suitable for recording subtle acoustic signal differences.

553 In sum, our data show that speech recordings made online and in the laboratory when  
554 participants wear surgical face masks are suitable when investigating VOT production in voiceless  
555 and voiced plosives of multilinguals speaking true-voicing and aspirating languages.

556 **L. Production of the English /i:/–/ɪ/ vowel contrast across conditions**

557 Participants produced the English vowels /i:/ and /ɪ/ with distinct F1, F2, and duration, with /i:/  
558 having a higher (smaller F1) and more frontal (larger F2) position and longer duration than /ɪ/, thus  
559 confirming our predictions and supporting previous findings that native speakers of languages  
560 lacking the /i:/–/ɪ/ contrast produce these vowels with distinct F2 (Georgiou, 2022b) and/or  
561 duration (Cebrian, 2007; Cebrian et al., 2021; Georgiou, 2022b). Unlike previous studies, the present  
562 study—with its larger participant sample (55 in the present study, 10 in Georgiou, 2022b; 30 in  
563 Cebrian, 2007; 43 in Cebrian et al., 2021)—detected production differences between /i:/ and /ɪ/  
564 across all three measures. This may be attributed to the present study being well-powered and thus  
565 able to detect small spectral and temporal differences in vowel production.

566 Against our prediction, the F1 difference between /i:/ and /ɪ/ was larger when participants wore a  
567 face mask (mean difference 0.176 bark) compared to control (mean difference 0.096 bark). In fact, a  
568 post hoc test failed to find an F1 difference between /i:/ and /ɪ/ in the control condition. We  
569 speculate that participants compensated for the communicative restrictions imposed by the face  
570 mask by hyperarticulating, which may have enhanced the F1 difference between /i:/ and /ɪ/ when  
571 participants wore face masks. Our finding that participants produced both vowels with longer  
572 duration when they wore face masks compared to control supports the hyperarticulation  
573 assumption. Previous work reporting enhanced intelligibility of face-masked speech (Cohn et al.,  
574 2021; Pycha et al., 2022; Zellou et al., 2023) further supports that people may be hyperarticulating  
575 when wearing face masks. Overall, when wearing surgical face masks, participants produced both  
576 /i:/ and /ɪ/ with smaller F2, corresponding to a more posterior position compared to control. This  
577 finding is against our prediction, which was based on Georgiou’s (2022a) finding that participants  
578 produce /i:/ with numerically larger F2 when they wear surgical face masks. However, although  
579 Georgiou showed that wearing surgical face masks affects vowel production, his results varied by

580 vowel. Some Cypriot Greek vowels were produced with larger F2 (significant difference for /e/ &  
581 /u/; numerical difference for /i/ & /o/) and others with smaller F2 (/a/) relative to a control  
582 measure without face masks. The more posterior position observed in the present study may result  
583 from the surgical face mask acting as a physical barrier at the front of the mouth, thus pushing the  
584 position of the front vowels /i:/ and /ɪ/ to a slightly posterior position. This is in line with research  
585 reporting a reduced vowel space size when participants wear oxygen face masks (Bond et al., 1989).  
586 A reduced vowel space size means that the F2 of front vowels, such as /i:/ and /ɪ/, becomes  
587 smaller, which is what we observed in the present study.

588 Finally, we observed a shorter vowel duration in the online condition compared to control but  
589 neither F1 nor F2 differed between the online and control conditions. The shorter vowel duration in  
590 online testing is in line with the shorter VOT duration of voiceless plosives in all languages in online  
591 testing discussed above and provides further support for our argument that the formal laboratory  
592 environment imposed more pronunciation effort than online testing, which may affect temporal  
593 properties of speech production. The lack of detectable F1 or F2 differences between the online and  
594 control conditions was unpredicted given the Calder et al. (2022) and Zhang et al. (2021) findings of  
595 smaller F1 and smaller F2 (Zhang et al., 2021) or larger F2 (Calder et al., 2022) in vowels recorded  
596 online using the Zoom cloud meeting application. These differences between the present study and  
597 the Calder et al. and Zhang et al. studies may be related to the different online testing tools used in  
598 the present study (jsPsych/JATOS) and in Calder et al. and Zhang et al. (Zoom cloud meeting  
599 application). Importantly, the Zoom cloud meeting application as used by Calder et al. and Zhang et  
600 al. had a different sampling rate (32 kHz) than their in-person recording devices (44.1 kHz), which  
601 may have contributed to their observed differences between recording conditions. In the present  
602 study, both online and laboratory-based recordings were made at 44.1 kHz, and it is possible that  
603 these identical recording settings minimized between-condition differences.



604 Our data suggest that speech recordings made online and in the laboratory when participants wear  
605 surgical face masks are suitable when investigating the production of the nonnative /i:/–/I/  
606 contrast.

607 **M. The Spanish/Basque vowel space size across conditions**

608 Relative to the control condition, participants' vowel space was smaller when tested in the face mask  
609 condition and larger when tested in the online condition. The reduced vowel space in the face mask  
610 condition was predicted given the previous research finding that wearing oxygen face masks was  
611 associated with a smaller vowel space (Bond et al., 1989). The reason for this smaller vowel space  
612 may be due to the face mask restricting the jaw (and consequently F1) and the length of the vocal  
613 tract (and consequently F2). The assumption of face masks being associated with decreased F2 is  
614 also in line with our finding that participants produced the English front vowels /i:/ and /I/ with  
615 smaller F2 in the face mask condition, indicating a more posterior place of articulation. A reduced  
616 vowel space is associated with less clear and less intelligible speech (Bradlow and Bent, 2002). Our  
617 finding of a smaller vowel space when participants wear surgical face masks, therefore, directly  
618 relates to previous research, which found that speech produced with face masks may be less  
619 intelligible than speech produced without face masks (Atcherson et al., 2017; Corey et al., 2020;  
620 Goldin et al., 2020; Magee et al., 2020).

621 The larger vowel space in the online condition was not unexpected, as we assumed that the vowel  
622 space size differs between speech recorded online and control. However, given the lack of previous  
623 research on this topic, we were unable to predict whether online testing would result in a smaller or  
624 larger vowel space. We assumed that the driving force behind differences in vowel space size  
625 between online testing and control may be related to the use of various recording devices in online  
626 testing. When examining Figure 4, which shows the production of English /i:/ and /I/, there  
627 appears to be considerably more variability between participants' F1 production recorded online

628 compared to control. In the vowel space size analysis, however, we observe similar between-  
629 participant variability in the online and control conditions (Figure 5). This may be because we  
630 measured the vowel space size as the FCR, a measure which reduces between-participant variability  
631 (Sapir et al., 2010). Therefore, the larger vowel space in the online condition does not appear to  
632 result from greater between-participant variability. To test whether the larger vowel space in online  
633 testing may result from larger within-participant variability, which may have pushed the formant  
634 means to more extreme positions, we computed compactness scores for each vowel by condition  
635 and participant. These compactness scores were computed as the standard deviation of the mean of  
636 F1 multiplied by the standard deviation of the mean of F2 multiplied by  $\pi$ , assuming that vowel  
637 categories are elliptical (Kartushina and Frauenfelder, 2014). Surprisingly, vowels in the online  
638 condition were the most compact, followed by the control condition and face mask condition  
639 ( $M_{\text{Online\_Compactness}} = 1.20$ ;  $M_{\text{Control\_Compactness}} = 1.35$ ;  $M_{\text{FaceMask\_Compactness}} = 1.44$ ). It appears, then, that the  
640 larger vowel space size in the online condition (relative to control) does not emerge from between-  
641 or within-participant variability. As an alternative explanation, we propose that participants may have  
642 experienced more psychological stress in the formal laboratory environment than when they  
643 performed the online experiment in their homes. Psychological stress has been found to be  
644 associated with a smaller vowel space size (Karlsson et al., 2000). The present study did not include  
645 measures of the L3-English vowel space and future research can investigate if the vowel space in a  
646 language with relatively low proficiency is similarly or even more strongly affected by differences in  
647 the testing environment as the native language(s). If psychological stress is the driving force behind a  
648 reduced vowel space in laboratory-based testing, it is possible that the reduction is even larger in a  
649 low(er) proficiency language because any stress level is likely enhanced by having to speak in a less  
650 proficient language. However, there is no evidence that the hypothesized psychological stress affects

651 the temporal property VOT differently across languages, as we observed longer VOT in the control  
652 condition than in the online condition in Spanish, Basque, and English alike.

653 In summary, online testing in the home environment without an experimenter may have led to more  
654 natural speech production resulting in shorter VOT duration in all languages, shorter vowel duration  
655 of English /i:/ and /ɪ/, and a larger but more representative vowel space size than observed in the  
656 formal control condition conducted in the laboratory.

## 657 **V. CONCLUSIONS**

658 Testing participants in the laboratory while they wear surgical face masks or recording their speech  
659 online appear to be valid options when investigating phonetic detail in trilinguals' speech  
660 production. Across conditions, we observed the predicted phonetic differences between trilinguals'  
661 languages or within trilinguals' least proficient language. However, small phonetic differences  
662 emerged between conditions. Wearing surgical face masks was associated with shorter aspiration in  
663 English voiceless plosives, a larger F1 difference between English /i:/ and /ɪ/, smaller F2 and  
664 longer duration in English /i:/ and /ɪ/, and a smaller Spanish/Basque vowel space, all compared to  
665 control. When participants wear surgical face masks, two competing forces appear to be at play. On  
666 the one hand, surgical face masks shorten the vocal tract and restrict the articulators. A shortened  
667 vocal tract can explain the observed shorter VOT in English (aspirated) voiceless plosives and the  
668 lower F2 in the English vowels /i:/ and /ɪ/. The combination of a shortened vocal tract and  
669 restriction of the articulators can also explain the smaller vowel space in Spanish/Basque. On the  
670 other hand, participants seem to compensate for the limitations imposed by surgical face masks by  
671 hyperarticulating, which can explain the larger F1 difference and longer duration in English vowels  
672 when participants wore face masks.

673 Online testing was associated with shorter VOT in voiceless plosives in Spanish, Basque, and  
674 English, shorter vowel duration in English /i:/ and /ɪ/, and a larger Spanish/Basque vowel space,  
675 all compared to control. Overall, online testing may make participants feel more at ease, resulting in  
676 a more natural—and more ecologically valid—speaking style, which may have led to shorter VOT in  
677 voiceless plosives, shorter vowel duration, and a larger vowel space compared to control. Future  
678 studies still need to investigate how masked and online studies might differentially affect languages  
679 with different properties (e.g., different vowel space density). Nevertheless, we conclude that testing  
680 trilinguals' production of isolated words while they wear surgical face masks in the laboratory or  
681 record their speech online using jsPsych (de Leeuw, 2015) and JATOS (Lange et al., 2015) are  
682 suitable options for within-participant designs.

### 683 **ACKNOWLEDGMENTS**

684 This work was supported by institutional grants from the Basque Government [BERC 2022–  
685 2025 program] and the Spanish State Research Agency [BCBL Severo Ochoa excellence  
686 accreditation CEX2020-001010/AEI/10.13039/501100011033] awarded to the BCBL. This project  
687 has also received funding from the European Union's H2020 research and innovation program  
688 [Marie Skłodowska-Curie grant agreement No 843533 awarded to AS]; the European Research  
689 Council (ERC) under the European Union's Horizon 2020 research and innovation program [grant  
690 agreement No 819093 to CDM]; the Spanish State Research Agency [BES-2017-082500 to CS;  
691 PID2020-113926GB-I00 to CDM; PID2021-123578NA-  
692 I00/AEI/10.13039/501100011033/FEDER, UE, & FJC2020-044978-I to AS]; and by the Basque  
693 Government's Department of Education [Predoctoral training program for research staff  
694 PRE\_2021\_2\_0006 awarded to TT].

### 695 **AUTHOR DECLARATIONS**

696 **Conflict of Interest**

697 We have no conflicts to disclose.

## 698 **Ethics Approval**

699 The present study has been approved by the Basque Center on Cognition, Brain and Language's  
700 Ethics Committee prior to data collection (approval code 230222ML). Informed consent for  
701 participation and collection of speech data were obtained from all participants.

## 702 **DATA AVAILABILITY**

703 The data that support the findings of this study are openly available in the Open Science Framework  
704 at <http://doi.org/10.17605/OSF.IO/XYH3K>.

705 <sup>1</sup>At least two Basque varieties spoken in France employ aspiration (Zuberoan: Gaminde et al., 2002;  
706 Mounole, 2004; Mixean: Egurtzegi and Carignan, 2020). Here, we focus on Standard Basque spoken  
707 in Gipuzkoa/Spain, for which no aspiration has been found (Souganidis et al., 2022).

708 <sup>2</sup>See supplementary material at <https://doi.org/10.1121/10.0020064> for stimulus materials; for  
709 cognate rate measures; for results tables; and for linear mixed-effects model on prevoicing duration.

710 <sup>3</sup>We also ran a linear mixed-effects model testing for prevoicing duration differences across  
711 languages and conditions. This model did not detect any differences in prevoicing duration  
712 (supplementary material2).

## 713 **REFERENCES**

714 Anwyl-Irvine, A., Massonnié, J., Flitton, A., Kirkham, N., and Evershed, J. K. (2020). “Gorilla in our  
715 midst: An online behavioral experiment builder,” *Behav. Res.* **52**, 388–407.

716 Asadi, S., Cappa, C. D., Barreda, S., Wexler, A. S., Bouvier, N. M., and Ristenpart, W. D. (2020).

717 “Efficacy of masks and face coverings in controlling outward aerosol particle emission from  
718 respiratory activities.” *Sci. Rep.* **10**, 15665.

- 719 Atcherson, S. R, Mendel, L. L., Baltimore, W. J., Patro, C., Lee, S., Pousson, M., and Spann, M. J.  
720 (2017). “The effect of conventional and transparent surgical masks on speech understanding  
721 in individuals with and without hearing loss,” *J. Am. Acad. Audiol.* **28**, 58–67.
- 722 Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015). “Fitting linear mixed-effects models using  
723 lme4,” *J. Stat. Softw.* **67**, 1–48.
- 724 Boersma, P., and Weenink, D. (2019). “Praat: Doing phonetics by computer (Version 6.1.08)  
725 [computer program],” <http://www.praat.org/>.
- 726 Boersma, P., and Weenink, D. (2021). “Praat: Doing phonetics by computer (Version 6.1.40)  
727 [computer program],” <http://www.praat.org/>.
- 728 Bond, Z. S., Moore, T. J., and Gable, B. (1989). “Acoustic-phonetic characteristics of speech  
729 production in noise and while wearing an oxygen mask,” *J. Acoust. Soc. Am.* **85**, 907–912.
- 730 Bradlow, A. R., and Bent, T. (2002). “The clear speech effect for non-native listeners,” *J. Acoust.*  
731 *Soc. Am.* **112**, 272–284.
- 732 Brysbaert, M. (2021). “Power considerations in bilingualism research: Time to step up our game,”  
733 *Bilingualism* **24**, 813–818.
- 734 Bulgin, J., De Decker, P., & Nycz, J. (2010). “Reliability of formant measurements from lossy  
735 compressed audio,” in British Association of Academic Phoneticians Colloquium, March  
736 29–31, London, UK. Available at  
737 [https://research.library.mun.ca/684/1/Bulgin\\_De\\_Decker\\_Nycz\\_2010.pdf](https://research.library.mun.ca/684/1/Bulgin_De_Decker_Nycz_2010.pdf)
- 738 Cabrelli Amaro, J., and Wrembel, M. (2016). “Investigating the acquisition of phonology in a third  
739 language - a state of the science and an outlook for the future,” *Int. J. Multiling.* **13**, 395–409.
- 740 Calder, J., Wheeler, R., Adams, S., Amarelo, D., Arnold-Murray, K., Bai, J., Church, M., Daniels, J.,  
741 Gomez, S., Henry, J., Jia, Y., Johnson-Morris, B., Lee, K., Miller, K., Powell, D., Ramsey-  
742 Smith, C., Rayl, S., Rosenau, S., and Salvador, N. (2022). “Is Zoom viable for sociophonetic

- 743 research? A comparison of in-person and online recordings for vocalic analysis. *Ling.*  
744 *Vanguard* 20200148.
- 745 Cebrian, J. (2007). “Old sounds in a new contrast: L2 production of the English tense-lax vowel  
746 distinction,” *Proc. 16<sup>th</sup> Int. Congress Phonetic Sci.* (pp. 1637–1640).
- 747 Cebrian, J., Gorba, C., and Gavaldà, N. (2021). “When the easy becomes difficult: Factors affecting  
748 the acquisition of the English /i:/-/ɪ/ contrast,” *Front. Commun.* **6**, 660917.
- 749 Champely, S., Ekstrom, C., Dalgaard, P., Gill, J., Weibelzahl, S., Anandkumar, A., Ford, C., Volcic,  
750 R., and De Rosario, H. (2020). “Package ‘pwr: Basic functions for power analysis (Version  
751 1.3-0),” [https://CRAN.R-project.org/package = pwr](https://CRAN.R-project.org/package=pwr)
- 752 Cohn, M., Pycha, A., and Zellou, G. (2021). “Intelligibility of face-masked speech depends on  
753 speaking style: Comparing casual, clear, and emotional speech,” *Cognition* **210**, 104570.
- 754 Corey, R. M., Jones, U., and Singer, A. C. (2020). “Acoustic effects of medical, cloth, and  
755 transparent face masks on speech signals,” *J. Acoust. Soc. Am.* **148**, 2371–2375.
- 756 de Bruin, A., Carreiras, M., and Duñabeitia, J. A. (2017). “The BEST dataset of language  
757 proficiency,” *Front. Psy.* **8**, 522.
- 758 de Leeuw, J. R. (2015). “jsPsych: A JavaScript library for creating behavioral experiments in a web  
759 browser,” *Behav. Res.* **47**, 1–12.
- 760 de Leeuw, J. R., and Motz, B. A. (2016). “Psychophysics in a web browser? Comparing response  
761 times collected with JavaScript and Psychophysics Toolbox in a visual search task,” *Behav.*  
762 *Res.* **48**, 1–12.
- 763 Duñabeitia, J.A., Crepaldi, D., Meyer, A.S., New, B., Pliatsikas, C., Smolka, E., and Brysbaert, M.  
764 (2018). “MultiPic: A standardized set of 750 drawings with norms for six European  
765 languages,” *Quarterly J. Exp. Psy.* **71**, 808–816.

- 766 Egurtzegi, A., and Carignan, C. (2020). “An acoustic description of Mixean Basque,” *J. Acoust. Soc.*  
767 *Am.* **147**, 2791–2802.
- 768 Fairs, A., and Strijkers, K. (2021). “Can we use the internet to study speech production? Yes we can!  
769 Evidence contrasting online versus laboratory naming latencies and errors,” *PLoS One* **16**,  
770 e0258908.
- 771 Flege, J. E. (1987). “The production of ‘new’ and ‘similar’ phones in a foreign language: Evidence for  
772 the effect of equivalence classification,” *J. Phon.* **15**, 47–64.
- 773 Flege, J. E. (1991). “Age of learning affects the authenticity of voice-onset time (VOT) in stop  
774 consonants produced in a second language,” *J. Acoust. Soc. Am.* **89**, 395–411.
- 775 Flipsen, P., and Lee, S. (2012). “Reference data for the American English acoustic vowel space,”  
776 *Clin. Ling. Phon.* **26**, 926–933.
- 777 Fox, R. A., and Jacewicz, E. (2017). “Reconceptualizing the vowel space in analyzing regional dialect  
778 variation and sound change in American English,” *J. Acoust. Soc. Am.* **142**, 444–459.
- 779 Freeman, V., and De Decker, P. (2021). “Remote sociophonetic data collection: Vowels and  
780 nasalization over video conferencing apps,” *J. Acoust. Soc. Am.* **149**, 1211–1223.
- 781 Gaminde, I., Hualde, J. I., and Salaberria, J. (2002). “Zubereraren herskariak: Azterketa akustikoa,”  
782 (“Zuberoa’s plosives: An acoustic study”), *Lapurdum* **7**, 221–236.
- 783 Geiss, M., Gumbsheimer, S., Lloyd-Smith, A., Schmid, S., and Kupisch, T. (2022). “Voice onset time  
784 in multilingual speakers: Italian heritage speakers in Germany with L3 English,” *Stud.*  
785 *Second Lang. Acquis.* **44**, 435–459.
- 786 Georgiou, G. P. (2022a). “Acoustic markers of vowels produced with different types of face masks,”  
787 *Appl. Acoust.* **191**, 108691.
- 788 Georgiou, G. P. (2022b). “The acquisition of /ɪ/–/i:/ is challenging: Perceptual and production  
789 evidence from Cypriot Greek speakers of English,” *Behav. Sci.* **12**, 469.



- 790 Goldin, A., Weinstein, B., and Shiman, N. (2020). “How do medical masks degrade speech  
791 perception?” *Hearing Rev.* **27**, 8–9.
- 792 Hansen Edwards, J. H., and Zampini, M. L. (2008). *Phonology and Second Language Acquisition*  
793 (John Benjamins, Philadelphia).
- 794 Hayes-Harb, R., and Barrios, S. (2021). “The influence of orthography in second language  
795 phonological acquisition,” *Lang. Teach.* **54**, 297–326.
- 796 Hilbig, B. B. (2016). “Reaction time effects in lab- versus web-based research: Experimental  
797 evidence,” *Behav. Res.* **48**, 1718–1724.
- 798 Hualde, J. I. (1991). *Basque Phonology* (Routledge, London).
- 799 Karlsson, I., Banziger, T., Dankovicová, J., Johnstone, T., Lindberg, J., Melin, H., Nolan, F. and  
800 Scherer, K. (2000). “Speaker verification with elicited speaking styles in the VeriVox  
801 project,” *Speech Commun.* **31**, 121–129.
- 802 Kartushina, N., and Frauenfelder, U. H. (2014). “On the effects of L2 perception and of individual  
803 differences in L1 production on L2 pronunciation,” *Front. Psychol.* **5**, 1246.
- 804 Kirkham, S., and McCarty, K. M. (2021). “Acquiring allophonic structure and phonetic detail in a  
805 bilingual community: The production of laterals by Sylheti-English bilingual children,” *Int. J.*  
806 *Bilingual.* **25**, 531–547.
- 807 Kumle, L., Vö, M. L., and Draschkow, D. (2021). „Estimating power in (generalized) linear mixed  
808 models: An open introduction and tutorial in R,” *Behav. Res.* **53**, 2528–2543.
- 809 Kuznetsova, A., Brockhoff, P. B., and Christensen, R. H. B. (2017). “lmerTest Package: Tests in  
810 linear mixed effects models,” *J. Stat. Softw.* **82**, 1–26.
- 811 Ladefoged, P. and Johnson, K. (2010). *A Course in Phonetics* (Cengage Learning, Boston).
- 812 Lalonde, K., and Werner, L. A. (2019). “Infants and adults use visual cues to improve detection and  
813 discrimination of speech in noise,” *J. Speech Lang. Hear. Res.* **62**, 3860–3875.

- 814 Lange, K., Kühn, S., and Filevich, E. (2015). ““Just another tool for online studies” (JATOS): An  
815 easy solution for setup and management of web servers supporting online studies,” *PLoS*  
816 *One* **10**, e0130834.
- 817 Lemhöfer, K., and Broersma, M. (2012). “Introducing LexTALE: A quick and valid Lexical Test for  
818 Advanced Learners of English,” *Behav. Res.* **44**, 325–343.
- 819 Lenth, R. (2022). “emmeans: Estimated marginal means, aka least-squares means (Version 1.7.4-1),”  
820 [https://CRAN.R-project.org/package = emmeans](https://CRAN.R-project.org/package=emmeans)
- 821 Lisker, L., and Abramson, A. (1964). “A cross-language study of voicing in initial stops: Acoustical  
822 measurements,” *Word* **20**, 384–422.
- 823 Lüdecke, D., Ben-Shachar, M., Patil, I., Waggoner, P., and Makowski, D. (2021). “Performance: An  
824 R package for assessment, comparison and testing of statistical models,” *JOSS* **6**, 3139.
- 825 Magee, M., Lewis, C., Noffs, G., Reece, H., Chan, J. C. S., Zaga, C. J., Paynter, C., Birchall, O., Rojas  
826 Azocar, S., Ediriweera, A., Kenyon, K., Caverlé, M. W., Schultz, B. G., and Vogel, A. P.  
827 (2020). “Effects of face masks on acoustic analysis and speech perception: Implications for  
828 peri-pandemic protocols,” *J. Acoust. Soc. Am.* **148**, 3562–3568.
- 829 Mathôt, S., Schreij, D., and Theeuwes, J. (2012). “OpenSesame: An open-source, graphical  
830 experiment builder for the social sciences,” *Behav. Res.* **44**, 314–324.
- 831 Mounole, C. (2004). “Zubererazko herskarien azterketa akustikoa” (“The acoustic analysis of the  
832 plosives of Zuberoa”), *Anuario Del Seminario De Filología Vasca "Julio De Urquijo"* **38**  
833 207–248.
- 834 Osborne, D. M., and Simonet, M. (2021). “Foreign-language phonetic development leads to first-  
835 language phonetic drift: Plosive consonants in native Portuguese speakers learning English  
836 as a foreign language in Brazil,” *Languages* **6**, 112.

- 837 Peterson, G. E., and Barney, H. L. (1952). "Control methods used in a study of the vowels," J.  
838 Acoust. Soc. Am. **24**, 175–184.
- 839 Pettinato, M., Tuomainen, O., Granlund, S., and Hazan, V. (2016). "Vowel space area in later  
840 childhood and adolescence: Effects of age, sex and ease of communication," J. Phon. **54**, 1–  
841 14.
- 842 Piazza, G., Martin, C. D., and Kalashnikova, M. (2022). "The acoustic features and didactic function  
843 of foreigner-directed speech: A scoping review," J. Speech Lang. Hear. Res. **65**, 2896–2918.
- 844 Pierrehumbert, J. B., Bent, T., Munson, B., Bradlow, A. R., and Bailey, J. M. (2004). "The influence  
845 of sexual orientation on vowel production," J. Acoust. Soc. Am. **116**, 1905–1908.
- 846 Pycha, A., Cohn, M., and Zellou, G. (2022). "Face-masked speech intelligibility: the influence of  
847 speaking style, visual information, and background noise," Front. Commun. **7**, 874215.
- 848 R Core Team (2022). R: A language and environment for statistical computing (Version 4.2.0) (R  
849 Foundation for Statistical Computing, Vienna, Austria).
- 850 Rattanasone, N. X., Burnham, D., and Reilly, R. G. (2013). "Tone and vowel enhancement in  
851 Cantonese infant-directed speech at 3, 6, 9, and 12 months of age," J. Phon. **41**, 332–343.
- 852 Robb, M., Gilbert, H., and Lerman, J. (2005). "Influence of gender and environmental setting on  
853 voice onset time," Folia Phon. Logopaed. **57**, 123–133.
- 854 RStudio Team (2022). RStudio: Integrated Development Environment for R (Version IDE  
855 2022.02.2+485) (RStudio, Boston).
- 856 Saeidi, R., Niemi, T., Karpelin, H., Pohjalainen, J., Kinnunen, T., & Alku, P. (2015). "Speaker  
857 recognition for speech under face cover," in Proceedings 16th Annual Conference of the  
858 International Speech Communication Association, September 6–10, Dresden, Germany, pp.  
859 1012–1016.

- 860 Sapir, S., Ramig, L. O., Spielman, J. L., and Fox, C. (2010). “Formant Centralization Ratio (FCR): A  
861 proposal for a new acoustic measure of dysarthric speech,” *J. Speech Lang. Hear. Res.* **53**,  
862 114–125.
- 863 Schad, D. J., Vasishth, S., Hohenstein, S., and Kliegl, R. (2020). “How to capitalize on a priori  
864 contrasts in linear (mixed) models: A tutorial,” *J. Memory Lang.* **110**, 104038.
- 865 Shue, Y.-L. (2010). “The voice source in speech production: Data, analysis and models,” PhD thesis,  
866 University of California, Los Angeles.
- 867 Skodda, S., Grönheit, W., and Schlegel, U. (2012). “Impairment of vowel articulation as a possible  
868 marker of disease progression in Parkinson’s Disease,” *PLoS One* **7**, e32132.
- 869 Souganidis, C., Molinaro, N., and Stoehr, A. (2022). “Bilinguals produce language-specific voice  
870 onset time in two true-voicing languages: The case of Basque-Spanish bilinguals,” *Ling.*  
871 *Appr. Bilingualism* (published online).
- 872 Stanley, J. (2022). “barktools: Functions to help when working with Barks (Version 0.2.0),”  
873 <http://joestanley.github.io/barktools>
- 874 Stoehr, A., Benders, T., van Hell, J. G., and Fikkert, P. (2017). “Second language attainment and first  
875 language attrition: The case of VOT in immersed Dutch–German late bilinguals,” *Second*  
876 *Lang. Res.* **33**, 483–518.
- 877 Stoehr, A., Jevtović, M., de Bruin, A., and Martin, C. D. (2023). “Phonetic and lexical crosslinguistic  
878 influence in early Spanish-Basque-English trilinguals”, *Lang. Learn.* (published online).
- 879 Toscano, J. C., and Toscano, C. M. (2021). “Effects of face masks on speech recognition in multi-  
880 talker babble noise,” *PLoS One* **16**, e0246842.
- 881 Traunmüller, H. (1990). “Auditory scales of frequency representation,” *J. Acoust. Soc. Am.* **88**, 97–  
882 100.

- 883 Tremblay, A., and Ransijn, J. (2020). “Package ‘LMERConvenienceFunctions’. Model selection and  
884 post-hoc analysis for (G)LMER models (version 3.0),” [https://CRAN.R-](https://CRAN.R-project.org/package=LMErConvenienceFunctions)  
885 [project.org/package = LMErConvenienceFunctions](https://CRAN.R-project.org/package=LMErConvenienceFunctions)
- 886 Vogt, A., Hauber, R., Kuhlen, A. K., and Rahman, R. A. (2022). “Internet-based language  
887 production research with overt articulation: Proof of concept, challenges, and practical  
888 advice,” *Behav. Res.* **54**, 1954–1975.
- 889 Volaitis, L. E., and Miller, J. L. (1992). “Phonetic prototypes: Influence of place of articulation and  
890 speaking rate on the internal structure of voicing categories,” *J. Acoust. Soc. Am.* **92**, 723–  
891 735.
- 892 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis* (Springer, New York).
- 893 Zellou, G., Pycha, A., and Cohn, M. (2023). “The perception of nasal coarticulatory variation in  
894 face-masked speech,” *J. Acoust. Soc. Am.* **153**, 1084–1093.
- 895 Zhang, C., Jepson, K., Lohfink, G., and Arvaniti, A. (2021). “Comparing acoustic analyses of speech  
896 data collected remotely,” *J. Acoust. Soc. Am.* **149**, 3910–3916.
- 897 Zwicker, E. (1961). “Subdivision of the audible frequency range into critical bands  
898 (Frequenzgruppen),” *J. Acoust. Soc. Am.* **33**, 248.