

More on the articulation of devoiced [u] in Tokyo Japanese: effects of surrounding consonants

Jason A. Shaw^a & Shigeto Kawahara^b

^aDepartment of Linguistics, Yale University,
New Haven, CT 06520, USA

^bThe Institute of Cultural and Linguistic Studies, Keio University,
Minato-ku, Tokyo 108-8345, Japan

Running head: Articulation of devoiced [u] in Japanese

Corresponding author: jason.shaw@yale.edu +1 203-432-8289 (J.A. Shaw).

Keywords: Japanese, devoicing, EMA, articulation, deletion, phonetic interpolation, gestural coordination, syllable contact

Abstract

Several aspects of high vowel devoicing in Tokyo Japanese have been extensively studied. One aspect of the phenomenon that remains understudied is the lingual articulation of devoiced vowels, including whether devoiced vowels retain their lingual gesture. [Shaw & Kawahara \(2018b\)](#) addressed this question using EMA (Electro-Magnetic Articulography), finding optional but categorical deletion patterns: some vowels retained a full lingual target, just like their voiced counterparts, whereas other vowels showed trajectories that are best modelled as targetless, i.e., linear interpolation between the surrounding vowels. Extending this finding, as well as being inspired by various phonetic and phonological considerations, the current study explores the hypothesis that this probabilistic deletion of devoiced high vowels may be modulated by the identity of the surrounding consonants. A new follow-up EMA-based experiment with an extended stimulus set replicates the core finding of [Shaw & Kawahara \(2018b\)](#) that Japanese devoiced [u] sometimes lacks a tongue body raising gesture. The current results moreover show that surrounding consonants do indeed affect the probability of tongue dorsum targetlessness. We found that deletion of devoiced vowels is affected by the place of articulation of the preceding consonant, with deletion more likely following a coronal fricative than a labial fricative. Additionally, we found that the manner combination of the flanking consonants, fricative-fricative vs. fricative-stop also has an effect, at least for some speakers; however, unlike the effect of C1 place, the direction of the manner combination effect varies across speakers with some deleting more often in fricative-stop environments and others more often in fricative-fricative environments.

1 Introduction

1.1 General background

Vowels that are adjacent to voiceless obstruents are sometimes produced without vocal fold vibration—i.e. as voiceless—the phenomenon generally referred to as “vowel devoicing.” This pattern is observed systematically across many genetically-unrelated languages, including but not limited to Cheyenne ([Rhodes, 1972](#); [Vogel 2021](#)), French ([Smith, 2003](#)),

30 Greek (Dauer 1980), Korean (Jun et al. 1998), (Andean) Spanish (Delforge 2008), Us-
31 panteko (Bennett 2020), Uzbek (Sjoberg 1963), and Turkish (Jannedy 1995). Tokyo
32 Japanese also exhibits such devoicing of high vowels, and Japanese is arguably the best
33 studied language in this respect (Fujimoto 2015 for a recent review).

34 A general characterization of the high vowel devoicing pattern in Tokyo Japanese is
35 that high vowels are devoiced between two voiceless obstruents, as well as after a voice-
36 less obstruent and before a pause, although as we will review below, the likelihood of
37 devoicing is affected by various other factors, both linguistic and social. Previous studies
38 have explored this devoicing process from various perspectives, each bearing upon some
39 important issues in phonetic and phonological theory, including how devoicing is imple-
40 mented in terms of the laryngeal gesture (Fujimoto et al. 2002; Sawashima 1971; Yosh-
41 ioka 1981), how the precise environment affects the likelihood of devoicing (Maekawa &
42 Kikuchi 2005; Tsuchida 1997), its categorical or gradient nature (Nielsen 2015; Tanner
43 et al. 2019), its interaction with lexical pitch accent (Kuriyagawa & Sawashima 1989;
44 Maekawa 1990; Vance 1987) and other prosodic properties (Kilbourn-Ceron & Son-
45 deregger 2018), its consequences (or lack thereof) for prosodic reorganization (Kondo
46 1997, 2001; Kawahara & Shaw 2018), its perceptual consequences (Cutler et al. 2009;
47 Ogasawara 2013; Sugito & Hirose 1988; Whang 2019), its role in child-directed speech
48 (Fais et al. 2010; Martin et al. 2014), its acquisition patterns (Imaizumi & Hayashi 1995;
49 Imaizumi et al. 1995) as well as the influence of social factors on this pattern (Imai 2004;
50 Imaizumi et al. 1995). There is no doubt that these studies have revealed important as-
51 pects of this devoicing phenomenon in Japanese, and we understand its nature much better
52 than 50 years ago.

53 However, despite the accumulation of studies on vowel devoicing, one aspect that is
54 heavily under-addressed regarding the pattern of high vowel devoicing in Japanese—and
55 any other languages that exhibit vowel devoicing, for that matter—is the question of how
56 the lingual gesture is implemented for the devoiced vowels. This issue is related to the
57 question of whether these devoiced vowels are phonologically deleted or not; if the high
58 vowels are phonologically deleted, then we would expect them to lack any lingual gesture.
59 If the process at issue is phonologically devoicing rather than deletion, on the other hand,
60 we may expect that their lingual gestures are retained. Vance (2008), which is the most

61 recent and updated textbook on Japanese phonetics and phonology, indicates that this issue
62 is not settled. The current paper reports an experiment that addresses this issue by analyz-
63 ing articulatory kinematics via Electromagnetic Articulography. This paper moreover tests
64 a new, specific hypothesis that deletion probability may be modulated by the surrounding
65 consonantal environment. We will start with the overview of the relevant literature on this
66 topic, which leads us to examine this specific hypothesis.

67 1.2 Are devoiced vowels in Japanese deleted?

68 Since devoiced vowels lack a periodic energy source, it is difficult, if not entirely impossi-
69 ble, to infer from their acoustic profiles alone whether devoiced vowels retain their lingual
70 gestures or not. There are some studies which addressed this question via impressionistic
71 observations. Kawakami (1977) argues that vowels delete completely in some environ-
72 ments but not others, but he offers no experimental evidence for this claim. Vance (1987)
73 examines the hypothesis that devoiced high vowels in Japanese are entirely deleted but
74 ultimately rejects this hypothesis. Kondo (1997, 2001) argues that devoiced high vowels
75 are deleted based on a phonological consideration: vowel devoicing in consecutive sylla-
76 bles is often inhibited (though see Nielsen 2015), and Kondo (1997, 2001) attributes this
77 observation to a constraint against triconsonantal clusters. The underlying assumption of
78 this analysis is that devoiced vowels are deleted, resulting in consonant clusters.

79 On the other hand, Tsuchida (1997) and Kawahara (2015) point out that devoiced
80 vowels count toward a bimoraic requirement in foot-based morphophonological truncation
81 patterns (Poser, 1990), arguing that these vowels do not delete phonologically. Like these
82 two authors, Hirayama (2009) demonstrates that devoiced vowels behave just like voiced
83 vowels in the Japanese *haiku* poetry pattern, which is mora-based (Vance, 1987).

84 In line with this view, Jun and her colleagues advanced an explanation of high vowel
85 devoicing (in Korean) in terms of “gestural overlap” (Browman & Goldstein 1992a), ac-
86 cording to which the articulatory gesture of high vowels is overlapped in time by the
87 laryngeal glottal abduction gesture of surrounding consonants (Jun & Beckman, 1993; Jun
88 et al., 1998). In this gestural overlap view, Japanese phonology does not delete the de-
89 voiced high vowels; the high vowels are merely rendered inaudible because of the glottal

90 abduction gesture that coincides in time with the vocalic gesture. This is analogous to the
91 famous case of English *perfect memory*, in which the word-final [t] in *perfect* can be made
92 inaudible due to gestural overlap with the following [m], even when the [t]’s tongue tip
93 gesture remains intact (Browman & Goldstein, 1992a). Similar to this case in English, it
94 is conceivable that lingual gestures of devoiced high vowels are present, but are merely
95 rendered inaudible because of the overlapping glottal abduction gesture. In this gestural-
96 overlap scenario, it is also possible that lingual gestures are reduced, rather than remaining
97 completely intact, assuming speakers invest less articulatory energy into sounds that are
98 difficult to perceive and hence may not contribute much to lexical retrieval (e.g. Hall et al.
99 2018; Jaeger & Buz 2018).

100 Recently we have witnessed a rise of studies addressing this question—whether de-
101 voiced high vowels are deleted or not—using instrumental techniques. Beckman and her
102 colleagues, based on the inspection of spectrograms, argue that devoiced vowels are phys-
103 ically not present (Beckman, 1982; Beckman & Shoji, 1984), suggesting that the pattern
104 should be characterized as deletion, although they also suggest that it may make sense to
105 characterize the pattern as devoicing, not deletion, from the psycholinguistic perspective;
106 i.e. Japanese speakers feel that “vowels are there” even when they are actually deleted
107 (cf. Dupoux et al. 1999, 2011; Whang 2019). This is possibly because of coarticulatory
108 influences of vowels on flanking consonants that remain even when typical acoustic cues
109 to the vowel are absent. It is known, for example, that consonant identity influences vowel
110 quality in perceptual epenthesis (Durvasula et al. 2018; Kilpatrick et al., 2020).

111 Faber & Vance (2010) offer some acoustic evidence for the hypothesis that vowel de-
112 voicing is best characterized as gestural overlap of laryngeal gestures in Japanese (Jun &
113 Beckman, 1993; Jun et al. 1998). Jannedy (1995) and Bennett (2020) entertain a similar
114 hypothesis for devoiced vowels in Turkish and Uspanteko, respectively. Whang (2018)
115 measured COG during devoiced vowels in Japanese and argues that some devoiced vow-
116 els in Japanese are in fact deleted, while others are not. More specifically, Whang (2018)
117 argues that deletion is more likely in the environment where the quality of those vowels
118 can be recovered from surrounding consonants; e.g. after [F], only [u]¹ is possible, while

¹Here and throughout the paper, we use the IPA symbol [u] to denote the high non-front vowel in Japanese. The exact phonetic nature of this vowel, as well as how to transcribe it, is a contentious issue

119 after [s] both [u] and [i] are possible (see also Whang 2019).

120 However, generally speaking, there are limits on how much we can conclude about the
121 articulatory gestures from their resulting acoustic signals (see e.g. Browman & Goldstein
122 1989; Guenther et al. 1999; Munson et al. 2010; Perkell et al. 1993). It is thus impor-
123 tant that we address the nature of the lingual gesture of devoiced high vowels through
124 observation of articulatory movement. To that end, Matsui (2017) used EPG (Electro
125 PalatoGraphy) to examine the linguo-palatal contact pattern of devoiced syllable [s]y, and
126 showed that the pattern remains very constant across the syllable; i.e. there does not seem
127 to be a clear change in the linguo-palatal contact pattern from [s] to [y], implying that
128 the lingual gesture of the devoiced [y] is absent. Nakamura (2003) on the other hand re-
129 ports that vestiges of lingual gestures of devoiced vowels can be found in his EPG data.
130 Although these two results, which seem to conflict with each other, are telling, there are
131 limits on how much we can conclude about tongue body movement—primary correlates
132 of vowel gestures (Browman & Goldstein 1992b; Johnson et al. 1993)—from EPG data
133 in general, since EPG only registers contact with the palate. Funatsu & Fujimoto (2011)
134 used ElectroMagnetic Midsagittal Articulography (EMMA) to study articulatory gestures
135 of devoiced [i], showing that the articulatory gesture of [i] is comparable between voiced
136 [i] and devoiced [i]. This study however used one speaker and one pair of items (/kide/
137 vs. /kite/) with four repetitions, and offers no quantitative comparisons between the two
138 voicing conditions.

139 The most extensive study on this topic—the presence/absence of lingual gestures of
140 devoiced vowels in Japanese—to date is that of Shaw & Kawahara (2018b), who used
141 EMA (ElectroMagnetic Articulography) to study the articulatory nature of devoiced [y]s
142 of six naive speakers of Tokyo Japanese, and the current paper can be considered as a
143 direct follow-up of Shaw & Kawahara (2018b).

144 Shaw & Kawahara (2018b) analyzed four dyads to compare the articulatory trajectories
145 of CVC sequences, in which one member of each dyad contains a voiced vowel and the
146 other a devoiced vowel. The four dyads were: (1) [F_syoku] vs. [Fuzoku], (2) [s_syutaisee] vs.

even in the contemporary literature (Vance 2008). We will return to this issue in the method section (3.2),
where we justify our choice of phonetic parameters used to assess the deletion of this vowel.

147 [sudaika], (3) [katsutoki] vs. [katsudoo] and (4) [masutaa] vs. [masuda]² Their strategy,
148 reviewed in further detail below in §3.2 is to compare the articulatory trajectory of [CuC]
149 with respect to that of [CuC] and [C;C], the latter of which is characterized by linear
150 interpolation between the surrounding vowels (Choi, 1995; Cohn, 1993; Keating, 1988;
151 Pierrehumbert & Beckman, 1988). Their conclusion in a nutshell is that some productions
152 contain no articulatory target, while others show lingual targets that are no different from
153 voiced vowels; i.e. they found a pattern of optional but categorical deletion. Moreover,
154 they found some variation with respect to how often each item showed devoiced vowels
155 without lingual targets: devoiced vowels were more likely to be targetless between [ʃ] and
156 [t] ([sutaisee]) than between [f] and [s] ([fusoku]). This asymmetry was consistent across
157 the speakers (see also discussion in §3.2).

158 An intriguing hypothesis that emerges from this result is that vowel deletion probability
159 may be systematically modulated via surrounding consonant environments—Japanese [y]s
160 may be more likely to lack a lingual gesture between a fricative and a stop than between
161 two fricatives. We expand in the next subsection why this is an interesting and plausible
162 hypothesis to entertain, although we also note at this point that the results by Shaw &
163 Kawahara (2018b) are based on just one dyad per each phonological condition.

164 1.3 The current hypothesis

165 The general hypothesis pursued in this study is that the probability of [y] lacking its lingual
166 gesture—which we equate with the probability of phonological deletion for the sake of
167 exposition here (see Shaw & Kawahara 2018b)—is modulated by surrounding consonantal
168 environment. A more specific hypothesis is that [y]s are more likely to be phonologically
169 deleted when surrounded by a fricative and a stop than when surrounded by two fricatives.
170 As mentioned above, one reason to entertain this hypothesis is the results reported by Shaw
171 & Kawahara (2018b), who show that targetless [y]s were more likely in [sutaisee] than in
172 [fusoku]. However, it is hard to know whether or not their findings are generalizable to
173 other items with similar phonological properties, because their results are based on one

²Glosses: (1) shortage vs. attachment, (2) subjectivity vs. theme song, (3) when winning vs. activity and (4) master vs. PERSONAL NAME.

174 dyad per each phonological condition.

175 Nevertheless, this hypothesis dovetails with an observation by Starr & Shih (2017),
176 who found that devoiced vowels are often skipped in the text-setting of Japanese songs,
177 and this is especially so when they are surrounded by a fricative and a stop. Their observa-
178 tion may suggest that Japanese composers are sensitive to the higher likelihood of vowel
179 deletion in this environment. The higher likelihood of deletion after a fricative and be-
180 fore a stop is also compatible with the general cross-linguistic observation about prosodic
181 wellformedness, namely, syllable contact laws (Murray & Vennemann, 1983; Vennemann
182 1988)—languages generally prefer sonority fall to sonority plateau/rise across a syllable
183 boundary. To the extent that Japanese is also sensitive to such prosodic wellformed-
184 ness conditions (cf. Berent et al. 2007, 2008), we may expect Japanese high vowels to
185 delete more often in the environment which would result in a fricative-stop cluster than
186 a fricative-fricative cluster. To view it from the opposite perspective, if it can be shown
187 that Japanese speakers delete high vowels in accordance with syllable contact law, it may
188 imply that speakers of Japanese, generally considered to be a “CV-language” disallowing
189 hetero-organic consonant clusters, are sensitive to wellformedness conditions on conso-
190 nant clusters (see Berent et al. 2007, 2008 for related ideas), possibly because they can
191 extrapolate sonority-based patterns from limited data (Daland et al., 2011).

192 There are other reasons to entertain the current hypothesis. Previous studies have
193 shown that devoicing (not necessarily deletion) is more likely between a fricative and a
194 stop than between two fricatives (see e.g. Fujimoto 2015; Hirayama 2009; Maekawa &
195 Kikuchi 2005; Martin et al. 2014; Tsuchida 1997). Let us suppose that devoicing and
196 deletion are on the same “reduction continuum” ³. Then, everything else being equal, we
197 may expect deletion to be more likely in the environment where devoicing is more likely
198 in the first place. This leads us to expect that deletion is more likely between a fricative
199 and a stop, because devoicing is more likely in this environment.

200 A recent acoustic study by Whang et al. (2020) suggests that devoicing and deletion

³In Kagoshima Japanese, word-final high vowels—those that are devoiced—undergo phonological dele-
tion, which feeds other phonological changes of consonants in the word-final syllables (Haraguchi 1984;
Kaneke & Kawahara 2002; Kibe, 2001). It seems natural to consider deletion as the extreme end of the
reduction continuum, and that devoicing is one step in the continuum before deletion (see Haraguchi 1984;
McCarthy 2008; Tsuchida 1997).

201 may both be characterized as enhancement strategies of the larynx abduction gesture. [Fu-
202 jimoto et al. (2002) as well as Sawashima (1971) show that devoiced vowels in Japanese
203 involve an active abduction gesture, and thus there is a sense in which speakers are ac-
204 tively signalling “voicelessness.” According to Whang et al. (2020), vowel devoicing in
205 fact raises COG of the aperiodic energy of surrounding obstruents, possibly due to wider
206 glottal aperture and increased airflow, and deletion of the tongue dorsum raising gesture
207 for [u] further raises that COG. This hypothesis too leads us to expect that devoicing and
208 deletion should work in tandem with each other, as deletion can enhance the auditory cue
209 to devoicing. To the degree that devoicing is more likely after a fricative and before a stop
210 than between two fricatives (see above), deletion may show the same probabilistic pattern.

211 All of these considerations—prosodic wellformedness, reduction, enhancement of devoicing—
212 converge on the same prediction: deletion should be more likely when it results in a
213 fricative-stop sequence than when it results in a fricative-fricative sequence. Existing ev-
214 idence from Shaw & Kawahara (2018b) is consistent with this conclusion; however, the
215 evidence to date is rather thin.

216 To recap, the current experiment was set out to examine the general hypothesis that
217 vowel deletion probability is modulated by surrounding consonants. The more specific
218 hypothesis is that deletion is most likely between a fricative and a stop, and less likely be-
219 tween two fricatives. The experiment also serves as an attempt to replicated the basic find-
220 ings of Shaw & Kawahara (2018b)—devoiced [u]s in Japanese are optionally deleted—
221 with a much extended set of stimuli.

222 2 Experimental methods

223 The current experiment measured and analyzed the tongue dorsum trajectories of devoiced
224 [u], using EMA (Electro Magnetic Articulatograph). Most of the experimental details fol-
225 low those of Shaw & Kawahara (2018b). One distinct characteristic of this approach that
226 we would like to highlight at this stage is that it assesses the presence of an articulatory
227 target on a token-by-token basis, rather than analyzing averaged contours. This strategy is
228 important because analyzing averaged contours cannot distinguish two different phonolog-
229 ical hypotheses: reduction vs. optional deletion (Cohn, 2006; Kawahara et al. to appear;

230 Shaw & Kawahara (2018a). Lingual gestures of devoiced vowels, even when not phono-
231 logically deleted, are conceivably reduced in magnitude, since the vowel gestures are not
232 as audible due to devoicing and do not contribute much to lexical access anyway (e.g. Hall
233 et al. 2018; Jaeger & Buz 2018; see also Lindblom 1990). Interpreting any difference
234 between voiced vowels and devoiced vowels as deletion would therefore be hasty.

235 On the other hand, as Shaw & Kawahara (2018b) found, devoiced vowels can retain
236 their full lingual gestures, showing comparable movement trajectories to voiced vowels,
237 but they can also sometimes be deleted. Averaging over cases of full targets and cases of
238 categorical deletion can lead to an erroneous conclusion that the overall pattern supports
239 the reduction hypothesis (Cohn, 2006).

240 This specific problem can be illustrated by a comparison of two recent studies. Kawa-
241 hara et al. (to appear) developed a token-by-token analysis of the F0 patterns of the dataset
242 recorded and analyzed by Ishihara (2011). The averaged-based analysis by the latter con-
243 cluded that pitch accent after wh-elements in Japanese is reduced. On the other hand, a
244 token-by-token reanalysis by Kawahara et al. (to appear) shows that at least some speakers
245 show a mixture of full target and deletion. This comparison shows that when both deletion
246 and reduction are theoretically-justifiable hypotheses, it is important that we distinguish
247 between them through a token-by-token analysis.

248 In addition to avoiding this general problem of resorting to average-based analyses,
249 the current analysis has a virtue of analyzing the entire articulatory trajectories; in the
250 current analysis, no aspects of speech signals within the analysis window are given special
251 status, eschewing the potential danger of missing important aspects of dynamic speech
252 (Cho 2016; Mücke et al. 2014; Vatikiotis-Bateson et al. 2014).

253 2.1 Participants

254 Seven native speakers of Tokyo Japanese (4 male) participated in the current experiment.
255 They were all born in Tokyo, lived there at the time of their participation in the study, and
256 had spent no more than 3 months outside of the Tokyo region. Procedures were explained
257 to participants in Japanese by a research assistant, who was also a native speaker of Tokyo
258 Japanese. All participants were naive to the purpose of the experiment. Participants were

259 compensated for their time and local travel expenses. Data from one speaker had to be
260 excluded, because we were unable to record as many repetitions as other speakers. This
261 speaker was originally coded as Speaker 6; their data is not discussed further below. No
262 speakers who participated in Shaw & Kawahara (2018b) participated in this study, since
263 one of the aims was to examine whether the results of Shaw & Kawahara (2018b) can be
264 generalized to other speakers.

265 2.2 Stimuli

266 Following Shaw & Kawahara (2018b), the major target of our analysis is tongue dorsum
267 height in the trajectory of $V_1C_1V_2C_2V_3$ sequence, in which V_2 represents the devoiced
268 vowels in question—justification of this analytical choice is offered below in §3.2. The
269 primary question is whether we would observe a clear rise in tongue dorsum height from
270 V_1 to V_2 and a fall from V_2 to V_3 . V_3 was therefore always a non-high vowel in our
271 stimuli. The target vowels ($V_2=[u]s$) were always word-initial, and V_1 was the last vowel
272 of the preceding word in the frame sentence, [e].

273 At the time of stimulus design, four conditions were included in order to thoroughly
274 explore the effects of surrounding consonant types: fricative-stop (FS), fricative-fricative
275 (FF), stop-stop (SS), and stop-fricative (SF), consisting of 18 dyads shown in Table 1. All
276 the stimuli were existing words in Japanese, where the members on the left were expected
277 to undergo devoicing. Each dyad constituted near minimal pairs, in which one member
278 contained C_1VC_2 sequence where both consonants are voiceless and the other member
279 contained a minimally different C_1VC_2 sequence in which C_2 is voiced, hence V is not
280 expected to devoice.

Table 1: The list of stimuli recorded in the EMA experiment. S=Stop; F=Fricative. See footnote 5 for glosses. Accent is shown by a following apostrophe.

FS	FF
/Futon/ vs. /Fudou/	/Fusoku/ vs. /Fuzoku/
/Futan/ vs. /Fu'dan/	/Fusai/ vs. /Fuzai/
/Futa/ vs. /Fuda/	/Fusagaru/ vs. /Fuzake'ru/
/Sutaisei/ vs. /Suda'ika/	/Susai/ vs. /Suzai/
/Sutou/ vs. /Sudou/	/Su'sa/ vs. /Su'zan/
/Sutokou/ vs. /Sudo'uken/	/Su'so/ vs. /Suzou/
SS	SF
/kutakuta/ vs. /kudaranu/	/kusami/ vs. /kuzai/
/kutaba'ru/ vs. /kudasa'ru/	/kusari/ vs. /kuzawa/
/kutanijaki/ vs. /kuda'nſita/	/kusaka'ri/ vs. /kuzakitſo/

281 Choosing existing words with the appropriate segmental compositions did not en-
 282 able us to control for accentedness within each pair. For example, /Futan/ is unaccented,
 283 whereas /Fu'dan/ is accented on the initial syllable. However, Tsuchida (1997) and Martin
 284 et al. (2014) show that accent placement has little effect on devoicing patterns among con-
 285 temporary speakers of Japanese. Durational differences between accented and unaccented
 286 syllables are minimal in Japanese (Beckman, 1986), which if substantial, may affect de-
 287 voicability/deletability. For these reasons, we judged this difference to be non-crucial⁴

288 The current study focused on [u] instead of examining both [u] and [i], both of which
 289 are known to devoice. This is partly because the current study is a direct follow-up of Shaw
 290 & Kawahara (2018b), who also examined only [u], and also because we needed enough
 291 repetitions to execute the computational analysis that was planned (see 3.2 for details).
 292 Examining the lingual gesture of devoiced [i] warrants a new set of studies.

293 After the recording, we came to the conclusion that the conditions in which C₁ is a stop

⁴Shaw & Kawahara (2018b) did not perfectly control for accent between two members within a dyad ei-
 ther, although in their design, [u] is either accented or unaccented within each dyad, i.e. no direct comparison
 was made between accented [u] and unaccented [u].

294 (=the SS and SF conditions) could not be reliably analyzed for the following reason. At
295 the time of stimulus design, we decided that C₁ had to be [k], because [p] is not allowed in
296 the native vocabulary (Ito & Mester 1995), and [t] is affricated before high vowels (Vance
297 2008). However, since we were interested in the tongue dorsum height of (devoiced) [u],
298 it was not possible to objectively discern control of tongue dorsum height associated with
299 [k] from tongue dorsum height associated with [u]. For this reasons, this paper focuses on
300 the comparison between FS condition and the FF condition.⁵

301 2.3 Procedure

302 Each participant produced 14-15 repetitions of the 36 target words in the carrier phrase:
303 “okkee X to itte” (Ok, say X), where X is a stimulus word. Participants were instructed to
304 speak as if they were making a request of a friend. This was to ensure that the speakers did
305 not speak too formally or too slowly, which may inhibit vowel devoicing in the first place.

306 This resulted in a corpus of 3,204 tokens (14 or 15 repetitions ÷ 36 words ÷ 6 speak-
307 ers). Words were presented in Japanese script (composed of hiragana, katakana and kanji
308 characters as required for natural presentation) and fully randomized.

309 2.4 Equipment

310 We used an NDI Wave ElectroMagnetic Articulograph system sampling at 100 Hz to cap-
311 ture articulatory movement. NDI wave 5DoF sensors (receiver coils) were attached to
312 three locations on the sagittal midline of the tongue, and on the lips, jaw (below the lower
313 incisor), nasion and left/right mastoids. The most anterior sensor on the tongue, henceforth
314 TT, was attached less than one cm from the tongue tip (see Figure 1). The most posterior
315 sensor, henceforth TD, was attached as far back as was comfortable for the participant. A
316 third sensor, henceforth TB, was placed on the tongue body roughly equidistant between
317 the TT and TD sensors. Sensors were attached with attached with a combination of sur-

⁵The glosses for the items that were analyzed are as follows. FF: blanket vs. not moving, burden vs. usual, top vs. amulet, subjectivity vs. main theme, FOOD NAME vs. hand-moving, Tokyo Highway vs. lead; FS: shortage vs. attachment, debt vs. absence, filled vs. joke, organize vs. research, chair vs. abacus, main claim vs. sake-making.

318 gical glue and ketac dental adhesive. Acoustic data were recorded simultaneously at 22
319 KHz with a Schoeps MK 41S supercardioid microphone (with Schoeps CMC 6 Ug power
320 module).

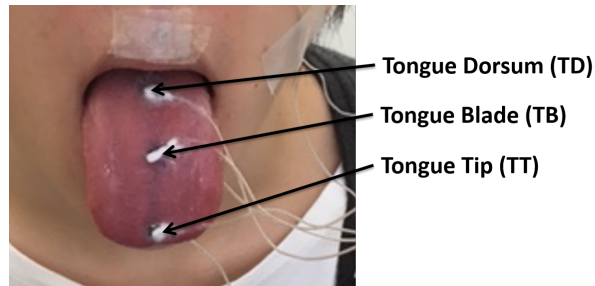


Figure 1: Illustration of the sensor placement (reproduced from [Shaw & Kawahara 2018b](#)).

321 **2.5 Stimulus display**

322 Words were displayed on a monitor positioned 25cm outside of the NDI Wave magnetic
323 field. Stimulus display was controlled manually using an Eprime script. This setup allowed
324 for online monitoring of hesitations, mispronunciations and disfluencies. These were rare,
325 but when they occurred, items were marked for repeated presentation by the experimenter.
326 These items were then re-inserted into the random presentation of remaining items. This
327 method ensured that we recorded at least 14 fluent tokens of each target item.

328 **2.6 Post-processing**

329 Following the main recording session, we also recorded the bite plane of each participant
330 by having them hold a rigid object, with three 5DoF sensors attached to it, between their
331 teeth. Head movements were corrected computationally after data collection with refer-
332 ence to three sensors on the head, the left/right mastoid and nasion sensors, and the three
333 sensors on the bite plane. The head corrected data was rotated so that the origin of the
334 spatial coordinates corresponds to the occlusal plane at the front teeth.

335 **3 Data analysis**

336 **3.1 Data processing**

337 The wav files recorded in the experiment were submitted to forced alignment, using FAVE⁶
338 Textgrids from forced alignment were hand-corrected and, during this process, the target
339 vowels were coded for voicing. Most vowels in devoicing environments were in fact de-
340 voiced, as evident from visual inspection of the spectrogram and waveform. However,
341 some tokens in the devoicing environment exceptionally retained clear signs of glottal
342 vibration. These vowels were coded as voiced, and excluded from the following computa-
343 tional analysis. The supplementary materials, available at DOI 10.17605/OSF.IO/PGRVZ,
344 provide example spectrograms of voiced and devoiced tokens and a list of all exclusions.

345 Articulatory data corresponding to each token were extracted based on the textgrids.
346 The data were smoothed using the robust smoothing algorithm (Garcia, 2010) and, sub-
347 sequently, visualized in MVIEW, a Matlab-based program to analyze articulatory data
348 (Tiede 2005). Within MVIEW, the consonant gestures flanking the target vowel were
349 parsed using the `findgest` algorithm. `findgest` identifies gestures semi-automatically
350 based upon the velocity signal in the movement toward and away from gestural targets. An
351 illustrative example is provided in Figure 2. The consonant gestures were used to define
352 a temporal interval for further analysis.⁷ Tokens with missing data in the target interval
353 were excluded from further analysis. Some tokens had velocity peaks that were not large
354 enough to clearly parse out movement related to the consonants. If a token was missing a
355 gesture parse for either consonant, it was excluded from further analysis. A total of 239
356 tokens were excluded for this reason. The resulting data set consisted of 2,431 tokens for
357 analysis, which had clearly distinguishable consonantal gestures flanking the target vowel.

⁶<https://github.com/JoFrhwld/FAVE/wiki/Using-FAVE-align>

⁷The onset of movement of the consonants occurs at a similar time as the maximum tongue height of the preceding vowel. We choose to define the temporal interval for analysis based on the onset of consonant movement instead of, e.g., the maximum TD height in the vicinity of the consonant, primarily because the results presented here are situated in a bigger project which includes also how the reduction/deletion of vowels influences the coordination of flanking consonants.

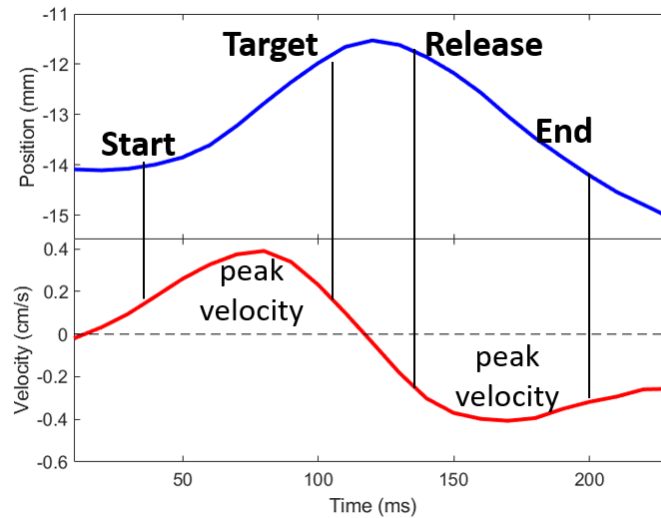


Figure 2: A sample articulatory trajectory and how the articulatory landmarks were identified using `findgest`.

3.2 Computational analyses

The temporal interval spanning from the onset of movement of C_1 , the consonant preceding the target vowel, and the offset of movement of C_2 , the consonant following the target vowel, was subjected to further analysis. To address the question of whether devoiced [u] has an articulatory target, we focused on tongue height, instead of tongue retraction or lip gestures, both of which have been questioned as reliable articulatory correlates of this vowel in contemporary Japanese (Isomura, 2009; Nogita et al., 2013; Shaw & Kawahara 2018a; Vance 2008). Like Shaw & Kawahara (2018b), the analysis focused on the movements of the TD sensor (see Figure 1), the most posterior sensor on the tongue, which is typically used to detect vowel gestures (Browman & Goldstein, 1992b; Johnson et al. 1993).

Figure 3 shows sample trajectories of a voiced vowel (left), a devoiced vowel with a clear tongue dorsum raising during [u] (middle), and a devoiced vowel without a very clear movement in terms of tongue dorsum height (right). The top panels show the audio

372 signal. The second panels from the top show tongue dorsum articulatory trajectories,
 373 which are the primary target of our analyses. For reference the third and fourth panels
 374 show trajectories related to the flanking consonants. The token in the right panel does not
 375 appear to have a clear tongue dorsum raising gesture during [y], whereas the [y]
 376 in the middle panel does seem to have a clear raising gesture. The challenge is to go
 377 beyond such impressionistic classifications and to establish an objective method to classify
 378 whether devoiced vowels show a tongue dorsum raising gesture or not.

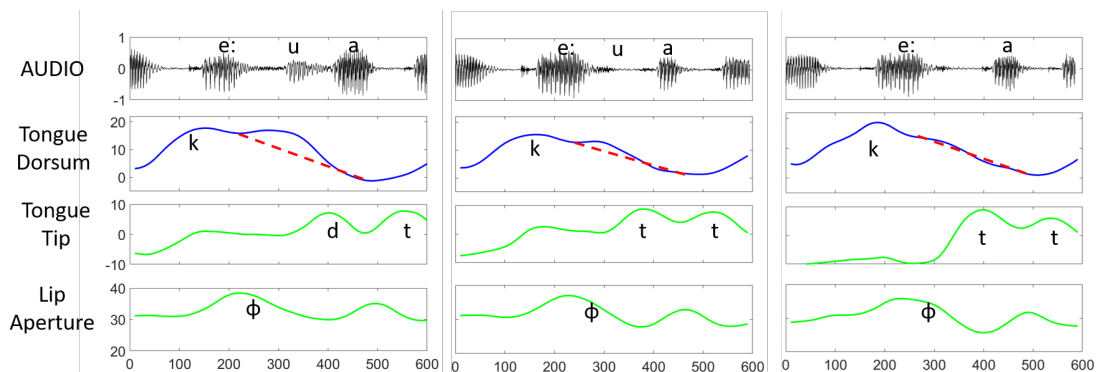


Figure 3: Sample EMA trajectories. The top panels show audio signals. The second panels show the tongue dorsum movement. The dotted red line is a linear interpolation from the preceding vowel to the following vowel.

379 To do so, we applied the approach described and motivated in detail in [Shaw & Kawahara \(2018a, b\)](#), schematically illustrated in Figure 4. This computational methodology
 380 was developed to assess the presence/absence of a lingual vowel target of devoiced vowels
 381 in articulatory trajectories. The approach is general enough that it has been extended to
 382 other types of continuous phonetic data, including nasal reduction in Ende ([Brickhouse &
 383 Lindsey, 2020](#)), pitch accent eradication in Japanese ([Kawahara et al., to appear](#)), and tone
 384 reduction in Mandarin Chinese ([Zhang et al., 2019](#)).

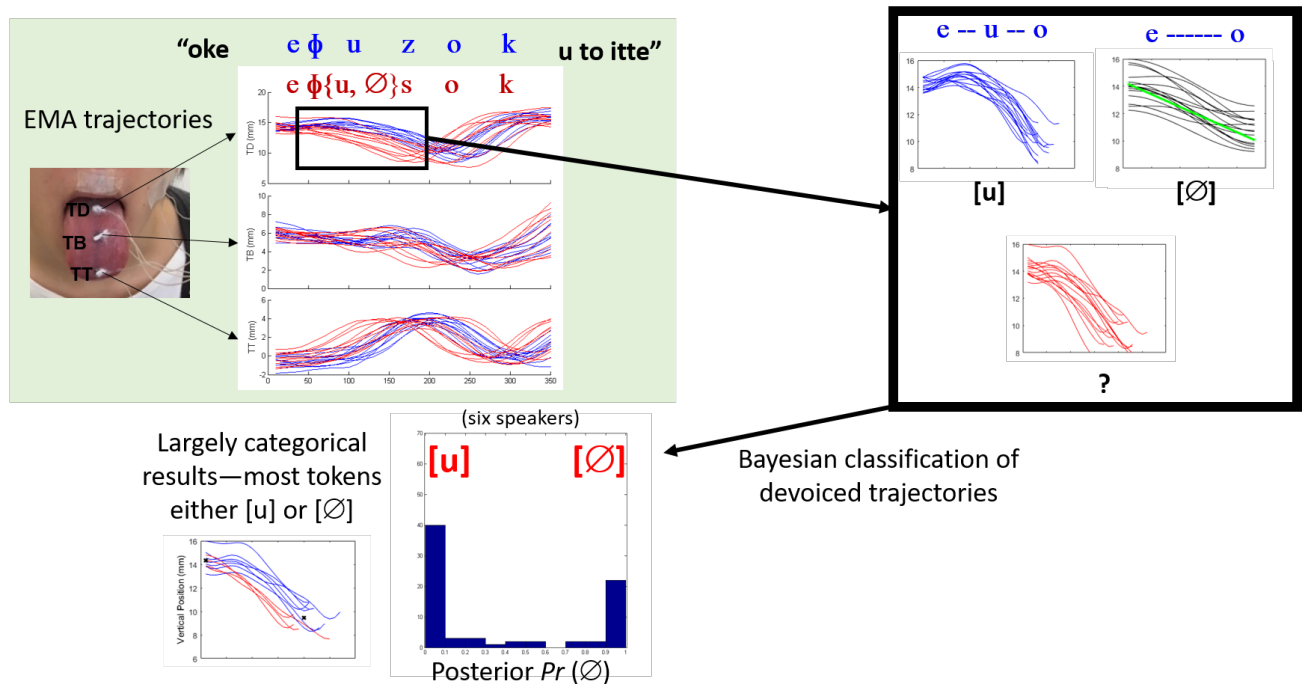


Figure 4: Summary of simulation and classification procedure developed and defended in [Shaw & Kawahara \(2018b\)](#).

386 The target interval spans from the preceding vowel to the following vowel (see the
 387 left upper panel of Figure 4). For example, for the word [Fuzoku], the analysis window
 388 starts from [e] in the carrier sentence and includes the main target CVC ([Fuz]) and the
 389 following vowel [o]. The question of interest is whether given the vowel sequence [e]-[u]-
 390 [o], we would observe a tongue dorsum raising gesture, when [u] is devoiced. When [u]'s
 391 tongue dorsum gesture is undoubtedly present, as in the case for voiced [u], we should
 392 observe a clear raising gesture (the left panel of Figure 3). On the other hand, if the vowel
 393 gesture is deleted, we expect articulatory trajectories that interpolate between [e] and [o]
 394 (represented as a green straight line in the right upper panel of Figure 4). Since articulatory
 395 movements, as behavioral data more generally, are always noisy actuations of intentions,
 396 the challenge is to develop an objective method with which we can assess whether each

397 articulatory contour of a devoiced [y] is better characterized as target-present or target-
398 absent (the upper right box in Figure 4). The computational toolkit developed by Shaw &
399 Kawahara (2018a,b) allows us to address this question on a token-by-token basis.

400 The first step in this computational method is to analyze the articulatory trajectories
401 in a low-dimensional space, by making use of Discrete Cosine Transform (DCT) (e.g.
402 Jain 1989). Through DCT, a signal is transformed into the sum of cosine components of
403 gradually increasing frequency. This transformation is similar to Fourier transform in that
404 timeseries data—here, the articulatory trajectory—is represented in frequency space, i.e.,
405 as cosines of varying frequency and magnitude. Unlike Fourier transform, DCT uses only
406 cosines instead of a combination of sines and cosines and there is no imaginary compo-
407 nent. Additionally, DCT has compression properties (Jain, 1989), like Principal Compo-
408 nent Analysis (PCA)—the articulatory trajectory within the analysis window can often be
409 represented with a small number of DCT components. Because speech articulators are
410 relatively slow, high frequency components are not needed to represent their controlled
411 movement, a point which we demonstrate below.

412 The numerical expression of DCT is provided in Equations (1) and (2): n is the po-
413 sitional signal, L is the length of the window (in samples), k is the number of the DCT
414 coefficient, which ranges from 1 to L , y is the magnitude of each coefficient, and w is a
415 weight. DCT coefficients can be positive or negative and their absolute value represents
416 the magnitude of their contribution to spatial modulation of the signal. For the first DCT
417 coefficient, the numerator in the scope of the cosine is zero, which means that it equals
418 1 for every sample n in the trajectory. These are summed, and when multiplied by the
419 relevant weight ($\frac{1}{L}$), they yield a quantity that is related to the average of the trajectory
420 (if the weight was $\frac{1}{L}$, then it would be the average). This first cosine coefficient serves as a
421 baseline, c.f. the intercept in a linear regression. As k increases beyond one, the resulting
422 cosines gradually increase in frequency, $k = 2$ yields a cosine that completes one quarter
423 of its cycle within the signal, $k = 3$, yields a half cycle and so on (see Figure 6). DCT
424 produces $k = L$ components, so the number of cosine components depends on the length
425 of the signal. However, the magnitude of the higher frequency components may be quite
426 small for signals of slow moving articulators.

$$y(k) = w(k) \sum_{n=1}^N \cos \frac{3(2n-1)(k-1)}{2L} \quad k = 1; 2; \dots; L \quad (1)$$

where

$$w(k) = \begin{cases} 1 & k = 1 \\ r \frac{2}{L} & 2 \leq k \leq L \end{cases} \quad (2)$$

427 DCT has a known inverse function, iDCT, which can be used to simulate trajectories
428 from DCT components (= Equations (3) and (4)).

$$x(n) = \sum_{k=1}^L w(k)y(k) \cos \frac{3(2n-1)(k-1)}{2L} \quad n = 1; 2; \dots; L \quad (3)$$

where

$$w(k) = \begin{cases} 1 & k = 1 \\ r \frac{2}{L} & 2 \leq k \leq L \end{cases} \quad (4)$$

429 We make use of iDCT to assess how many DCT components are necessary to faith-
430 fully represent the actual articulatory trajectories. We do this by fitting DCT components
431 to a set of trajectories and then resynthesizing using iDCT with progressively more DCT
432 components. In this way, we can observe how increasing the number of DCT components
433 improves the precision of the representation. Figure 5 shows representative results, from
434 one speaker and one item ([ʃutokou] produced by Speaker 7). The improvement from 1
435 DCT component to 2 is substantial, as is the improvement from the 2 components to 3
436 components. With four components the correlation between the raw trajectories and the
437 iDCT-simulated trajectories reaches $r = 0.99$. In our case, only a small number of DCT

438 components (3 or 4) are required to faithfully represent articulatory trajectories over the
439 target VCVCV window. This result is similar to past studies, which have modelled tra-
440 jectories of similar duration and linguistic complexity using either 3 (Shaw & Kawahara
441 2018a) or 4 (Shaw & Kawahara 2018b; Kawahara et al. to appear) components.

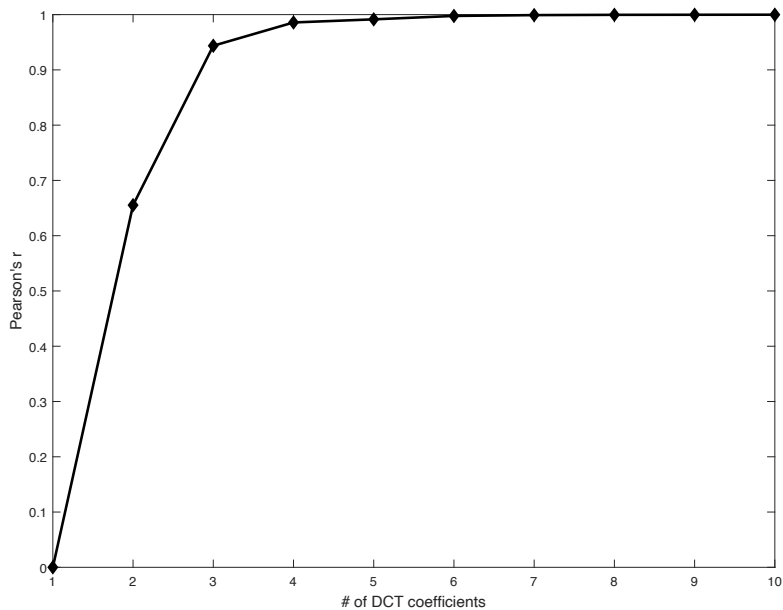


Figure 5: The increase in Pearson coefficients between the number of DCT components and the correlation between actual trajectories and simulated trajectories.

442 We can also use iDCT to illustrate how each DCT component contributes to the rep-
443 resentation of the articulatory trajectory. The top panel of Figure 6 shows the average
444 articulatory trajectories for each item of the dyad, [ʃutokoo] (left) vs. [ʃudooken] (right).

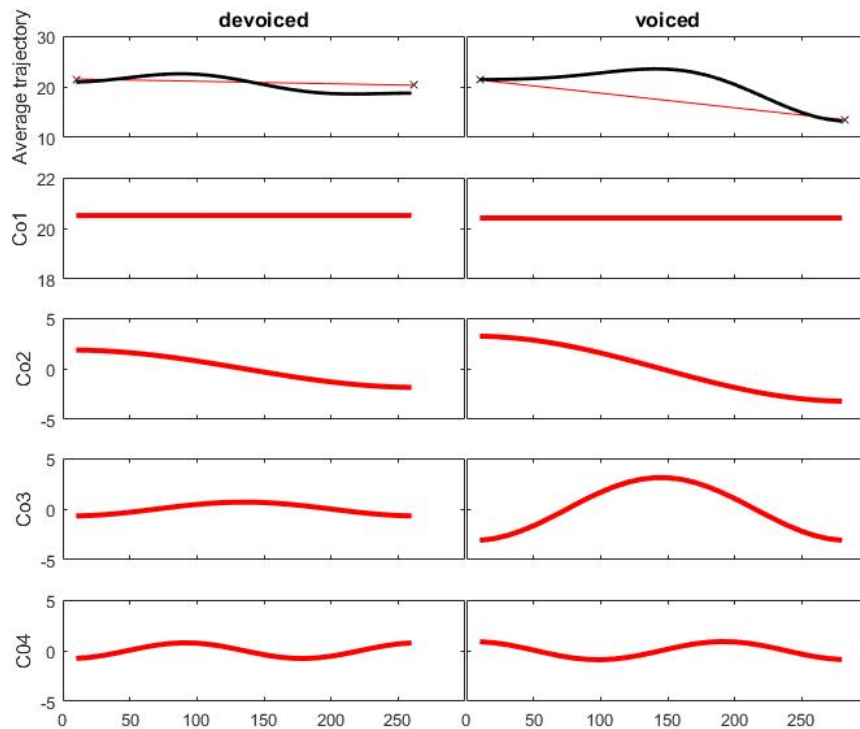


Figure 6: A sample comparison between the four DCT components of articulatory trajectories of devoiced and voiced tokens (averaged). The top panel shows the signal, with the 'x' marking the average height at the beginning and end of the trajectories and the line between the 'x's indicating linear interpolation.

445 Given this dyad, we can observe that the average change in tongue dorsum height over
 446 time, shown in the top panel, is noticeably different between devoiced and voiced items.
 447 For the voiced item (right), the tongue dorsum rises in the middle of the trajectory for
 448 [u]. For the devoiced item, there is less variation in the positional signal over time. For
 449 reference, the "x"s in the top panel show the average position at the start and end of the
 450 analysis window. The straight line connecting the x-points is equivalent to a linear inter-
 451 polation of spatial position across the analysis window. The panels below the trajectory

452 show the contribution of each DCT component to spatial modulation of the signal. The
453 duration of the simulated iDCT is based on the average duration of the tokens.

454 Comparison across devoiced and voiced items reveals similar modulations for the first
455 coefficient (Co1) and the second coefficient (Co2). The main difference is in the third
456 (Co3) and fourth (Co4) coefficients. Co3 picks up on the large rise for [u] in the voiced
457 case.⁸ The magnitude of the rise contributed by Co3 is greatly reduced for the devoiced
458 item compared to the voiced item. Finally, the fourth DCT coefficient (Co4) is also quite
459 different between voiced and devoiced items but it has only a small effect on spatial posi-
460 tion overall.

461 The next step is to assess whether the devoiced item contains a vowel target or not.
462 To do this we set up stochastic generators of our competing hypotheses, which we use for
463 Bayesian classification. The “target present” hypothesis is based on the voiced member
464 of each dyad. Specifically, since we have multiple repetitions of each item, we can cal-
465 culate a distribution over each DCT component. The normal distribution is characterized
466 by a mean value and a standard deviation. Thus, the mean and standard deviation of each
467 DCT component characterizes a normal probability distribution function. For the “target
468 absent” case, we adopt the common assumption that, in the absence of phonological spec-
469 ification, the trajectory will interpolate between surrounding targets (Choi, 1995; Cohn
470 1993; Keating, 1988; Pierrehumbert & Beckman, 1988). We therefore construct prob-
471 ability distributions for the “target absent” hypothesis that capture a realistically noisy
472 interpolation. For each token of a devoiced item, we fit DCT components to the straight
473 line connecting the position at the onset and offset of the analysis window.⁹ The average of
474 these components defines the probability distributions for the “target absent” hypothesis.
475 The standard deviation for the distributions is computed from the devoiced trajectories in
476 the same manner as for the voiced item. Consequently, the probability distributions that
477 characterize the “target absent” hypothesis are defined by linear interpolation (means of
478 the distribution) and the variability around each DCT component in the data. An example

⁸We note however that it is not necessarily the case that each DCT coefficient has to have a meaningful linguistic interpretation; neither is it the case that we have reasons to believe that Co3 is solely responsible for representing the tongue dorsum raising gesture of [u].

⁹See Pierrehumbert (1980) and Myers (1998) for cases of non-linear interpolation. We will reexamine this analytical choice of ours in 5.3

479 of the resulting distributions is provided in Figure 7. The horizontal axis is the value of
480 the coefficient, i.e., y in Equation (1), and the vertical axis is probability.

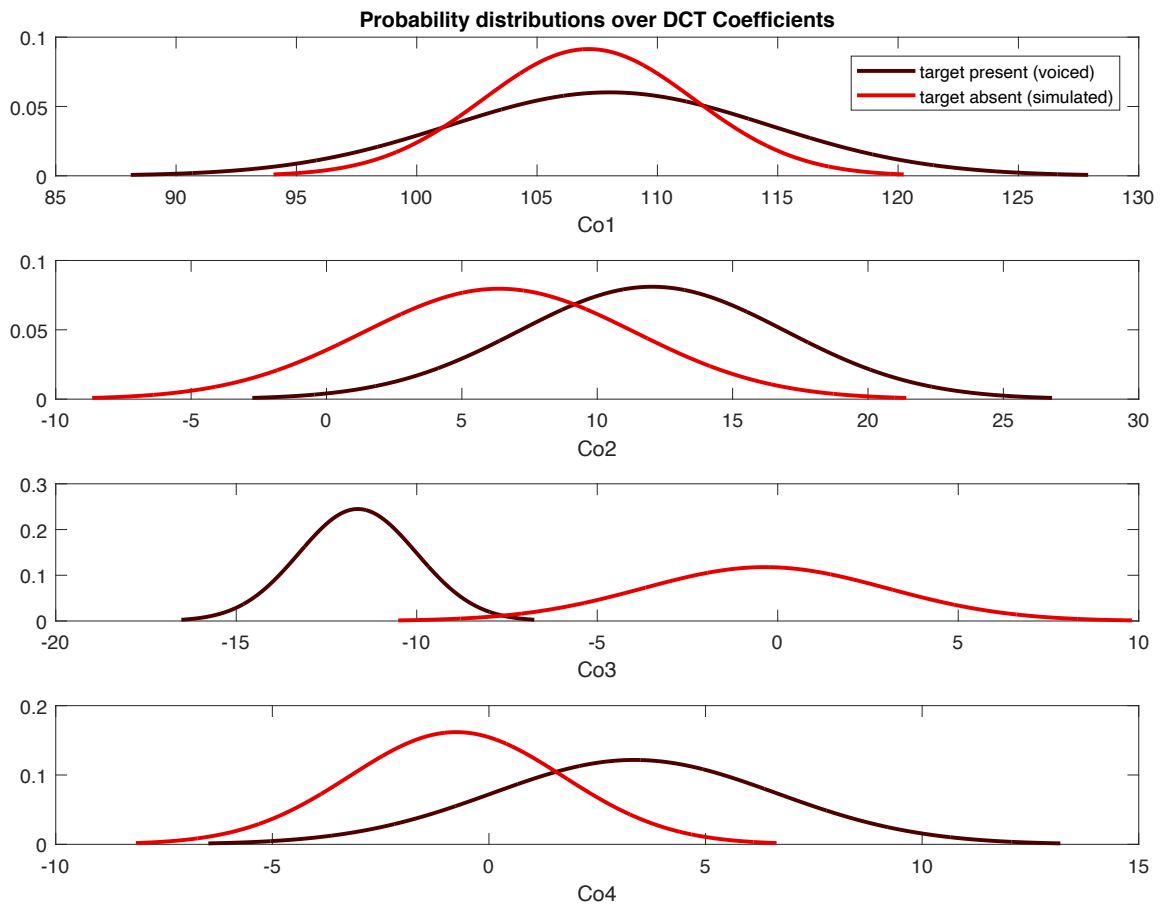


Figure 7: Probability distributions for DCT coefficients for the two competing hypotheses. The “target present” condition is based on the voiced vowels. The “target absent” condition is based on linear interpolation and the level of variability in the devoiced vowels.

481 We observe that the distributions for Co1 between the two conditions overlap heavily.
482 For Co2, there is a small difference between the “target present” distributions, based on
483 voiced vowels, and “target absent” distribution, based on linear interpolation. The largest

484 difference appears to lie in Co3. Naturally, the mean of the “target absent” distribution
 485 is very close to zero, and the same goes for Co4. This is because there is no rise for the
 486 straight line fit connecting the positional signal at the onset and offset of the analysis win-
 487 dow. The “target absent” Co3 distribution is also more variable than the corresponding
 488 “target absent” distribution—this difference reflects greater variability across devoiced to-
 489 kens than voiced tokens in whether the trajectory showed a rise characteristic of a vowel
 490 or not.

491 As the final step of the computational analysis, for each devoiced token, we determined
 492 the posterior probability of a vowel height target, based on Bayesian classification of the
 493 tongue dorsum trajectory (=Equation (5)). The posterior probability of the targetless hy-
 494 pothesis given the set of DCT coefficients (the left term of the Equation) is expressed as
 495 the prior probability of the targetless hypothesis—always set to be 0.5 in the current analy-
 496 sis, i.e, a uniform prior—multiplied by the product of the conditional probabilities of each
 497 DCT coefficient given the targetless hypothesis (i.e. linear interpolation), normalized by
 498 the denominator term. The classifier was trained on the distributions described above (see
 499 Figure 7) for voiced tokens, which unambiguously contain a vowel target, and a noisy null
 500 hypothesis, defined as linear interpolation across the target interval.

$$p(T|C_{o_1}; \dots; C_{o_n}) = \frac{p(T) \prod_{i=1}^n p(C_{o_i}|T)}{\prod_{i=1}^n p(C_{o_i})} \quad (5)$$

501 To summarize, the approach described in this subsection assigns a probability of target
 502 absence to each token. It does so by considering the probability that the token follows a
 503 linear interpolation as opposed to the trajectory of voiced vowels.

504 4 Results

505 Figure 8 shows the posterior probability of target absence for each condition by each
 506 speaker. The figures are violin plots which show the distribution of posterior probabilities
 507 of target absence. Points around the high y-axis region are tokens with a high probability
 508 of target absence, i.e., lingual movements that can be characterized as linear interpolation
 509 through the devoiced portion of the signal. Those at the bottom of the y-axis are tokens

510 that have a high probability of a vowel target, i.e., lingual articulations that resemble the
 511 voiced tokens. Those in the middle range are intermediate between target present and
 512 target absent, indicating a spatially reduced vowel target.

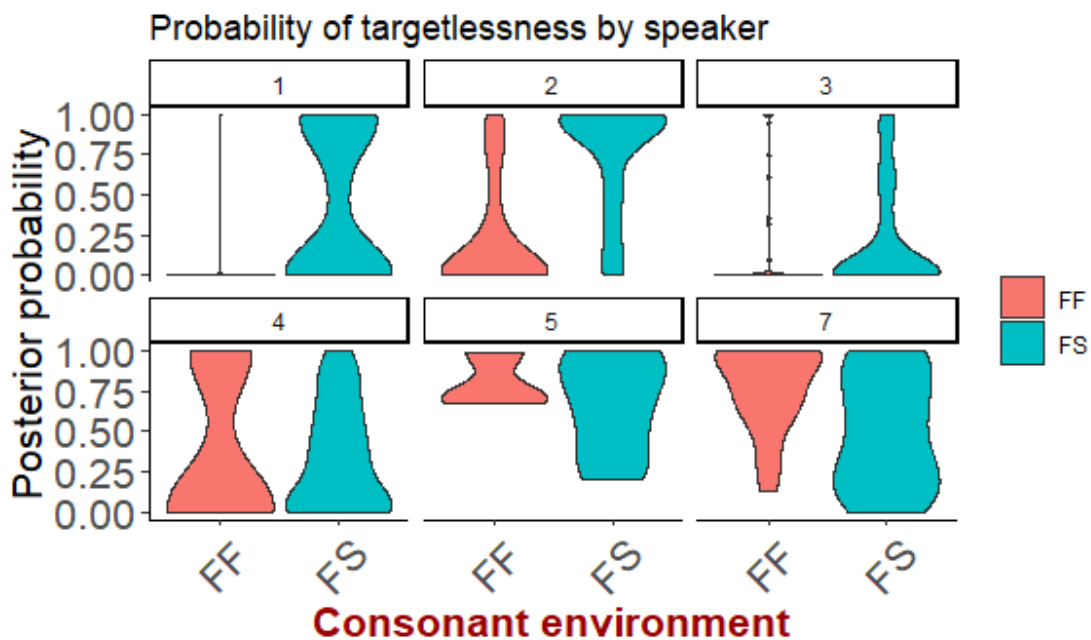


Figure 8: Posterior “target absent” probability for each condition by speaker. FF=Fricative-Fricative; FS=Fricative-Stop.

513 We observe that, as with [Shaw & Kawahara \(2018b\)](#), the distribution of posterior prob-
 514 abilities is bimodal. Across speakers, there tends to be a large probability mass at the high
 515 end of the probability scale (e.g., FS items for speaker 2 and speaker 5), at the low end
 516 of the probability scale (e.g., FF items for speaker 2, all items for speaker 3, FS items
 517 for speaker 4), or both (e.g., FS items for speaker 1, FF items for speaker 4). In many
 518 conditions, items skew towards the high and low ends of the scale. This is not to say that
 519 there are no intermediate items, which we take to be reduced. There are several cases with
 520 probability mass in the middle range, e.g. the FF condition for Speakers 5 and 7. Overall,
 521 however, the by-speaker view shows a tendency to either fully retain the lingual gesture

522 or entirely lose it. The one possible exception is FF items for speaker 5, the only plot of
523 12 in Figure 8 which does not have the majority of the probability mass at one end of the
524 scale. This result replicates the findings by Shaw & Kawahara (2018b) with a new set
525 of speakers and an expanded set of stimuli. Recall that the study by Shaw & Kawahara
526 (2018b) examined only four dyads; the current results are based on twelve dyads.

527 How the flanking consonants influenced targetless probability varied between speak-
528 ers. Speaker 1 showed almost no targetless tokens in the FF condition, but showed some
529 targetless tokens in the FS condition. This pattern—more targetlessness in the FS condi-
530 tion than in the FF condition—accords well with the prediction laid out in 1.3. Speaker 2
531 shows a similar, and perhaps clearer, pattern; this speaker showed rather consistent target-
532 present production in the FF condition, but typically deleted tongue dorsum raising gesture
533 in the FS condition. The pattern exhibited by Speaker 3 is less clear, but is also consistent
534 with the hypothesis presented in 1.3: almost no targetless tokens in the FF condition,
535 but greater probability of targetlessness in the FS condition. These three speakers thus
536 confirmed the hypothesis that we formulated in 1.3.

537 However, not all speakers behaved as we hypothesized. Speaker 5, especially in the
538 FF condition, seems to show some tokens whose posterior probabilities are in the middle
539 range—those tokens that are neither clearly targetless nor have a full target. Speakers 4 and
540 7, especially the latter, showed a pattern that is opposite from what is predicted from the
541 considerations discussed in 1.3—more targetless tokens in the FF condition than in the
542 FS condition. Thus, looking across the six speakers, we observe speaker-specific variation
543 in whether FF or FS environments conditions more deletion of the tongue dorsum raising
544 gesture.

545 Figure 9 shows the results by item. From this plot we can see some variability across
546 items as well. For example, [Fusagaru], the only verb in the item list, shows the lowest
547 probability of targetlessness. Many words show fairly sharp bi-modal patterns, with some
548 tokens showing high probability of targetlessness and others showing high probability of
549 full targets with few intermediate tokens. This bi-modal pattern applies especially clearly
550 to [Futa], [Futan], [Futon], [\$usa], and [\$utokou]. In contrast, most tokens of [\$usai] are
551 intermediate, with few extreme probabilities in either direction.

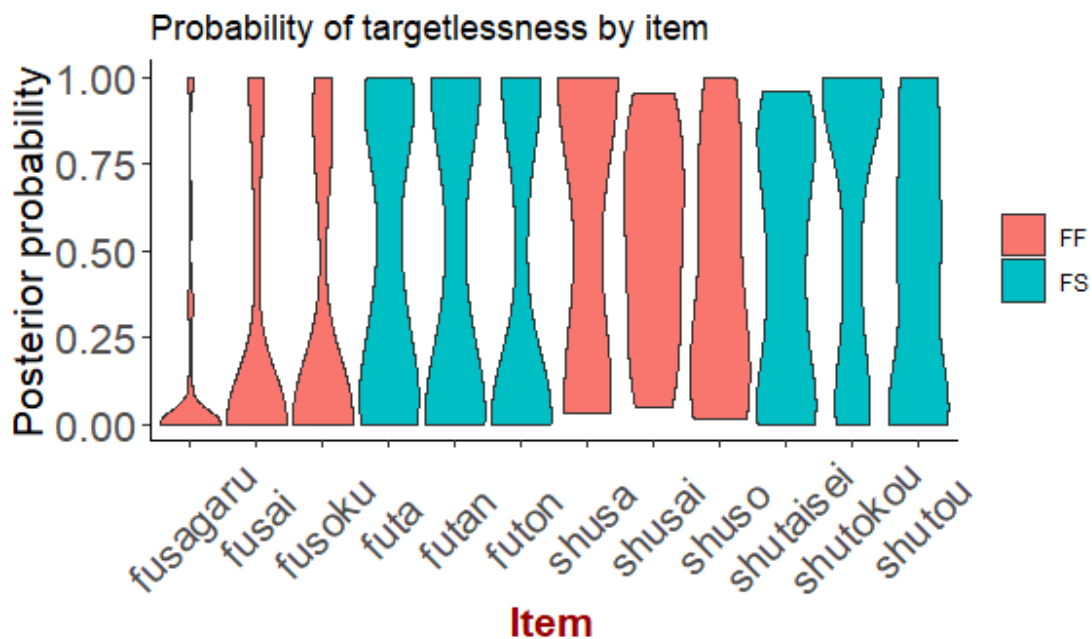


Figure 9: Posterior “target absent” probability by item. “f” and “sh” are used in the figure in place of [F] and [ʃ], respectively.

552 To assess the overall results statistically, we fit a series of nested linear mixed effects
 553 models in (6). The results of model comparisons appear in Table 2. The baseline model,
 554 m_0 , was compared to m_1 ; then m_2 and m_3 , which have the same number of parameters,
 555 were compared to m_1 . Finally, m_4 was compared to m_3 . The dependent variable was
 556 the posterior probability of deletion. Since probabilities are bounded dependent variables
 557 (upper bound of 1; lower bound of 0), we also ran the same models on arcsin-transformed
 558 probabilities. The same pattern of results came out of both raw and transformed probab-
 559 ities. For reasons of space we report results based on the non-transformed probabilities.
 560 The key fixed effect of interest was the consonant environment, coded as a two level fac-
 561 tors, FF vs. FS (“*Cond*”). Speakers and items were treated as random intercepts.

$$\begin{aligned}
m0 &: \text{post } \frac{3}{4} (1jspeaker) + (1jitem) & (6) \\
m1 &: \text{post } \frac{3}{4} (1 + Condjspeaker) + (1jitem) \\
m2 &: \text{post } \frac{3}{4} Cond + (1 + Condjspeaker) + (1jitem) \\
m3 &: \text{post } \frac{3}{4} C1 + (1 + Condjspeaker) + (1jitem) \\
m4 &: \text{post } \frac{3}{4} C1 \& Cond + (1 + Condjspeaker) + (1jitem)
\end{aligned}$$

Table 2: Summary of model comparisons.

	df	AIC	BIC	logLik	deviance	χ^2	χ^2 df	p
<i>m0</i>	4	464.7	481.7	-228.3	456.7	–	–	–
<i>m1</i>	6	402.6	428.1	-195.3	390.6	66.07	2	< .001
<i>m2</i>	7	404.1	433.8	-195.0	390.1	0.53	1	<i>n:s:</i>
<i>m3</i>	7	400.4	430.1	-193.2	386.4	4.25	1	< .05
<i>m4</i>	9	403.7	441.9	-192.8	385.7	4.95	3	<i>n:s:</i>

562 The baseline model, *m0*, includes only the random effects. The next model, *m1*, adds
563 a by-speaker random slope for the fixed effect, i.e. surrounding consonants (FF vs. FS) to
564 this model. The by-speaker random slope improved the model significantly. This result
565 indicates that speakers show different sensitivities to the consonantal environments. As
566 we observed in Figure 8 some speakers (e.g. Speakers 1 and 2) show less deletion in FF
567 than FS environments, while others (Speakers 4 and 7) show the opposite pattern.

568 Because the effect of consonant environment differs by speaker, the average effect of
569 consonantal environment is not predictive. These statistical comparisons support what
570 we observed in Figure 8: different speakers are sensitive to consonantal environment in
571 different ways.

572 We also ran models that included the C_1 type ([F] vs. [S]) and the interaction between
573 C_1 and consonant environment (“*Cond*”) as fixed factors. The addition of C_1 led to im-
574 provement over *m1*, and was marginally significant within the model ($b = 0.098; t =$
575 $2.136; p = 0.055$), indicating that deletion probability is slightly higher when C_1 is [S]

576 than when C1 is [F]. The interaction between C₁ and consonant environment (“*Cond*”)
 577 did not lead to further improvement, indicating that the effect of C₁ is not dependent on
 578 the consonant sequence. Thus, our best fitting model, m3, includes a (“*Cond*”) as random
 579 effect but not as a fixed effect.

580 Figure 10 shows the by-speaker random slopes for our best fitting model. The x-axis
 581 shows the estimate for FS sequences. As we observed in the violin plots of probabilities
 582 (=Figure 8), Speakers 1 and 2 have positive estimates, indicating that deletion is more
 583 likely in FS sequences than in FF sequences. Moreover, the confidence intervals around
 584 the estimate do not overlap with zero. Additionally, as we also observed above, Speaker
 585 7 shows the opposite pattern. This speaker has a negative estimate, which also does not
 586 overlap with zero, indicating significantly higher probability of targetlessness in FF se-
 587 quences than in FS sequences. The other speakers have estimates that overlap with zero,
 588 indicating an effect that is not statistically significant.

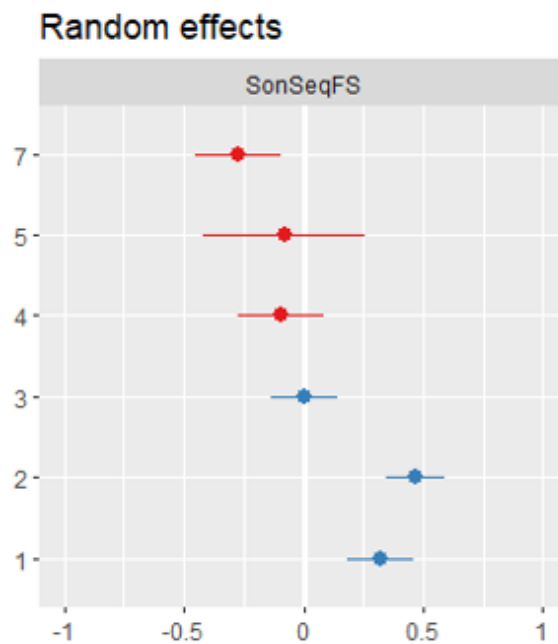


Figure 10: By-speaker random slopes for the effect of sonority sequencing (= *Cond*). The estimate is for the FS condition, relative to FF.

589 In summary, consonant environment had a significant impact on deletion probability,
590 but the direction of the effect was not uniform across speakers. Some speakers showed
591 consistently more deletion in FS, as predicted, others showed more deletion in FF, or no
592 effect of consonant context.

593 **5 Discussion**

594 **5.1 Summary**

595 The current experiment replicated the core finding of Shaw & Kawahara (2018b) with a
596 new set of speakers and an extended set of stimuli. The posterior probability of vowel
597 presence/absence showed a bimodal distribution for many speakers (see, Figure 8) and
598 items (see, Figure 9). One mode was centered on the low end, near zero probability of
599 vowel absence. These devoiced vowel tokens were produced with tongue height trajec-
600 tories very similar to voiced vowels. The other mode of the distribution was centered on the
601 high end, indicating that the tongue height trajectory resembled our noisy null hypothesis,
602 a linear interpolation between flanking vowel targets. These modes of the posterior prob-
603 ability distribution represent endpoints on a continuum from a full target to no detectable
604 vowel target. A mono-modal distribution centered between 0 and 1 would have provided
605 evidence for consistent vowel reduction, i.e., a vowel height target of reduced magnitude.
606 Although we did also see some tokens with intermediate probabilities, the variation clus-
607 tered more around the high and low ends of the scale, a similar pattern reported in Shaw
608 & Kawahara (2018b).

609 The results also revealed some systematic patterns in how the flanking consonants
610 influence deletion probability. The design of the study featured conditions contrasting
611 devoiced vowels intervening between fricative-fricative (FF) sequences and fricative-stop
612 (FS) sequences. The original hypothesis developed in 1.3 is that was that we would
613 observe more deletion in FS sequences than in FF sequences. Recall that, to the extent
614 that we can conceive of deletion as an extreme instantiation of devoicing, either in terms
615 of reduction or enhancement, we would expect targetless tokens to be more likely in the FS
616 condition than in the FF condition, because devoicing is more likely in this environment.

617 Syllable contact laws (Murray, 1988; Murray & Vennemann, 1983), if Japanese speakers
618 are sensitive to them, also predict this pattern. Our hypothesis was also motivated by an
619 empirical observation. Shaw & Kawahara (2018b) found that, even though the speakers
620 in the study differed substantially in their individual rates of vowel deletion, all speakers
621 deleted devoiced vowels more often in [\$utaisei], resulting in a FS consonant sequence,
622 than in [Fusoku], resulting in a FF sequence.

623 The current study revealed inter-speaker variability with respect to the prediction laid
624 out in 1.3: some speakers showed more targetless tokens in the FS condition than in the
625 FF condition (Speakers 1 and 2), as we initially hypothesized, and some speakers showed
626 the opposite pattern (Speaker 7, and to a less clear extent, Speaker 4).

627 Our items in the FF and FS condition both featured two fricatives, [F] and [s]. Al-
628 though we did not predict this differences, there was a significant effect of fricative, with
629 higher deletion probability following [s] than [F]. Moreover, this effect is significant in
630 a group analysis while consonant sequence was only significant as a by-subject random
631 slope. Quite possibly, the observed difference in deletion probability between [\$utaisei]
632 and [Fusoku] in past work as well is attributable not to the consonant manner sequence,
633 FF vs, FS, but to the identity of the initial consonant.

634 5.2 Time and target undershoot in DCT representations

635 Our approach to analyzing time-varying kinematic data in terms of discrete hypotheses
636 makes use of a low parameter stochastic representational space. Time varying signals, in
637 this case tongue dorsum height trajectories, are represented as modulations of frequency
638 components, using DCT. The DCT coefficients effectively represent the signal with high
639 precision but without directly encoding the temporal duration of the trajectories. Instead,
640 time is indirectly encoded in the frequencies of the DCT components. The representation
641 of time is indirect because it comes in the form of what frequencies are represented in each
642 component, which is dependent on the analysis window.

643 We represented all trajectories in this study using just four DCT components. Since
644 the frequency of the DCT components vary as a function of the length (in samples) of a
645 trajectory (see (1)), they have the potential to indirectly encode the duration of the tra-

646 jectory. For example, past work has shown that DCT representations alleviate the need
647 to represent temporal duration independently. For example, [Watson & Harrington \(1999\)](#)
648 compared several methods of representing time-varying formants, including DCT repre-
649 sentations, in a study of Australian vowels. They showed that adding vowel duration to
650 the representation of Australian vowels improved machine classification in many cases.
651 When Australian vowels were represented by measurements of formants at percentages of
652 total vowel duration, vowel duration was needed as an additional factor to reach a high-
653 level of classification accuracy. This is because several Australian vowel pairs have very
654 similar (possibly indistinguishable) vowel quality but differ in duration ([Cox & Fletcher](#)
655 [2017](#)). However, when [Watson & Harrington \(1999\)](#) represented the same vowels with
656 DCT components only, vowel duration did not improve classification accuracy. Two DCT
657 components fit to the first and second formants were sufficient to classify all 19 Australian
658 vowels, including vowels differentiated primarily by duration.

659 Since DCTs can represent both the spatial modulation and the temporal duration of a
660 signal, we cannot know if one of these dimensions or the other had a dominating influence
661 on our classification results. Although high vowel devoicing in Tokyo Japanese occurs
662 at both fast and slow speech rates ([Fujimoto 2015](#)), we do not know if vowel deletion
663 is likewise rate independent. Conceivably, the probability of detecting a vowel move-
664 ment decreases at fast rates due to target undershoot ([Lindblom, 1963](#); [Moon & Lindblom](#)
665 [1994](#)). To investigate this, we evaluated the correlation between the duration of our target
666 intervals, as a measure of local speech rate, and the posterior probability of deletion. Fig-
667 ure [11](#) shows a scatter plot of these two variables. There was a weak negative correlation
668 ($r = 0.11; p < .05$), indicating that the probability of targetlessness decreases at slower
669 speech rates (longer duration).

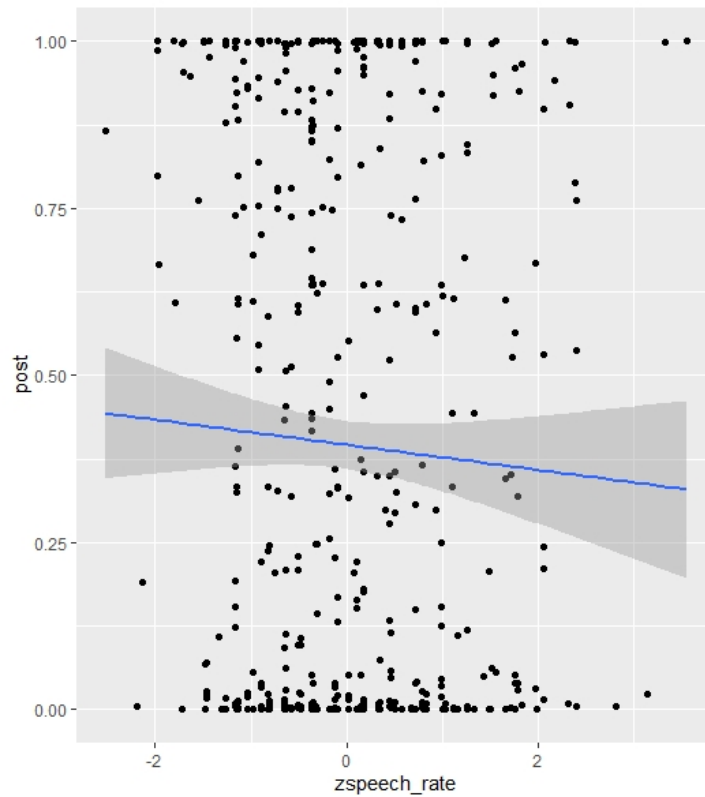


Figure 11: Correlation between speech rate, represented by a Z-scored of target trajectory duration (x-axis) and the posterior probability of targetless (y-axis).

670 To further investigate the influence that speech rate might have on our deletion prob-
 671 ability results, we subsetted the data into relatively short and relatively long tokens. Our
 672 short-ish tokens were those that were less than one standard deviation from the mean to-
 673 ken duration; our long-ish tokens were those that were greater than one standard deviation
 674 from the mean. This subsetting procedure produced 74 tokens (14.4% of the data) for the
 675 short group and 76 tokens (14.8%) for the long group. We looked at the distribution of
 676 long and short tokens across speakers and found that all speakers produced some tokens
 677 that fell into the long group and some that fell into the short group. The mean duration of
 678 the CV interval in the short group was 228 ms. The mean duration of the CV interval in

679 the long group was 362 ms. Figure 12 compares the posterior probability of deletion for
680 the long (slow local speech rate) and short (fast local speech rate) data subsets. Consistent
681 with the weak correlation between speech rate and targetlessness across the entire corpus,
682 we see a slight increase in targetlessness probability for the short data subset. This is the
683 case for both FF and FS consonant manner sequences. Notably, however, a substantial
684 number of tokens still show a high probability of targetlessness at slow speech rates. This
685 indicates that while increased speech rate may contribute to targetlessness, based on the
686 diagnostic methods employed here, there are still tokens that approximate a linear inter-
687 polation trajectory even at the slowest speech rates in the data set. This indicates that, like
688 high vowel devoicing, vowel deletion, or at least extreme reduction of the tongue dorsum
689 height target, also occurs at slow end of natural speech rate variation. This result implies
690 that whether or not to retain a tongue dorsum gesture is under speakers' control, rather
691 than an automatic consequence of fast speech¹⁰

692 5.3 Minimal paths for targetless trajectories

693 One of the challenges of assessing whether the tongue dorsum height target is completely
694 absent or just heavily reduced is that there are no unequivocal FF or FS sequences in
695 Japanese that could serve as a baseline for assessing whether pronunciation of /FuF/ and
696 /FuS/ deviate enough from these underlying forms to conclude that they are indeed [FF]
697 and [FS]. Our approach to this challenge is to simulate tongue dorsum trajectories that
698 interpolate between vowels, V1 and V2, in /V1CCV2/. Our simulations in this paper are
699 based on two assumptions: (1) first, movements take the minimal path between targets and
700 (2) second, like all biological signals, there will be variability in the movement trajectory.
701 We calculated the minimal path as a linear interpolation between vowel targets and we
702 modelled variability as random deviations from the minimal path. The magnitude and
703 structure of the random deviations are based on the devoiced tokens in our corpus. In this
704 way, the variability injected into our simulations has the same item-specific and speaker-
705 specific properties of our corpus. The difference between the vowel-absent class, as we
706 simulated it, and the devoiced tokens in our corpus, is that the tongue-dorsum trajectory in

10

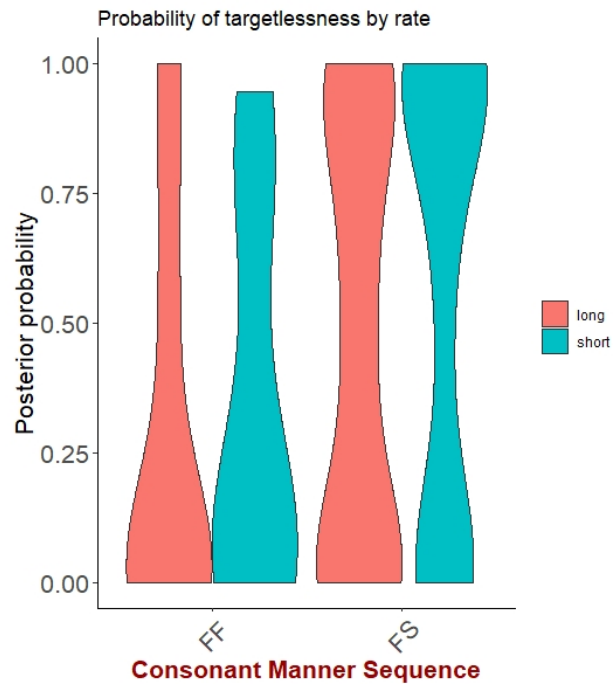


Figure 12: Posterior probabilities of the short-ish subset and long-ish subset.

707 the vowel-absent class is always guided by the minimal distance between V1 and V2. The
 708 degree to which the actual tongue dorsum trajectories in our devoiced tokens also follow a
 709 realistically noisy actuation of the the minimal distance path or whether they instead move
 710 towards an elevated tongue dorsum height target for [u] is represented in the results of our
 711 Bayesian classification. A substantial number of tokens were classified as belonging to the
 712 minimal distance path.

713 Our decision to simulate the vowel-absent tongue dorsum trajectory as taking the path
 714 of minimal distance between flanking targets is intended to be a theory-neutral decision.
 715 It is also possible to apply our method of analysis by simulation and classification with
 716 different theoretical assumptions about what the vowel-absent trajectory should look like.
 717 Here, we consider the predictions of Task Dynamics (Saltzman & Munhall 1989) as im-
 718 plemented in the Task Dynamics Application (TADA: Nam et al. 2004, 2012). One prop-
 719 erty of this model is that articulators that are not under direct phonological control (i.e. by

720 a gesture, in the sense of Articulatory Phonology: [Browman & Goldstein 1986](#) et seq.) at
721 a particular time are driven to a rest position by a neutral attractor. Because of the neutral
722 attractor, there are conditions under which articulators will not necessarily follow the min-
723 imal path between targets. Instead, articulators will return to a neutral position until they
724 are brought under control by another gesture. To explore how TADA predictions for the
725 vowel-absent case might differ from linear interpolation for the items in our study, we ran
726 a series of TADA simulations.

727 The first TADA simulation compares [eFta] and [eFuda]. There are a number of manip-
728 ulable parameters in TADA, and variation in some of these parameter settings has been hy-
729 pothesized to capture cross-language variation, i.e., language-specific phonetics ([Iskarous
730 et al. 2012](#)). To minimize researcher degrees of freedom ([Roettger 2019](#)), we used default
731 TADA gestural parameters whenever reasonable for Japanese. For the [eFta] vs. [eFuda]
732 comparison, we used default parameters for [e], [f] for [F], [t], [d], and [a]. The only ges-
733 ture that required manipulation to approximate Japanese-specific phonetics was [u]. The
734 default [u] in TADA produces a much longer vowel, 300 ms, than is typical in Japanese,
735 and it produces a vowel with lip protrusion. To adapt the gesture parameter settings for
736 Japanese [u], which is much shorter, ca. 50 ms (e.g. [Shaw & Kawahara 2019](#)), and lacks
737 lip protrusion (e.g. [Vance 2008](#)), we eliminated the lip protrusion gesture and shortened
738 the activation duration of the tongue body gesture. The gesture parameter values for all
739 simulations are provided in the supplementary materials.

740 Figure [13](#) compares the trajectories for [eFta] and [eFuda] simulated by TADA. The top
741 panel shows the simulated waveform. The bottom three panels show kinematic trajectories
742 in the vertical dimension for the tongue dorsum, tongue tip and lower lip. The tongue
743 dorsum trajectory for [eFta] has a mid-level plateau for [e], in the temporal window from
744 0 to 250 ms, and then falls to [a]. The tongue dorsum trajectory for [eFuda] starts with a
745 similar plateau for [e] but then rises for [u]. The peak of the rise comes near the end of the
746 voicing period for the vowel and remains rather high during the [d] before falling for [a].
747 The data simulated with TADA are qualitatively quite similar to our experimental data. For
748 comparison with representative tokens from the experimental data, see Figure [3](#). For this
749 particular case, our theory-neutral choice of linear interpolation for “vowel-absent” tokens
750 is quite similar to the TADA simulations, which also show a roughly linear trajectory. It

751 should be noted, however, that this linearity is not a general prediction of TADA. It follows
 752 in part from the properties of our stimulus items. The progression of vowel height targets
 753 from mid, [e], to low, [a], does not involve a neutral attractor driving the tongue dorsum
 754 height away from the minimal path between these vowels. For items such as [eFta], there
 755 would be little difference between using linear interpolation between flanking vowels and
 756 using TADA simulations, with default gesture parameters.

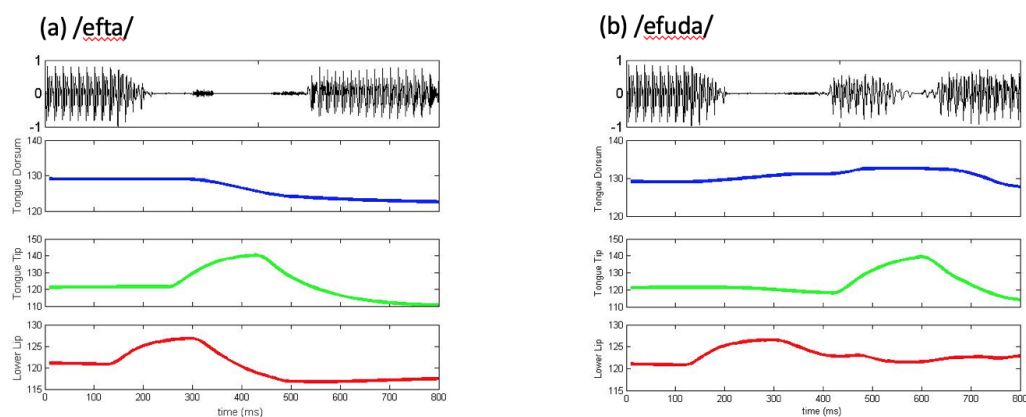


Figure 13: TADA simulations of [eFta] and [eFuda]

757 We now move on to [eʃta] and [eʃuda]. Figure 14 shows TADA simulations of these
 758 items. The top two panels show simulation results with default gesture parameters for
 759 all segments except for [u], which used the same Japanese-specific parameters described
 760 above. Of relevance is that the default gestures for [ʃ] include both a tongue body gesture
 761 and a tongue tip gesture. For Japanese, our materials were not designed to assess the
 762 presence/absence of a tongue body gesture for the fricative, [ʃ], directly (see S5.5 for an
 763 indirect attempt). The Japanese fricative has different acoustic and articulatory properties
 764 from the English post-alveolar fricative, but it is unclear whether the difference is due
 765 to the tongue body gesture or to other aspects of fricative production, including a labial
 766 component, tongue-tip constriction area, or relative degree of tongue grooving. Because of
 767 this uncertainty, we also ran TADA simulations with the fricative unspecified for a tongue
 768 body gesture. This result is shown in the bottom panel of Figure 14

769 When [ʂ] was simulated without a tongue body gesture, the difference in tongue dor-
 770 sum trajectories between [eʂta] and [eʂuda] is nearly identical to the difference found for
 771 [eʃta] and [eʃuda]. That is, the tongue dorsum height trajectory follows a roughly linear
 772 path from [e] to [a] in [eʂta] but it rises for [eʂuda]. However, when [ʂ] is specified with
 773 a tongue body gesture, then we see a rise in the tongue dorsum height trajectory in [eʂta],
 774 which disrupts the linearity of the transition from [e] to [a], even in the absence of [u].

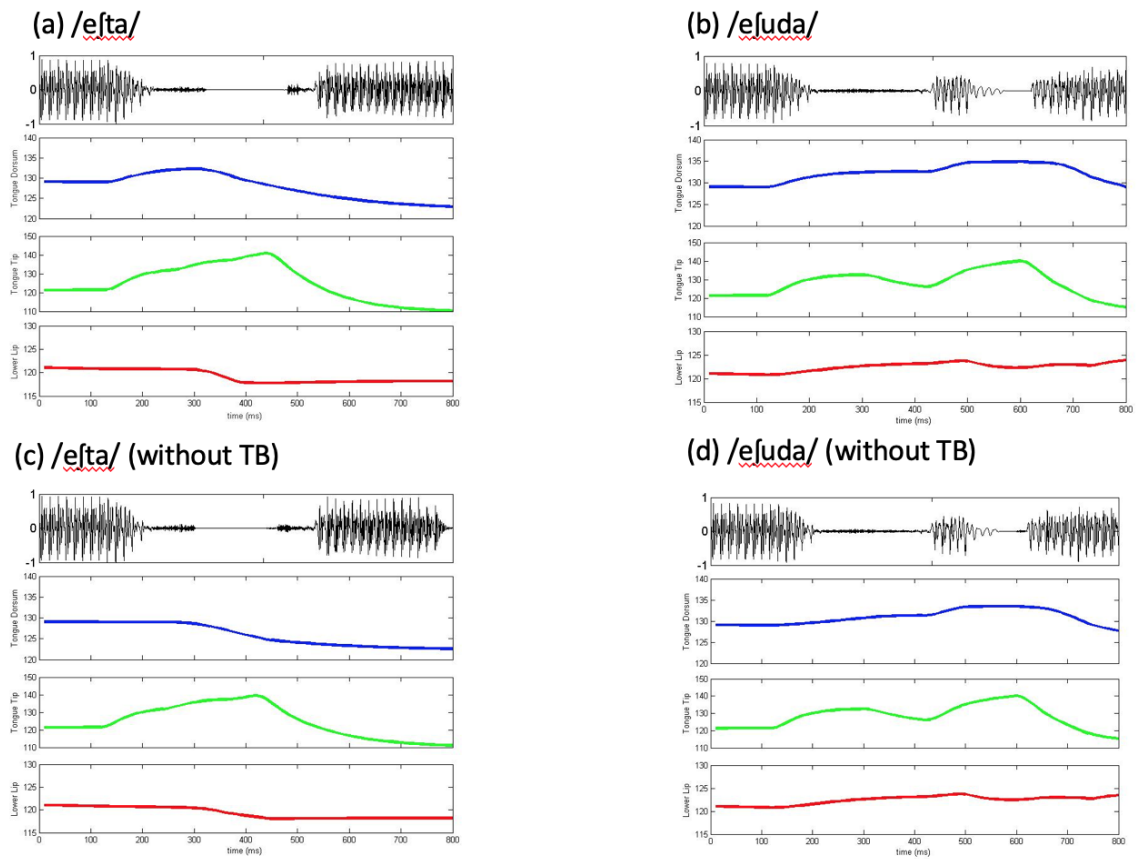


Figure 14: TADA simulations of [eʂta] and [eʂuda] with (top row) and without (bottom row) TB gesture.

775 The case of [ʂ] specified with a tongue body gesture allows us to consider how using a

776 theory-specific alternative to the minimal path assumption might influence our results. If
777 we used (a stochastic version of) the TADA simulation trajectory for [e\$ta] and [eFta] as the
778 basis for our Bayesian classification (instead of linear interpolation), we would introduce
779 a bias in deletion likelihood towards the [\$] environment over the [F] environment. This is
780 because, to detect a vowel in the [F] environment, the trajectory would only have to rise
781 above the linear trajectory in the TADA simulation (Figure 13, [eFta] panel). However, to
782 detect a vowel in the [\$] environment, the trajectory would have to rise above not just the
783 linear trajectory between vowels but also above the magnitude of the tongue body gesture
784 for [\$]. Deviations from minimal path would still be classified as deletion, if the magnitude
785 of the deviation did not exceed the tongue body magnitude for [\$]. In contrast, relative
786 to using a TADA baseline, if there actually is a tongue body gesture for [\$], the minimal
787 path method is biased towards finding more vowel deletion in the [F] environment than in
788 the [\$] environment. This is because increases in tongue body height, including those due
789 to [\$], will count as deviation from the minimal path, and push classification towards the
790 vowel present case.

791 Using the minimal path method, we observed significantly greater deletion in the [\$]
792 environment than in the [F] environment. If we had used a TADA-baseline with a tongue
793 body gesture for [\$], this result would probably have been even stronger. On the other
794 hand, if we had used a TADA baseline without a tongue body gesture for [\$], then there
795 is really not much difference between the minimal path method and a TADA baseline.
796 However, we reiterate that the similarity between TADA and minimal path is not a general
797 result—it is particular to the materials that we selected for this experiment. Additionally,
798 the above conclusions are based on default gesture parameters (with the exception of [u]),
799 which are appropriate for English, but might require fine-tuning in order to capture sys-
800 tematic differences across languages. Generally, there may be conditions under which a
801 minimal path baseline is inappropriate, or, at least, is inconsistent with the Task Dynamics
802 framework, as implemented in TADA.

803 With the above caveats in place, we conclude that the finding of more deletion in the
804 [\$] environment than in the [F] environment is likely robust to variation in how we might
805 simulate the vowel absent scenario. If there is indeed a tongue body gesture for [\$], the
806 minimal path method is biased against our finding, and yet it still emerged as statistically

807 significant.

808 **5.4 Tongue dorsum trajectories for voiced vowels**

809 In the last sub-section, we discussed how we simulated, for the purpose of classification,
810 trajectories lacking a vowel target. The other relevant factor in classifying devoiced trajec-
811 tories using our method is the trajectory of the corresponding voiced vowel. We defined a
812 separate classifier for each combination of speaker and item. This allows us to incorporate
813 any speaker-specific variation into the analysis. How a particular devoiced trajectory is
814 classified depends both on the degree to which it deviates from the minimal path as well as
815 the degree to which it deviates from the corresponding voiced vowel. Correspondence in
816 this case is based on the materials—we selected near minimal pairs matched on as many
817 relevant properties as possible. To facilitate appropriate generalization of our approach to
818 new data, we discuss some possible non-obvious implications of using a local (by speaker,
819 by item) baseline.

820 To illustrate the importance of the local baseline, we zoom in on a small subset of
821 the data, just the [F] environment tokens produced by Speaker 2. Recall that Speaker 2
822 was one of the speakers that produced a particularly sharp bimodal distribution in vowel
823 deletion probabilities and showed the predicted effect of consonant sequence (see Figure
824 8). Figure 15 shows three panels summarizing tongue dorsum trajectories for Speaker
825 2. The first panel shows the average tongue dorsum trajecory for voiced and devoiced
826 tokens. This was generated by fitting an SSANOVA, using the GSS package in R (Gu
827 2002), to the first 150 ms of each token. We choose 150 ms because it is the length of the
828 smallest analysis window for this speaker. The SSANOVA plot shows that, on average,
829 the devoiced trajectories are flatter than for the voiced trajectories. Note that this was not
830 the pattern for all speakers; Speaker 3, for example, showed very little difference between
831 voiced and devoiced trajectories. The second panel breaks down the devoiced tokens by
832 item. Looking across items, we see that [Fusagaru] seems to have the flattest trajectory.
833 From this figure, we might erroneously suspect that [Fusagaru] has the highest probability
834 of deletion. The third panel shows that this is absolutely not the case. In fact, for this
835 speaker, [Fusagaru] has the lowest posterior probability of deletion of any [F]-tokens. This

836 might seem puzzling. Why does [Fusagaru] have a low probability of deletion, given its
 837 relatively linear trajectory?

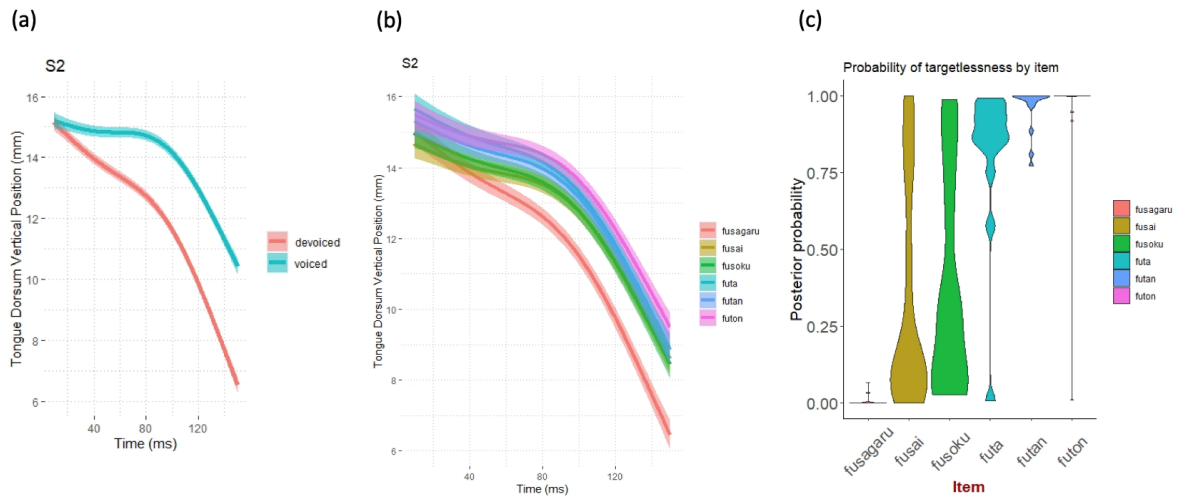


Figure 15: [F]-tokens for S2: (a) shows the average tongue dorsum height trajectory for voiced and devoiced vowels; (b) breaks down the devoiced trajectories by item; (c) shows posterior probability of vowel deletion by item.

838 The answer is in the patterning of the voiced vowel counterpart for the devoiced to-
 839 kens. Figure 16 plots [Fusagaru] along with its voiced vowel counterpart [Fuzakeru].
 840 The key observation is that the trajectory for [Fuzakeru], the voiced vowel counterpart to
 841 [Fusagaru] in our materials, also has a relatively flat trajectory. Because of this relatively
 842 flat baseline for the voiced vowel, the trajectory for [Fusagaru] does not have to depart
 843 very far from linearity to be classified as a vowel. The Speaker 2 voiced vowel baseline
 844 for [Futa] is quite different. As show in the right side of Figure 16, the tongue dorsum rises
 845 substantially for [Fuda], which serves as the voiced vowel baseline for assessing target-
 846 lessness in [Futa]. Given this baseline, a [Futa] token that shows only a minimal depart-
 847 ure from linearity will still have a higher probability of linearity than of a full vowel.

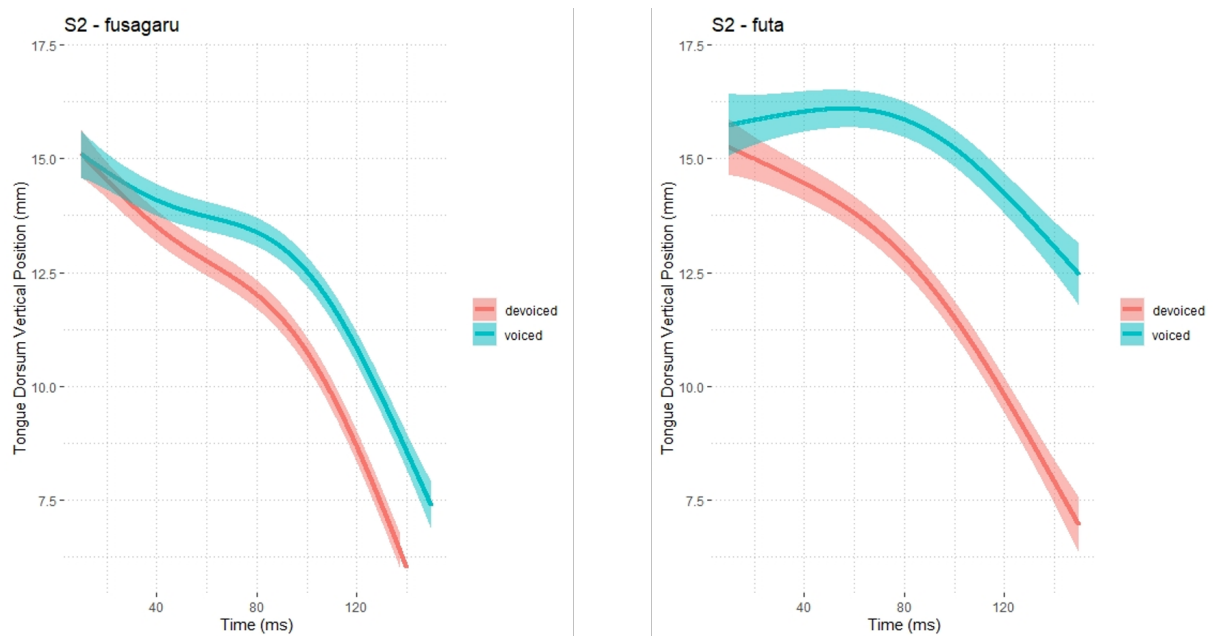


Figure 16: S02 tongue dorsum trajectories for two dyads: the left panel shows [Fusagaru] (devoiced) paired with [Fuzakeru] (voiced); the right panel shows [Futa] (devoiced) paired with [Fuda] (voiced) .

848 The case above serves to illustrate the role of the speaker- and item- specific baseline
 849 in our analytical approach. In assessing whether a given speaker produces a vowel, we
 850 pursue a very targeted machine learning approach that factors speaker-specific productions
 851 of baseline words in the analysis.

852 5.5 The effect of fricative place

853 We now return to the effect that fricative place of articulation had on vowel deletion proba-
 854 bility. For starters, we explore an indirect test of whether [ʃ] in Japanese has a tongue body
 855 gesture. As illustrated through TADA simulations (Figure 14), whether [ʃ] in Japanese
 856 has a tongue body gesture or not is an important consideration in interpreting our results.
 857 When we simulated [ʃ] without a tongue body gesture, then the tongue dorsum height

858 trajectory for [eFta] and [eSta] was very similar. As an indirect test of whether Japanese
859 [ʃ] has a tongue body gesture, we compare the distribution of DCT coefficients for all
860 voiced vowel tokens in our corpus. This includes all of the words with voiced vowels
861 that served as item-specific baselines for the devoiced items in both [F] and [ʃ] environ-
862 ments. Figure 17 compares the distributions. The distributions of all four DCT compo-
863 nents are heavily overlapped. Independent t-tests (Welch's two sample) show that differ-
864 ences are not significant for the first three DCT components: 1 ($t = 1.11; p = 0.267$),
865 2 ($t = 0.406; p = 0.685$), 3 ($t = 1.25; p = 0.214$). Only the fourth DCT component,
866 which explains only a small amount of variance in the trajectories (Figure 5), showed a
867 significant difference ($t = 4.87; p < .001$) across [F] and [ʃ]. Although this result cannot
868 be taken as conclusive evidence for the presence or absence of a tongue body gesture, it
869 does indicate that the trajectories, as represented by DCT coefficients in our classification
870 process, were quite similar across [F] and [ʃ]. This is despite the fact that [F] and [ʃ] tokens
871 were not completely balanced for vowel sequences and other properties (e.g. word length,
872 pitch accent placement, and vowel sequence).

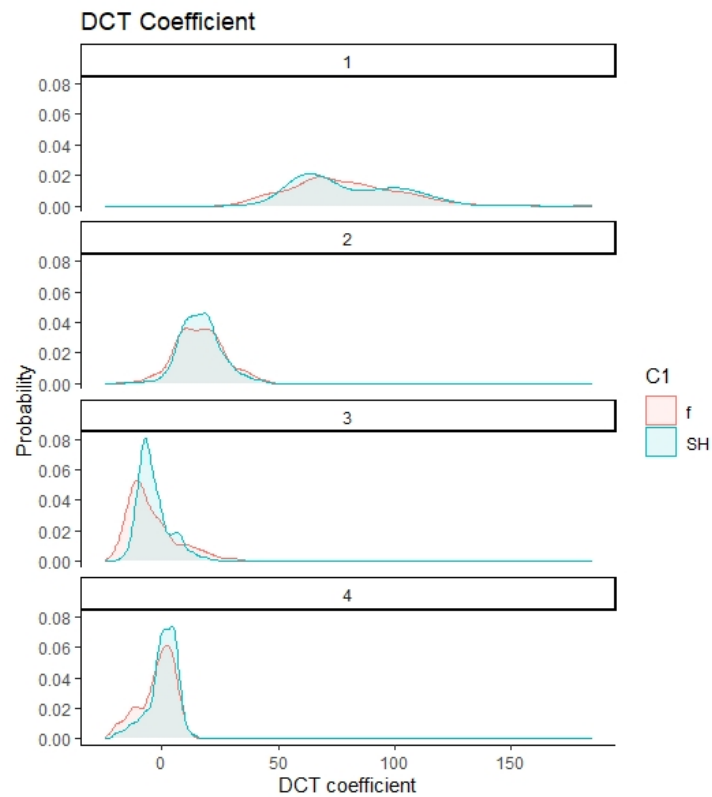


Figure 17: DCT distributions

873 Given the similarity of the DCT distribution of voiced vowel items across [ʃ] and [F],
 874 the difference between [ʃ]-initial items and [F]-initial items can be attributed to the tongue
 875 dorsum trajectory in the voiceless items. The trajectory of devoiced vowels is more likely
 876 to resemble a linear trajectory between flanking vowels when preceded by [ʃ] than when
 877 preceded by [F]. This result is independent of consonant sequence, i.e., FF vs. FS.

878 One possible explanation for the effect of fricative place on vowel deletion relates di-
 879 rectly to the goal of achieving vowel devoicing. While vowel devoicing does not serve a
 880 contrastive function, it does serve as a sociolinguistic marker of prestige in Tokyo Japanese
 881 (Imai, 2004), and there is evidence that it is under direct control, c.f., devoicing as a passive
 882 consequence of overlapping laryngeal gestures, as it may be in some cases of vowel de-

883 voicing in other languages (see [Fujimoto 2015](#) and other references cited in introduction).
884 One piece of support for the conclusion the devoicing is actively controlled in Japanese
885 comes from the observation of laryngeal gestures associated with voiceless stops ([Fujimoto 2015](#)).
886 When voiceless stops in Japanese precede voiced vowels, the peak opening
887 of the laryngeal gesture is timed to occur around the release of the supralaryngeal con-
888 striction, resulting in long-lag VOT. When a voiceless stop precedes a devoiced vowel,
889 in contrast, the laryngeal gesture of the voiceless stop temporally aligns with the vowel
890 midpoint and increases in magnitude substantially. In devoiced vowels, the laryngeal ab-
891 duction is greater than two times the magnitude of a voiceless stop preceding a voiced
892 vowel. The shift in the timing and magnitude of the laryngeal gesture indicates a gesture
893 reorganization that facilitates devoicing.

894 In contrast to voiceless stops, which show substantial temporal variation between la-
895 rylaryngeal and supra-laryngeal gestures, both in Japanese and in the world's languages, the
896 laryngeal and supra-laryngeal gestures of fricatives cannot be temporally displaced so eas-
897 ily. This has consequences for the kinematics of devoicing. In fricative environments,
898 devoicing is not achieved by adjusting the timing or magnitude of the glottal opening, at
899 least not in Tokyo Japanese. Instead, the timing and magnitude of the laryngeal gestures
900 for fricatives is similar when preceding both voiced and devoiced vowels ([Fujimoto, 2015](#)).
901 This means that devoicing following fricatives is achieved in some other way.

902 As an acoustic description, high vowel devoicing following fricative environments can
903 be characterized as a prolonging of the aperiodic energy of a fricative so that it extends
904 across the lingual articulation of the vowel. Articulatorily, maintenance of turbulent air-
905 flow is facilitated by narrowing the vocal tract. A key difference between [ʃ] and [f] is
906 vocal tract width, i.e. cross-sectional area. Since [ʃ] has a constriction in the vocal tract,
907 it naturally conditions a narrow channel that facilitates prolonged turbulent airflow. This
908 is not the case for [f]. Since there is no oral constriction, it is naturally more difficult to
909 sustain turbulent airflow. Specifically, the amount of airflow needed to generate turbulence
910 is a function of the width of the channel, so narrowing the channel means that turbulence
911 can be achieved with less airflow.

912 Raising the tongue dorsum for [u] narrows the vocal tract and therefore facilitates
913 devoicing, when devoicing is generated from the prolongation of aperiodic energy. Such

914 facilitation is likely more helpful in the environment following [F] than in the environment
915 following [ʃ], since [ʃ] already has a lingual constriction. Speakers may be less likely to
916 delete [u] following [F] (compared with [ʃ]) because deletion actually makes it harder to
917 maintain devoicing.

918 6 Conclusion

919 Despite extensive past research on high vowel devoicing in Japanese, one issue that has
920 remained open is whether the devoiced vowels are phonologically deleted or not. Follow-
921 ing a previous study on this topic (Shaw & Kawahara 2018b), the current EMA-based
922 experiment explored this question with an extended stimulus set, with a new hypothesis
923 that surrounding consonantal environments may modulate the deletion probability. The
924 current experiment replicated the core findings of Shaw & Kawahara (2018b); there was
925 a bimodal distribution in deletion probabilities for devoiced [y], with one mode represent-
926 ing vowels that fully retained their articulatory target and the other representing a linear
927 tongue dorsum trajectory between flanking vowels.

928 The lack of a tongue dorsum height target, if it is due to vowel deletion, will presum-
929 ably have phonological consequences for the language, including, at least, syllabification
930 and syllable-based phonological patterns, e.g. accent placement and truncation patterns
931 (as reviewed in the introduction). However, such evidence has not yet been identified.
932 This could be for a number of reasons. The vowel may be retained, even if it lacks a
933 tongue dorsum height target, affecting the phonetics in other dimensions. Possibilities
934 include duration, lip movements, and the relative timing of flanking gestures. Alterna-
935 tively, the vowel may be deleted while higher level prosodic structure, including moras
936 and syllables, are retained, a possibility explored in Kawahara & Shaw (2018).

937 Additionally, the current experiment found an effect of fricative place of articulation
938 on deletion probability—more deletion following [F] than [ʃ]—and individual differences
939 in sensitivity to surrounding consonantal environments. These results are of descriptive
940 importance, as we still know very little about the factors that condition variable phonolog-
941 ical deletion of devoiced vowels in Japanese or, for that matter, any other languages that
942 exhibit vowel devoicing. The current results highlight the importance of examining such

943 behavior both within- and across- speakers, as sensitivity to phonological factors may also
944 vary within a speech community.

945 **Statement of Ethics**

946 The current experiment was conducted with the approval of Western Sydney University
947 and Keio University (Protocol number: HREC 9482). A consent form was obtained from
948 each participant before the experiment.

949 **Conflict of Interest**

950 The authors declare no conflicts of interest.

951 **Author Contributions Statement**

952 Designing the experiment: JS and SK; data analysis: JS; discussion of the results: JS and
953 SK; writing up the paper: JS and SK.

954 **Acknowledgments**

955 Thanks to Chika Takahashi and Jeff Moore for help with the experiment, to Emily Grabowski
956 for coding the data for devoicing, to Noah Macey for parsing out gestures. Thanks also
957 go to Emily for work organizing the Matlab code for analysis into functions and a wrap-
958 per script, which are available on request. Portions of this study were presented at AMP
959 2018. This project is supported by JSPS grant #15F15715 to both authors. Thanks to two
960 anonymous reviewers for comments.

961 **References**

- 962 Beckman, Mary. 1982. Segmental duration and the ‘mora’ in Japanese. *Phonetica* 39.
963 113–135.
964 Beckman, Mary. 1986. *Stress and non-stress accent*. Dordrecht: Foris.

- 965 Beckman, Mary & Atsuko Shoji. 1984. Spectral and perceptual evidence for CV coarticulation in devoiced /si/ and /syu/ in Japanese. *Phonetica* 41. 61–71.
- 966
- 967 Bennett, Ryan. 2020. Vowel deletion as phonologically-condition gestural overlap in Usanteko. Talk presented at Keio-ICU LINC.
- 968
- 969 Berent, Iris, Tracy Lennertz, Jongho Jun, Miguel A. Moreno & Paul Smolensky. 2008. Language universals in human brains. *Proceedings of the National Academy of Sciences* 5321–5325.
- 970
- 971
- 972 Berent, Iris, Donca Steriade, Tracy Lennertz & Vered Vaknin. 2007. What we know about what we have never heard: Evidence from perceptual illusions. *Cognition* 104. 591–630.
- 973
- 974
- 975 Brickhouse, C.J. & Kate Lindsey. 2020. Investigating the phonetics-phonology interface with field data: Assessing phonological specification through acoustic trajectories. Poster presented at the 96th meeting of the Linguistics Society of America.
- 976
- 977
- 978 Browman, Catherine & Louis Goldstein. 1986. Towards an articulatory phonology. *Phonology Yearbook* 3. 219–252.
- 979
- 980 Browman, Catherine & Louis Goldstein. 1989. Articulatory gestures as phonological units. *Phonology* 6. 201–251.
- 981
- 982 Browman, Catherine & Louis Goldstein. 1992a. Articulatory phonology: An overview. *Phonetica* 49. 155–180.
- 983
- 984 Browman, Catherine & Louis Goldstein. 1992b. "targetless" schwa: An articulatory analysis. In G. Docherty & R. Ladd (eds.), *Papers in laboratory phonology II: Gesture, segment, prosody*, 26–56. Cambridge: Cambridge University Press.
- 985
- 986
- 987 Cho, Taehong. 2016. Prosodic boundary strengthening in the phonetics-prosody interface. *Language and Linguistic Compass* 10(3). 120–141.
- 988
- 989 Choi, John. 1995. An acoustic-phonetic underspecification account of marshallese vowel allophony. *Journal of Phonetics* 23. 323–347.
- 990
- 991 Cohn, Abigail. 1993. Nasalisation in English: Phonology or phonetics. *Phonology* 10. 43–81.
- 992
- 993 Cohn, Abigail. 2006. Is there gradient phonology? In Gisbert Fanselow, Caroline Fery, Matthias Schlesewsky & Ralf Vogel (eds.), *Gradience in grammar: Generative perspectives*, 25–44. Oxford: Oxford University Press.
- 994
- 995
- 996 Cox, F & Janet Fletcher. 2017. *Australian english pronunciation and transcription*. Cambridge: Cambridge University Press.
- 997
- 998 Cutler, Anne, Takashi Otake & James McQueen. 2009. Vowel devoicing and the perception of spoke japanese words. *Journal of the Acoustical Society of America* 1693.
- 999
- 1000 Dauer, Rebecca M. 1980. The reduction of unstressed high vowels in Modern Greek. *Journal of the International Phonetic Association* 10(1-2). 17–27.
- 1001
- 1002 Delforge, A.M. 2008. Gestural alignment constraints and unstressed vowel devoicing in

- 1003 andean spanish. *Proceedings of WCCFL* 26. 147–155.
- 1004 Dupoux, Emmanuel, Kazuhiko Kakehi, Yuki Hirose, Christophe Pallier & Jacques Mehler.
1005 1999. Epenthetic vowels in Japanese: A perceptual illusion? *Journal of Experimental*
1006 *Psychology: Human Perception and Performance* 25. 1568–1578.
- 1007 Dupoux, Emmanuel, Erika Parlato, Sónia Frota, Yuki Hirose & Sharon Peperkamp. 2011.
1008 Where do illusory vowels come from? *Journal of Memory and Language* 64(3). 199–
1009 210.
- 1010 Durvasula, K., H.-H. Huang, S. Uehara, Q. Luo & Y.-H. Lin. 2018. Phonology modu-
1011 lates the illusory vowels in perceptual illusions: evidence from Mandarin & English.
1012 *Laboratory Phonology*.
- 1013 Faber, Alice & Timothy Vance. 2010. More acoustic traces of “deleted” vowels in
1014 Japanese. In Mineharu Nakayama & Carles Quinn (eds.), *Japanese/korean linguistics*
1015 9, 100–113. CSLI.
- 1016 Fais, Laurel, Sachiyo Kajikawa, Shigeaki Amano & Janet F. Werker. 2010. Now you hear
1017 it, now you don’t: Vowel devoicing in japanese infant-directed speech. *Journal of Child*
1018 *Language* 37(2). 319–340.
- 1019 Fujimoto, Masako. 2015. Vowel devoicing. In Haruo Kubozono (ed.), *The handbook of*
1020 *Japanese language and linguistics: Phonetics and phonology* 167-214, Mouton.
- 1021 Fujimoto, Masako, Emi Murano, Seiji Niimi & Shigeru Kiritani. 2002. Difference in
1022 glottal opening pattern between tokyo and osaka dialect speakers: Factors contributing
1023 to vowel devoicing. *Folia Phoniatica et Logopaedica* 54(3). 133–143.
- 1024 Funatsu, Seiya & Masako Fujimoto. 2011. Physiological realization of japanese vowel
1025 devoicing. *Proceedings of Forum Acousticum* 2709–2714.
- 1026 Garcia, D. 2010. Robust smoothing of gridded data in one and higher dimensions with
1027 missing values. *Computational Statistics & Data Analysis* 54(4). 1167–1178.
- 1028 Gu, Chong. 2002. Smoothing spline anova models. R package.
- 1029 Guenther, F. H., C. Y. Espy-Wilson, S. E. Boyce, M. L. Matthies, M. Zandipour & J. S.
1030 Perkell. 1999. Articulatory tradeoffs reduce acoustic variability during American En-
1031 glish /r/ production. *Journal of the Acoustical Society of America* 105. 2854–2865.
- 1032 Hall, Kathleen Currie, Elizabeth Hume, Florian T. Jaeger & Andrew Wedel. 2018. The
1033 role of predictability in shaping phonological patterns. *Linguistics Vanguard* 4(S2).
- 1034 Haraguchi, Shosuke. 1984. Some tonal and segmental effects of vowel height in Japanese.
1035 In M. Aronoff & R. T. Oehrle (eds.), *Language sound structure: Studies in phonology*
1036 *presented to morris halle by his teacher and students*, 145–156. Cambridge: MIT Press.
- 1037 Hirayama, Manami. 2009. *Postlexical prosodic structure and vowel devoicing in Japanese*.
1038 University of Toronto Doctoral dissertation.
- 1039 Imai, Terumi. 2004. *Vowel devoicing in Tokyo Japanese: A variationist approach*. Michi-
1040 gan State University Doctoral dissertation.

- 1041 Imaizumi, Satoshi & A. Hayashi. 1995. Listener-adaptive adjustments in speech produc-
 1042 tion: Evidence from vowel devoicing. *Annual Bulletin Research Institute of Logopedics*
 1043 *and Phoniatics*, 29. 43–48.
- 1044 Imaizumi, Satoshi, A. Hayashi & T. Deguchi. 1995. Listener adaptive characteristics of
 1045 vowel devoicing in Japanese dialogue. *Journal of the Acoustical Society of America* 98.
 1046 768–778.
- 1047 Ishihara, Shinichiro. 2011. Japanese focus prosody revisited: Freeing focus from prosodic
 1048 phrasing. *Lingua* 121(13). 1870–1889.
- 1049 Iskarous, Khalil, Joyce McDonough & Douglas H Whalen. 2012. A gestural account of
 1050 the velar fricative in Navajo. *Laboratory Phonology* 3(1). 195–210.
- 1051 Isomura, Kazuhiro. 2009. *Onsei-wo oshieru [Teaching Japanese phonetics]*. Tokyo: Hit-
 1052 suji Shobo.
- 1053 Ito, Junko & Armin Mester. 1995. Japanese phonology. In John Goldsmith (ed.), *The*
 1054 *handbook of phonological theory*, 817–838. Oxford: Blackwell.
- 1055 Jaeger, Florian T. & Esteban Buz. 2018. Signal reduction and linguistic encoding. In
 1056 Eva M. Fernández & Cairns Helen Smith (eds.), *The handbook of psycholinguistics*,
 1057 38–81. John Wiley & Sons.
- 1058 Jain, Anil K. 1989. *Fundamentals of digital image processing*. Englewood Cliffs: Prentice
 1059 Hall.
- 1060 Jannedy, Stephanie. 1995. Gestural phasing as an explanation for vowel devoicing in
 1061 Turkish. *Ohio State University Working Papers in Linguistics* 45. 56–84.
- 1062 Johnson, Keith, Peter Ladefoged & Mona Lindau. 1993. Individual differences in vowel
 1063 production. *Journal of the Acoustical Society of America* 94(2). 701–714.
- 1064 Jun, Sun-Ah & Mary Beckman. 1993. A gestural-overlap analysis of vowel devoicing
 1065 in Japanese and Korean. Paper presented at the 67th annual meeting of the Linguistic
 1066 Society of America, Los Angeles.
- 1067 Jun, Sun-Ah, Mary Beckman & H.-J. Lee. 1998. Fiberscopic evidence for the influence
 1068 on vowel devoicing of the glottal configurations for Korean obstruents. *UCLA Working*
 1069 *Papers in Phonetics* 96. 43–68.
- 1070 Kaneko, Ikuyo & Shigeto Kawahara. 2002. Positional faithfulness theory and the emer-
 1071 gence of the unmarked: The case of Kagoshima Japanese. *ICU English Studies* 11 5.
 1072 18–36.
- 1073 Kawahara, Shigeto. 2015. A catalogue of phonological opacity in Japanese. *REPORTS of*
 1074 *the Keio Institute of Cultural and Linguistic Studies* 46. 145–174.
- 1075 Kawahara, Shigeto & Jason Shaw. 2018. Persistency of prosody. *Hana-bana: A Festschrift*
 1076 *for Junko Ito and Armin Mester*. .
- 1077 Kawahara, Shigeto, Jason A. Shaw & Shinichiro Ishihara. to appear. Assessing the
 1078 prosodic licensing of wh-in-situ in Japanese: A computational-experimental approach.

- 1079 *Natural Language and Linguistic Theory*.
- 1080 Kawakami, Shin. 1977. *Nihongo onsei gaisetsu [an overview of Japanese phonetics]*.
 1081 Tokyo: Ohuusha.
- 1082 Keating, Patricia A. 1988. Underspecification in phonetics. *Phonology* 5. 275–292.
- 1083 Kibe, Nobuko. 2001. Sound changes in Kagoshima dialect. *Journal of the Phonetic Society*
 1084 *of Japan* 5. 42–48.
- 1085 Kilbourn-Ceron, Oriana & Morgan Sonderegger. 2018. Boundary phenomena and vari-
 1086 ability in Japanese high vowel devoicing. *Natural Language and Linguistic Theory*
 1087 36(1). 175–217.
- 1088 Kilpatrick, Alexander, Shigeto Kawahara, Rikke Bungaard-Nielsen, Brett Baker & Janet
 1089 Fletcher. 2020. Japanese perceptual epenthesis is modulated by transitional probability.
 1090 *Language and Speech*.
- 1091 Kondo, Mariko. 1997. *Mechanisms of vowel devoicing in Japanese*: University of Edin-
 1092 burgh Doctoral dissertation.
- 1093 Kondo, Mariko. 2001. Vowel devoicing and syllable structure in Japanese. In
 1094 *Japanese/korean linguistics*, CSLI.
- 1095 Kuriyagawa, F & Masayuki Sawashima. 1989. Word accent, devoicing and duration of
 1096 vowels in Japanese. *Annual Bulletin of the Research Institute of Language Processing*
 1097 23. 85–108.
- 1098 Lindblom, Björn. 1963. Spectrographic study of vowel reduction. *Journal of the Acousti-*
 1099 *cal Society of America* 35. 1773–1781.
- 1100 Lindblom, Björn. 1990. Explaining phonetic variation: A sketch of the H&H theory. In
 1101 W. J. Hardcastle & A. Marchal (eds.), *Speech production and speech modeling*, 403–
 1102 439. Dordrecht: Kluwer.
- 1103 Maekawa, Kikuo. 1990. Production and perception of the accent in the consecutively
 1104 devoiced syllables in Tokyo Japanese. *Proceedings of ICSLP 1990*.
- 1105 Maekawa, Kikuo & H. Kikuchi. 2005. Corpus-based analysis of vowel devoicing in spon-
 1106 taneous Japanese: An interim report. In J. van de Weijer, Kensuke Nanjo & Tetsuo
 1107 Nishihara (eds.), *Voicing in Japanese*, 205–228. Berlin: de Gruyter.
- 1108 Martin, Andrew, Akira Utsugi & Reiko Mazuka. 2014. The multidimensional nature of
 1109 hyperspeech: Evidence from Japanese vowel devoicing. *Cognition* 132(2). 216–228.
- 1110 Matsui, Michinao. 2017. On the input information of the C/D model for vowel devoicing
 1111 in Japanese. *Journal of the Phonetic Society of Japan* 21(1). 127–140.
- 1112 McCarthy, John J. 2008. The gradual path to cluster simplification. *Phonology* 25(2).
 1113 271–319.
- 1114 Moon, S. J. & Björn Lindblom. 1994. Interaction between duration, context and speaking
 1115 style in English stressed vowels. *Journal of Acoustical Society of America* 96.1. 40–55.
- 1116 Mücke, Doris, Martine Grice & Taehong Cho. 2014. More than a magic moment—Paving

- 1117 the way for dynamics of articulation and prosodic structure. *Journal of Phonetics* 44.
1118 1–7.
- 1119 Munson, Benjamin, Jan Edwards, S.K. Shellinger, Mary E. Beckman & M.K. Meyer.
1120 2010. Deconstructing phonetic transcription: Covert contrast, perceptual bias, and an
1121 extraterrestrial view of Vox Humana. *Clinical linguistics & phonetics* 24(4-5). 245–260.
- 1122 Murray, Robert & Theo Vennemann. 1983. Sound change and syllable structure: Problems
1123 in Germanic phonology. *Language* 59. 514–28.
- 1124 Murray, Robert W. 1988. *Phonological strength and early Germanic syllable structure*.
1125 München: Wilhelm Fink Verlag.
- 1126 Myers, Scott. 1998. Surface underspecification of tone in chichewa. *Phonology* 15. 367–
1127 391.
- 1128 Nakamura, Mitsuhiro. 2003. The articulation of vowel devoicing: A preliminary analysis.
1129 *On-in Kenkyuu [Phonological Studies]* 6. 49–58.
- 1130 Nam, H., V. Mitra, M. Tiede, M. Hasegawa-Johnson, C. Espy-Wilson, E. Saltzman &
1131 L. Goldstein. 2012. A procedure for estimating gestural scores from speech acoustics.
1132 *Journal of the Acoustical Society of America* 132(6). 3980–3989.
- 1133 Nam, Hosung, Louis Goldstein, Elliot Saltzman & Dani Byrd. 2004. TADA: An enhanced,
1134 portable Task Dynamics model in MATLAB. *The Journal of the Acoustical Society of*
1135 *America* 115(5). 2430–2430.
- 1136 Nielsen, Kuniko. 2015. Continuous versus categorical aspects of Japanese consecutive
1137 devoicing. *Journal of Phonetics* 52. 70–88.
- 1138 Nogita, A., N. Yamane & S. Bird. 2013. The Japanese unrounded back vowel [W] is in
1139 fact rounded central/front [ɪ/Y]. Paper presented at the Ultrafest VI, Edinburgh.
- 1140 Ogasawara, Naomi. 2013. Lexical representation of Japanese vowel devoicing. *Language*
1141 *and Speech* 56(1). 5–22.
- 1142 Perkell, J. S., M. L. Matthies, M. A. Svirsky & M. I. Jordan. 1993. Trading relations
1143 between tongue body raising and lip rounding in production of the vowel /u/: A pilot
1144 motor equivalence study. *JASA* 93. 2948–2961.
- 1145 Pierrehumbert, Janet B. 1980. *The phonetics and phonology of English intonation*. MIT
1146 Doctoral dissertation.
- 1147 Pierrehumbert, Janet B. & Mary Beckman. 1988. *Japanese tone structure*. Cambridge:
1148 MIT Press.
- 1149 Poser, William. 1990. Evidence for foot structure in Japanese. *Language* 66. 78–105.
- 1150 Rhodes, Richard. 1972. Cheyenne vowel devoicing and transderivational constraints. *Work*
1151 *Papers of the Summer Institute of Linguistics, University of North Dakota Session* 16.
1152 52–55.
- 1153 Roettger, Timo B. 2019. Researcher degree of freedom in phonetic research. *Laboratory*
1154 *Phonology* 10. <http://doi.org/10.5334/labphon.147>

- 1155 Saltzman, E.L. & Kevin. G. Munhall. 1989. A dynamical approach to gestural patterning
1156 in speech production. *Ecological psychology* 1(4). 333–382.
- 1157 Sawashima, Masayuki. 1971. Devoicing of vowels. *Annual Bulletin of Research Institute*
1158 *of Logopedics and Phoniatics* 5. 7–13.
- 1159 Shaw, Jason & Shigeto Kawahara. 2018a. Assessing surface phonological specification
1160 through simulation and classification of phonetic trajectories. *Phonology* 35. 481–522.
- 1161 Shaw, Jason & Shigeto Kawahara. 2018b. The lingual gesture of devoiced [u] in Japanese.
1162 *Journal of Phonetics* 66. 100–119.
- 1163 Shaw, Jason & Shigeto Kawahara. 2019. Effects of surprisal and entropy on vowel duration
1164 in Japanese. *Language and Speech* 62(1). 80–114.
- 1165 Sjoberg, A.F. 1963. *Uzbek structural grammar*. Indiana University.
- 1166 Smith, Caroline L. 2003. Vowel devoicing in contemporary French. *Journal of French*
1167 *Language Studies* 13(2). 177–194.
- 1168 Starr, Rebecca L & Stephanie S. Shih. 2017. The syllable as a prosodic unit in Japanese
1169 lexical strata: Evidence from text-setting. *Glossa*.
- 1170 Sugito, Miyoko & Hajime Hirose. 1988. Production and perception of accented devoiced
1171 vowels in Japanese. *Annual Bulletin of Research Institute of Logopedics and Phoniatics*
1172 22. 19–36.
- 1173 Tanner, James, Morgan Sonderegger & Francisco Torreira. 2019. Durational evidence that
1174 Tokyo Japanese vowel devoicing is not gradient reduction. *Frontiers in Psychology*.
- 1175 Tiede, Mark. 2005. Mview. Software.
- 1176 Tsuchida, Ayako. 1997. *Phonetics and phonology of Japanese vowel devoicing*. Cornell
1177 University Doctoral dissertation.
- 1178 Vance, Timothy. 1987. *An introduction to Japanese phonology*. New York: SUNY Press.
- 1179 Vance, Timothy. 2008. *The sounds of Japanese*. Cambridge: Cambridge University Press.
- 1180 Vatikiotis-Bateson, Eric, Adriano Vilela Barbosa & Catherine T. Best. 2014. Articulatory
1181 coordination of two vocal tracts. *Journal of Phonetics* 44. 167–181.
- 1182 Vennemann, Theo. 1988. *Preference laws for syllable structure and the explanation of*
1183 *sound change: With special reference to German, Germanic, Italian, and Latin*. Berlin:
1184 Mouton de Gruyter.
- 1185 Vogel, Rachel. 2021. A unified account of two vowel devoicing phenomena: the case of
1186 Cheyenne. *Proceedings of Annual Meeting of Phonology*.
- 1187 Watson, Catherine I. & Jonathan Harrington. 1999. Acoustic evidence for dynamic for-
1188 mant trajectories in Australian English vowels. *Journal of the Acoustical Society of*
1189 *America* 106(458-468).
- 1190 Whang, James. 2018. Recoverability-driven coarticulation: Acoustic evidence from
1191 Japanese high vowel devoicing. *Journal of the Acoustical Society of America* 143.
1192 1159–1172.

- 1193 Whang, James. 2019. Effects of phonotactic predictability on sensitivity to phonetic detail.
1194 *Laboratory Phonology* 10(1).
- 1195 Whang, James, Jason Shaw & Shigeto Kawahara. 2020. Acoustic consequences of vowel
1196 deletion in devoicing environments. Talk presented at LabPhon 17.
- 1197 Yoshioka, H. 1981. Laryngeal adjustments in the production of the fricative consonants
1198 and devoiced vowels in Japanese. *Phonetica* 38. 236–251.
- 1199 Zhang, M., C. Geissler & Jason Shaw. 2019. Gestural representations of tone in Mandarin:
1200 Evidence from timing alternations. *Proceedings of the 19th International Congress of*
1201 *Phonetic Sciences* 1803–1807.