

1 **The lingual articulation of devoiced /u/ in Tokyo Japanese**

2

3 Jason Shaw* (Yale University)

4 Shigeto Kawahara (Keio University)

5

6 *Corresponding author: jason.shaw@yale.edu

7 New Haven, CT 06520, USA

8

9

Abstract

10 In Tokyo Japanese, /u/ is typically devoiced between two voiceless consonants. Whether the
11 lingual vowel gesture is influenced by devoicing or present at all in devoiced vowels remains
12 an open debate, largely because relevant articulatory data has not been available. We report
13 ElectroMagnetic Articulography (EMA) data that addresses this question. We analyzed both
14 the trajectory of the tongue dorsum across VC₁uC₂V sequences as well as the timing of C₁ and
15 C₂. These analyses provide converging evidence that /u/ in devoicing contexts is optionally
16 targetless—the lingual gesture is either categorically present or absent but seldom reduced.
17 When present, the magnitude of the lingual gesture in devoiced /u/ is comparable to voiced
18 vowel counterparts. Although all speakers produced words with and without a vowel height
19 target for /u/, the frequency of targetlessness varied across speakers and items. The timing
20 between C₁ and C₂, the consonants flanking /u/ was also effected by devoicing but to varying
21 degrees across items. The items with the greatest effect of devoicing on this inter-consonantal
22 interval were also the items with the highest frequency of vowel height targetlessness for
23 devoiced /u/.

1 **Keyword:** Japanese, devoicing, /u/, articulatory phonetics, EMA, phonetic interpolation,
2 Discrete Cosine Transform (DCT), Bayesian classifier, gestural coordination

3

4 **General background**

5 This paper examines the lingual articulation of devoiced /u/ in Tokyo Japanese. A classic
6 description of the devoicing phenomenon is that high vowels are devoiced between two
7 voiceless consonants and after a voiceless consonant before a pause (Fujimoto, 2015; Kondo,
8 1997, 2005; Tsuchida, 1997 among many others). This sort of description, high vowel
9 devoicing in a particular context, applies to vowels in numerous other languages including e.g.,
10 French (Cedergren & Simoneau, 1985; Smith, 2003), Greek (Dauer, 1980; Eftychiou, 2010)
11 Korean (Jun, Beckman, & Lee, 1998) and Uzbek (Sjoberg, 1963), but Tokyo Japanese is
12 arguably the best studied case of vowel devoicing.

13

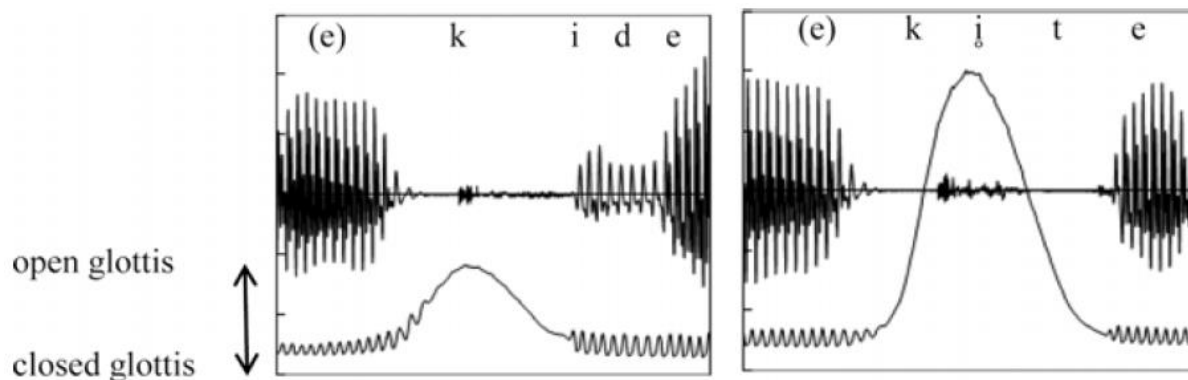
14 There is a large body of work on this phenomenon in Japanese, covering its phonological
15 conditions (e.g., Kondo, 2005; Tsuchida, 1997), its interaction with other phonological
16 phenomena like pitch accent (e.g., Kuriyagawa & Sawashima, 1989; Maekawa, 1990; Maekawa
17 & Kikuchi, 2005; Vance, 1987) and prosodic structure (Kilbourn-Ceron & Sonderegger, 2017),
18 its acoustic and perceptual characteristics (Beckman & Shoji, 1984; Faber & Vance, 2000;
19 Matsui, 2014; Nielsen, 2015; Sugito & Hirose, 1988), and studies of the vocal folds (Fujimoto,
20 Murano, Niimi, & Kiritani, 2002; Hirose, 1971; Sawashima, 1971; Tsuchida, 1997). Fujimoto
21 (2015) provides a recent, comprehensive overview of this research. While now we have a good
22 understanding of many aspects of high vowel devoicing in Tokyo Japanese, there is little data
23 available on the lingual gestures of high vowels when they are devoiced. The only study that
24 we are aware of is Funatsu & Fujimoto (2011), which used EMMA (ElectroMagnetic

1 Midsagittal Articulography) with concurrent imaging of the vocal fold using nasal endoscopy.
2 They found little difference between devoiced and voiced /i/ in terms of lingual articulation.
3 However, this experiment used only one speaker and one item pair (/kide/ vs. /kite/). The study
4 included four repetitions of each item, and offered no quantitative analyses of the data. Our
5 study is intended to expand on this previous work by reporting more data from more speakers
6 and more extensive quantitative analysis.

7

8 Why is it important to study the lingual gestures of devoiced vowels? There are a few lines of
9 motivation behind the current study. First, consider Figure 1, taken from Fujimoto et al.'s (2002)
10 study, which used nasal endoscopy to image the glottal gestures of high vowel devoicing in
11 Japanese.

12



13

14 **Figure 1: The degrees of glottal abduction in Japanese. The left panel: a voiceless stop**
15 **followed by a voiced stop, which has a single abduction gesture for /k/. The right panel:**
16 **a voiceless stop /k/ followed by a voiceless vowel and another voiceless stop /t/, which**
17 **also has a single abduction gesture. The magnitude of the abduction gesture in the right**
18 **panel is larger than twice the size of the abduction gesture in the left panel. Taken from**
19 **Fujimoto et al. (2002), cited and discussed in Fujimoto (2015).**

20

1 Figure 1 shows that a Japanese devoiced vowel has a single laryngeal gesture of greater
2 magnitude than a single consonant gesture, or even the sum of two voiceless consonant gestures
3 (c.f., Munhall & Lofqvist, 1992 for English which shows the latter pattern).¹ This observation
4 implies that Japanese devoiced vowels involve active laryngeal abduction, not simply overlap
5 of two surrounding gestures. This conclusion in turn implies that Japanese speakers exert *active*
6 laryngeal control over devoiced high vowels (cf. Jun & Beckman 1993 for an analysis that relies
7 on passive gestural overlap, to be discussed below). To the extent that Japanese speakers
8 actively control the laryngeal gesture for devoiced vowels, are lingual gestures of devoiced
9 vowels also actively controlled? There are competing views on this matter. On the one hand, it
10 seems logical that active control of a non-contrastive property (allophonic devoicing) would
11 imply active control of a contrastive property (tongue position in the vocal tract). On the other
12 hand, the way that devoicing operates physiologically in Japanese obliterates much of the
13 acoustic signature of lingual articulation. Speakers may not exert active control over aspects of
14 articulation that do not have salient auditory consequences.

15

16 The second line of motivation for the current study is the question of whether “devoiced”
17 vowels are simply devoiced or deleted. This issue has been discussed extensively in previous
18 studies of Japanese high vowel devoicing. Kawakami (1977: 24-26) argues that vowels delete
19 in some environment and devoice in others, but he offers no phonological or phonetic evidence.
20 Vance (1987) raised and rejected the hypothesis that high vowels in devoicing contexts are

¹ Munhall and Lofqvist (1992) investigate the timing and magnitude of laryngeal gestures in the consonants /s/ and /t/ in the sequences *kiss_s ted* spoken at different speech rates. At slow speech rates two distinct laryngeal gestures can be identified but at faster speech rates the gestures merge into one laryngeal gesture approximating the sum of the two smaller consonantal gestures.

1 deleted. Kondo (2001) argues that high vowel devoicing is actually deletion based on a
2 phonological consideration. Devoicing in consecutive syllables is often prohibited (although
3 there is much variability: Nielsen 2015), and Kondo argues that this prohibition stems from a
4 prohibition against complex onset or complex coda (i.e. *CCC). On the other hand, Tsuchida
5 (1997) and Kawahara (2015) argue that bimoraic foot-based truncation (Poser, 1990) counts a
6 voiceless vowel as one mora (e.g. [sɯto] from [sɯtoraiki] ‘strike’, *[stora]).² If /u/ was
7 completely deleted losing its mora, the bimoraic truncation should result in *[stora]. Hirayama
8 (2009) makes a similar phonological argument by showing that devoiced vowels' moras are just
9 as relevant for Japanese *haiku* poetry as moras in voiced vowels. However, just because moras
10 for the devoiced vowels remain does not necessarily mean that the vowel is present. The
11 adjacent consonant could conceivably host the mora and syllable—this hypothesis is actually
12 proposed by Matsui (2014), who argues that Japanese has consonantal syllables in this
13 environment (see Dell and Elmedlaoui, 2002 for similar analyses of Tashlhiyt Berber and
14 Moroccan Arabic). Thus, evidence for either deletion or devoicing from a phonological
15 perspective is mixed (see Fujimoto 2015: 197-198 for other studies addressing this debate).³

16
17 Previous acoustic studies show that on spectrograms, vowels leave no trace of lingual
18 articulation except for coarticulation on surrounding consonants, which lead them to conclude
19 that vowels are deleted (Beckman 1982; Beckman & Shoji 1984; Whang 2014). An anonymous

² Here and throughout we use the symbol [u] to refer to a broad phonetic transcription. The actual realizations of this vowel in our data in Tokyo Japanese more generally tend not to be as back or as rounded as [u] is strictly defined in the IPA. See Vance (2008:51) for a detailed description of Japanese /u/. We return to this point when discussing our specific hypotheses below.

³ Tsuchida (1997) argues that there is “phonological devoicing” as well as “phonetic devoicing”.

1 reviewer questions this finding reported in past work. In principle, a change in the sound source,
2 from modal voicing to turbulence, is independent of the resonance properties of the vocal tract
3 (e.g., Stevens, 1999). We might therefore expect to be able to identify formant structure in the
4 aperiodic energy characteristic of devoiced vowels. The reported absence of such structure in
5 past studies may follow from the particular location of turbulent energy sources excited in the
6 vocal tract preceding devoiced vowels in Japanese and their perseverative influence on
7 devoiced vowels. Most of the voiceless consonants preceding devoiced vowels in Japanese are
8 fricatives or affricates that involve turbulence generated by a narrow channel of air in the
9 anterior portion of the vocal tract.⁴ The formants produced by these consonants are resonances
10 of the cavity in front of the noise source, i.e., the front cavity. Using ElectroPalatoGraphy
11 (EPG), Matsui (2014) shows that the narrow channel characteristic of fricatives and affricates
12 persists across following devoiced vowels. The perseveration of the airflow channel across
13 vowels no doubt helps to sustain devoicing by maintaining high intraoral air pressure, but may
14 also contribute to the obliteration of spectral cues to vowel articulation. In contrast, when the
15 energy source is modal phonation at the glottis, higher formants (above F1) of vowels result
16 from (coupled) resonances of both the front and back cavities, which provide very different
17 acoustic signatures from resonance of the front cavity alone. These factors may contribute to
18 the claim that devoiced vowels show no acoustic traces of lingual articulation. Beckman (1982:
19 118, footnote 3) states that "deletion" is a better term physically, because "there is generally no
20 spectral evidence for a voiceless vowel", whereas "devoicing" is a better term psychologically,

⁴ Japanese lacks singleton /p/, except in some recent loanwords, and /t/ is affricated before high vowels, so the only stop consonant conditioning devoicing that does not involve turbulence generated in the anterior portion of the vocal tract is /k/.

1 because Japanese speakers hear a voiceless vowel even in the absence of spectral evidence (c.f.,
2 Dupoux, Kakehi, Hirose, Pallier, & Mehler, 1999). Beckman and Shoji (1984: 64) likewise
3 state that “[w]hen the waveform of a devoiced syllable is examined, however, neither its
4 spectral nor its temporal structure indicates the presence of a voiceless vowel.” These
5 statements embrace the "deletion" view, at least at the speech production level. Even if vowel
6 devoicing involves phonological deletion, or deletion of some component (feature, gesture) of
7 the vowel, it could be the case that its application is optional or variable, influenced by various
8 linguistic and sociological factors (Fujimoto, 2015; Imaizumi & Hayashi, 1995; Nielsen, 2015).

9
10 Not all phonetic studies have embraced the deletion view, however. The clearest instantiation
11 of an alternative view is the "gestural overlap theory" of high vowel devoicing (Faber & Vance,
12 2000; Jun & Beckman, 1993; Jun et al., 1998). In this theory, high vowel devoicing occurs
13 when glottal abduction gestures of the surrounding consonants overlap with the vowel's glottal
14 gesture (though cf. Figure 1). In this sense, the high vowel devoicing processes in Japanese (and
15 other languages like Korean) are "not...phonological rules, but the result of extreme overlap
16 and hiding of the vowel's glottal gesture by the consonant's gesture" (Jun & Beckman 1993:
17 p.4). This theory implies that there is actually no deletion—oral gestures remain the same, but
18 do not leave their acoustic traces because of devoicing.

19
20 To summarize, there is an active debate about whether the lingual gestures of “devoiced”
21 vowels in Japanese are present but inaudible due to devoicing or absent altogether, possibly
22 because of phonological deletion. Vance (2008), the most recent and comprehensive phonetic
23 textbook on Japanese, states that this issue is not yet settled. Studying lingual movements of
24 devoiced vowels will provide crucial new evidence. Recall that the only past study on this topic,
25 Funatsu & Fujimoto (2011), is based on a small number of tokens and one speaker. In this

1 paper, we expand the empirical base, reporting more repetitions (10-15) of ten real words
2 produced by six naive speakers, and we deploy rigorous quantitative methods of analysis, as
3 detailed by Shaw and Kawahara (submitted-a). While Shaw and Kawahara (submitted-a)
4 focused on motivating the computational methodology, this paper reports more data and
5 examines the consequences of vowel devoicing for the temporal organization of gestures. In
6 particular, we examine C-C timing across devoiced vowels in order to ascertain whether
7 devoicing impacts the gestural coordination of flanking consonants.

8 **Hypotheses**

9 Building on the previous studies reviewed in this section, we entertain four specific hypotheses
10 about the lingual articulation of devoiced vowels, stated in (1)

11

12 (1) Hypotheses about the status of lingual articulation in devoiced vowels

13 H1: **full lingual targets**—the lingual articulation of devoiced vowels is the same as
14 for voiced counterparts.

15 H2: **reduced lingual targets**—the lingual articulation of devoiced vowels is
16 phonetically reduced relative to voiced counterparts.

17 H3: **targetless**—devoiced vowels have no lingual articulatory target.⁵

18 H4: **optional target**—devoiced vowels are sometimes targetless (deletion is
19 optional, token-by-token).

⁵ We use the term “targetless” rather than “deletion”, because the latter term commits to (1) a surface representation, (2) a process mapping an underlying representation to a surface representation and (3) the identity of the underlying representation. Our experiment is solely about the surface representation, and hence the term “targetless” is better.

1

2 The passive devoicing hypothesis, or the gestural overlap theory (e.g. Jun & Beckman 1993),
3 maintains that there is actually no phonological deletion, and hence would predict that lingual
4 gestures would remain intact (=H1). This is much like Brownman and Goldstein's (1992)
5 argument that apparently deleted [t] in *perfect memory* in English keeps its tongue tip gesture.
6 This is also the conclusion that Funatsu & Fujimoto (2011) reached for devoiced /i/ in their
7 sample of EMMA data. Besides the small sample size mentioned above, another caveat is
8 concurrently collected nasal endoscopy may have promoted stable lingual gestures across
9 contexts, since retraction of the tongue body may trigger a gag-reflex.

10

11 Even if devoiced high vowels are not phonologically deleted, it would not be too surprising if
12 the lingual gestures of high vowels were phonetically reduced, hence H2 in (1). At least in
13 English, more predictable segments tend to be phonetically reduced (Aylett & Turk, 2004,
14 2006; Bell, Brenier, Gregory, Girand, & Jurafsky, 2009; Jurafsky, Bell, Gregory, & Raymond,
15 2001). In Japanese, the segmental identity of devoiced vowels is often highly predictable from
16 context (Beckman & Shoji, 1984; Whang, 2014). Due to devoicing, moreover, the acoustic
17 consequences of a reduced lingual gesture would not be particularly audible to listeners. Hence,
18 from the standpoint of effort-distinctiveness tradeoff (Hall, Hume, Jaeger, & Wedel, 2016;
19 Lindblom, 1990), it would not be surprising to observe reduction of oral gestures in high
20 devoiced vowels.

21

22 The phonological deletion hypothesis, which was proposed by various authors reviewed above,
23 predicts that there should be no lingual targets for devoiced vowels (=H3 in (1)). However, we
24 know that many if not all phonological patterns are variable, i.e., optional (e.g., Coetzee &
25 Pater, 2011), to some degree. There is an increasing body of evidence that phonological and

1 phonetic patterns are stochastic (Hayes & Londe, 2006; McPherson & Hayes, 2016;
2 Pierrehumbert, 2001). For example, Bayles, Kaplan and Kaplan (2016) have shown recently
3 that there is intra-speaker variation in French schwa production such that the same speaker may
4 produce a word with or without a schwa vowel. Even studies using high spatio-temporal
5 resolution articulatory data capable of picking up gradience have revealed that some
6 phonological patterns, e.g., place assimilation, can be optional within a speaker (Ellis &
7 Hardcastle, 2002; Kochetov & Pouplier, 2008). Therefore, we need to consider the possibility
8 that deletion of lingual gestures in devoiced vowels is optional, by assessing the
9 presence/absence of phonetic specification on a token-by-token basis (=H4).

10

11 In contrast to H1 and H2, H3/H4 present some thorny methodological challenges. Even when
12 we limit ourselves to one phonetic dimension, distinguishing categorical absence of phonetic
13 specification from heavy phonetic reduction is challenging. To pursue this challenge, we adopt
14 an assumption in the literature on phonetic underspecification (e.g., Keating, 1988) that the
15 absence of phonetic specification in some dimension results in phonetic interpolation between
16 flanking segments. Making use of this assumption, the essence of our approach is to assess the
17 linearity of the trajectory between flanking vowels, a method introduced in Shaw & Kawahara
18 (submitted-a). To implement this analysis with appropriate baselines on a token-to-token basis,
19 we setup Bayesian classifiers of our devoiced vowel tokens based on two categories of training
20 data: voiced vowels in similar contexts (to our devoiced test items) and linear trajectories
21 between flanking segments. The classifier returns the probability that each token belongs to the
22 voiced vowel trajectory as opposed to the linear interpolation trajectory. This provides us with
23 a rigorous quantitative approach to assessing both degree differences in phonetic reduction and
24 the absence of phonetic specification in a particular phonetic dimension.

25

1 Our analysis focuses on the phonetic dimension that is most characteristic of the target vowel.
2 For the case of /u/, the focus of this study, that dimension is tongue height. Although we follow
3 Vance (2008) and many other authors in the Japanese phonetics literature in using the symbol
4 /u/, the backness and rounding components of this vowel in Japanese are not precisely as the
5 IPA symbol implies. In Japanese, /u/ is rather central, as reported in the ultrasound study of
6 Nogita et al. (2013) and shown in the instructional MRI images of Isomura (2009). Vance
7 (2008: 55) calls Japanese /u/ the hardest of the Japanese vowels to describe, in part because of
8 the labial component, which involves compression instead of rounding. In contrast to /u/ in
9 other languages, e.g., English, where the longitudinal position of EMA sensors placed on the
10 lips gives a reasonable measure of rounding (Blackwood-Ximenes, Shaw, & Carignan, 2017),
11 lip compression associated with /u/ in Japanese is difficult to detect from sensors on the
12 vermilion borders of the lips. Hence, when it comes to differentiating trajectories of voiced
13 vowels from devoiced vowels, the labial and backness components of Japanese /u/ cannot be
14 expected to provide a particularly strong signal against the backdrop of natural variability in
15 vowel production. For this reason, we focus on tongue height and designed stimulus materials
16 in which the target vowel /u/ is always flanked by non-high vowels, e.g., /...eCuCo.../. If the
17 tongue body does not rise from its position for the mid-vowel /e/ to /u/ before falling to /o/, but
18 proceeds instead on a linear trajectory from /e/ to /u/, we would conclude that the token lacks a
19 vowel height target for /u/.

20

21 With respect to the targetless hypotheses (H3/H4), we acknowledge that lacking a phonetic
22 height target for /u/ is not quite the same as vowel deletion, as intended by some of the studies
23 reviewed above. Certainly, vowel deletion implies targetlessness in the height dimension but it
24 also implies that other phonetic dimensions are similarly targetless. It is conceivable that a
25 vowel that lacks a height target for /u/ is nonetheless specified in other phonetic dimensions,

1 e.g., a laryngeal gesture, tongue backness, lip compression, etc., in which case targetless is an
2 inappropriate designation for the vowel. On the other hand, our materials are well-designed to
3 falsify H3/H4 if it can be shown that the movement trajectory of the tongue does not
4 approximate linear interpolation (see Browman & Goldstein, 1992a for a comparable
5 approach). Our study is also capable of falsifying H1 by showing systematic differences
6 between tongue trajectories in voiced and voiceless tokens and H2 by showing either no
7 difference between lingual movements or a bimodal distribution representing presence and
8 absence of a height target. In this way, the limitations in focusing on the vowel height dimension
9 extend to our assessment of all four hypothesis. Our analysis of C1-C2 timing provides
10 additional evidence bearing on H1-H4.

11
12 Besides adjudicating between the hypotheses in (1), another line of motivation for this study
13 stems from the insight that devoiced vowels may provide on how laryngeal and supra-laryngeal
14 gestures are coordinated. As our review of the literature revealed, there is evidence that some
15 aspects of devoiced vowels—the laryngeal gesture and the airflow channel—may be
16 phonetically controlled to ensure devoicing. This suggests that passive devoicing may not be
17 the right synchronic analysis of the Japanese facts, but it is unclear what implications controlled
18 devoicing has for the organization of lingual gestures. The laryngeal gestures that contribute to
19 devoicing ostensibly originate not with the vowel but with the flanking consonants.
20 Nevertheless, EGG evidence in Figure 1 (Fujimoto et al., 2002) indicates that the timing and
21 magnitude of laryngeal gestures shift when vowel devoicing is at stake. When there is an
22 intervening high vowel, the laryngeal gestures of flanking voiceless consonants aggregate to
23 form one large-magnitude laryngeal abduction centered on the vowel. Aggregation of C₁ and
24 C₂ laryngeal gestures across the vowel may also require the oral gestures associated with the
25 consonants to come closer together, if the internal temporal integrity of the consonants, i.e. the

1 articulatory binding of oral and laryngeal gestures (Kingston, 1990), is to be maintained.
2 However, it may be that articulatory binding is violable. In this case, Japanese may sacrifice
3 the internal timing of consonants in order to maintain temporal coordination between supra-
4 laryngeal gestures. On this scenario, supra-laryngeal gestures maintain patterns of temporal
5 organization generally characteristic of VC₁uC₂V sequences, even as laryngeal gestures
6 aggregate to devoice /u/. To assess these possibilities, we also analyzed the timing of the
7 flanking consonants, C₁ and C₂, and the inter-consonantal interval, i.e., the interval from the
8 release of C₁ to the achievement of target of C₂ (specific measurements are defined below).

9
10 The conflict between maintaining consonant-internal timing between laryngeal and oral
11 gestures, on the one hand, and between maintaining inter-gestural timing across consonants and
12 vowels, on the other hand, is complicated by the possibility that the devoiced vowel is absent
13 (H3) at least in some tokens (H4). We assume, following, e.g., Smith (1995), that in a VC₁VC₂V
14 sequence in Japanese, C₂ is coordinated locally to the preceding vowel. If that vowel is absent,
15 yielding C₁C₂V, we expect C-C coordination, whereby C₂ is coordinated with C₁.

16
17 One way to assess the presence of a coordination relation is to evaluate the predicted covariation
18 between temporal intervals (Shaw et al., 2011). Our Japanese materials afford this opportunity.
19 To make the comparison between C-V and C-C timing more concrete, we express the
20 alternatives as coordination topologies (in the sense of Gafos, 2002; Gafos & Goldstein, 2012)
21 in Table 1 along with consequences for relative timing. The rectangles represent activation
22 durations for gestures. The dotted lines indicate variation in intra-gestural activation duration.
23 Under C-V coordination, we assume here that the start of consonant and vowel gestures are
24 coordinated in time, i.e., the gestures are in-phase (Goldstein, Nam, Saltzman, & Chitoran,
25 2009). Under this coordination regime, shortening of C₁ would expose more of the inter-

1 consonantal interval (ICI), predicting a negative correlation between C₁ duration and ICI. Under
 2 C-C coordination, on the other hand, the end of C₁ is coordinated with the start of C₂. Variation
 3 in activation duration for C₁ impacts directly when in time C₂ begins. Hence, under C-C
 4 coordination no trade-off between C₁ duration and ICI is predicted. By assessing the predicted
 5 covariation between C₁ and the inter-consonantal interval, we can adjudicate between C-V and
 6 C-C timing, potentially providing an independent argument for the presence/absence of the
 7 vowel (H3/H4). Evidence for C-V timing is also evidence that the vowel is present to some
 8 degree, whereas C-C timing is indirect evidence that the vowel is absent, or at least that it is not
 9 specified to the extent that it affects the coordination of flanking consonants. This is particularly
 10 useful since our analysis of vowel target presence/absence is limited to the height dimension.
 11 The additional analysis, although indirect, provides another way to assess whether the vowel is
 12 phonetically specified.

13

14 **Table 1: Schematic illustration of C1 duration variation under different coordination**
 15 **topologies.**

	C-V coordination	C-C coordination
Coordination Topology		
Temporal intervals		

16

17 Through the set of analyses described above, we aim to evaluate the lingual articulation of the
 18 devoiced vowel and the consequences of devoicing for the temporal organization of flanking
 19 segments, i.e., the consonants that define the devoicing environment.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Methodology

Speakers

Six native speakers of Tokyo Japanese (3 male) participated. Participants were aged between 19 and 22 years at the time of the study. They were all born in Tokyo, lived there at the time of their participation in the study, and had spent no more than 3 months outside of the Tokyo region. Procedures were explained to participants in Japanese by a research assistant, who was also a native speaker of Tokyo Japanese. All participants were naïve to the purpose of the experiment. They were compensated for their time and local travel expenses. In subsequent discussion we refer to the speakers as S01 through S06, numbered in the order in which they were recorded.

Materials

A total of 10 target words, listed in Table 2, were included in the experiment. These included five words containing /u/ in a devoicing context (second column) and a set of five corresponding words with /u/ in a voiced context (third column). The target /u/ is underlined in each word. Together, these 10 words constitute minimal pairs or near minimal pairs. The word pairs are matched on the consonant that precede /u/. They differ in the voicing specification of the consonant following /u/. The consonant following /u/ is voiceless in devoicing context words (second column) and voiced in voiced context words (third column). This study focused on /u/ and did not consider /i/ for several practical reasons, the most important one being that collecting both /u/ and /i/ tokens would mean reducing the repetitions for each target word and the analytical approach we planned (following Shaw and Kawahara, submitted-a) requires a

1 large number of repetitions per word. We focused on /u/ rather than /i/, because the former is
 2 more likely to be devoiced, as confirmed by the study of high vowel devoicing using the Corpus
 3 of Spontaneous Japanese (Maekawa & Kikuchi 2005; see also Fujimoto et al. 2015 and
 4 references cited therein).

5

6 **Table 2: stimulus items. W and K show the environments in which Kawakami (1971)**

7 **and Whang (2014) predict deletion.**

Comments	Devoicing/deletion	Voiced vowel
V deletion (K, W)	ϕ _u soku 不足 ‘shortage’	ϕ _u zoku 付属 ‘attachment’
V devoicing (K, W)	ʃ _u taiſe: 主体性 ‘subjectivity’	ʃ _u daika 主題歌 ‘theme song’
V deletion (K, W)	kats _u toki 勝つ時 ‘when winning’	kats _u do: 活動 ‘activity’
V devoicing (K) V deletion (W)	hak _u sai 白菜 ‘white cabbage’	jak _u zai 薬剂 ‘medicine’
V deletion (K, W)	mas _u taa マスター ‘master’	mas _u da 益田 ‘Masuda (a surname)’

8

9 The consonant preceding /u/ draws from the set: /s/, /k/, /ʃ/, /ts/, /ϕ/.⁶ All of these consonants
 10 may contribute to the high vowel devoicing environment, when they occur as C₁ in C₁V_[High]C₂
 11 sequences, but some of them are also claimed to condition deletion—not just devoicing—of the
 12 following vowel in the same environment. According to Kawakami (1977), /s/, /ts/, and /ϕ/
 13 condition deletion of the following /u/, while /k/ and /ʃ/ condition devoicing only, although
 14 recall that he offers no phonological or phonetic evidence. According to Whang (2014), vowel
 15 deletion occurs when the identity of the devoiced vowel is predictable from context. His
 16 recoverability-based theory predicts that /u/ will be deleted following /s/, /ts/, /ϕ/, and /k/ but

⁶ Although here and throughout we list /ts/ and /ϕ/ in slashes, we note that they are largely (but not entirely) predictable allophones: [ts] is the allophone of /t/ that occurs before /u/; [ϕ] is the allophone of /h/ that occurs before /u/.

1 that /u/ will be present (though devoiced) following /ʃ/. The predictions match Kawakami's
2 intuition for four of the five consonants in our stimuli (/s/, /ts/, /ʃ/, /ʒ/). The point of divergence
3 is the /k/ environment. Whang's theory predicts vowel deletion following /k/, whereas
4 Kawakami claims that the vowel is present (although devoiced) following /k/. The predictions
5 for the stimulus set from Whang (2014) are labelled as "(W)" in the first column of Table 2;
6 those due to Kawakami are labelled "(K)"; converging predictions are labelled as "(W,K)".

7

8 We did not include stimuli in which high vowels are surrounded by two sibilants, as it is known
9 that devoicing may be inhibited in this environment (Hirayama 2009; Fujimoto 2015; Maekawa
10 & Kikuchi 2005; Tsuchida 1997). In addition, if the vowel is followed by /h/, /ʃ/, or /ç/
11 devoicing may be inhibited (Fujimoto 2015). Our stimuli avoided this environment as well.

12

13 We avoided any words in which the vowel following the target vowel is also high, because
14 consecutive devoicing is variable (Fujimoto, 2015; Nielsen, 2015). We also chose near minimal
15 pairs in such a way that accent always matches within a pair. More specifically, /u/s in /hokusai/
16 and /jakuzai/ are both accented, meaning that the pitch fall begins at the /u/, but all the other
17 target /u/s are unaccented, and they carry low pitch. Although intonational accents can influence
18 vowel coarticulation, at least in English (Cho, 2004), Tsuchida (1997) shows that young
19 Japanese speakers, at the time of 1997, show no effects of pitch accent on devoicing, so that

1 controlling for accent is important but may not be crucial.⁷ All the stimulus words are common
2 words.⁸

3

4 The target words were displayed in the carrier phrase: *okee* _____ *to itte* ‘Okay say _____’.

5 The preceding word *okee* was chosen, as it ends with a vowel /e/, which differs in height from
6 the target vowel, /u/. Participants were instructed to speak as if they were making a request of
7 a friend.

8

9 Each participant produced 10-15 repetitions of the target words, generating a corpus of 690
10 tokens for analysis. We aimed to get 15 tokens from each speaker, but if a sensor came off in
11 the late stages of the experiment, after we had collected at least 10 repetitions of all target words,
12 we ended the session. Words were presented in Japanese script (composed of hiragana,

⁷ There are very little if any durational differences between accented and unaccented vowels (Beckman, 1986), which would otherwise potentially affect the deletability of /u/. Since accented /u/ is not longer than unaccented /u/, this is yet another reason not to be too concerned about the placement of accent.

⁸ The written frequencies of our stimulus items in the Balanced Corpus of Contemporary Written Japanese (BCCWJ), a 104.3 million word corpus of diverse written materials (Maekawa et al., 2014), are as follows: /ϕsoku/ (6,342), /ϕzoku/ (2,273), /ʃutaise:/ (4,480), /ʃudaika/ (1,584), /katsu/ (6,163), /katsudo:/ (31,440), /masutaa/ (1,517), /masuda/ (0), /hakusai/ (639), and /jakuzai/ (1,561). The only item with low frequency in the corpus is /masuda/, probably because it is a proper noun, although this word is not uncommon as a Japanese name.

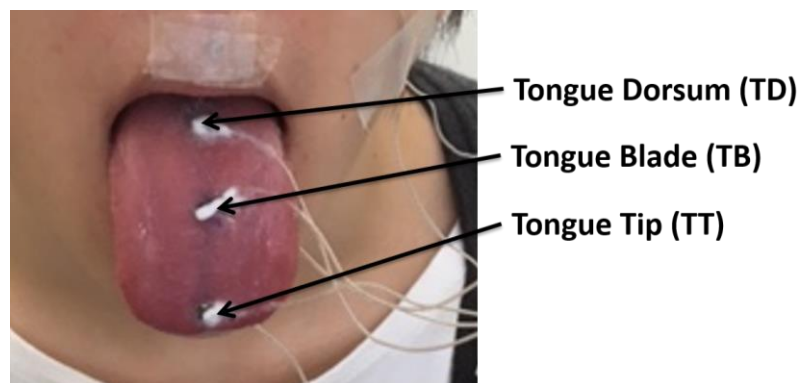
1 katakana and kanji characters as required for natural presentation) and fully randomized with
2 10 additional filler items that did not contain /u/.

3

4 ***Equipment***

5 The current experiment used an NDI Wave electromagnetic articulograph system sampling at
6 100 Hz to capture articulatory movement. The NDI wave tracks fleshpoints with an accuracy
7 typically within 0.5 mm (Berry, 2011). NDI wave 5DoF sensors were attached to three locations
8 on the sagittal midline of the tongue, and on the upper and lower lips near the vermilion border,
9 lower jaw (below the incisor), nasion and left/right mastoids. The most anterior sensor on the
10 tongue, henceforth TT, was attached less than one cm from the tongue tip. The most posterior
11 sensor, henceforth TD, was attached as far back as was comfortable for the participant, ~4.5-6
12 cm. A third sensor, henceforth TB, was placed on the tongue body roughly equidistant between
13 the TT and TD sensors. Figure 2 illustrates the location of the lingual sensors for one participant.
14 Acoustic data were recorded simultaneously at 22 kHz with a Schoeps MK 41S supercardioid
15 microphone (with Schoeps CMC 6 Ug power module).

16



17

18 **Figure 2: a representative illustration of lingual sensor placement for one participant.**

1 ***Stimulus display***

2 Words were displayed on a monitor positioned 25cm outside of the NDI Wave magnetic field.
3 Stimulus display was controlled manually using an Eprime script. This allowed for online
4 monitoring of hesitations, mispronunciations and disfluencies. These were rare, but when they
5 occurred, the experimenter repeated the trial. Participants were instructed to read the target
6 word in the carrier phrase fluently, as if providing instructions to friend, and told explicitly not
7 to pause before the target word. Each trial consisted of a short (500 ms) preview presentation
8 of the target word followed by presentation of the target word within the carrier phrase. The
9 purpose of the preview presentation was to further facilitate fluent reading of the target word
10 within the carrier phrase, since it is known that a brief visual presentation of a word facilitates
11 planning (Davis et al., 2015) and, in particular, to discourage insertion of phonological phrase
12 boundary between “okee (okay)” in the carrier phrase and the target word.

13

14 ***Post-processing***

15 Following the main recording session, we also recorded the occlusal plane of each participant
16 by having them hold a rigid object, with three 5DoF sensors attached to it, between their teeth.
17 Head movements were corrected computationally after data collection with reference to three
18 sensors on the head, left/right mastoid and nasion sensors, and the three sensors on the occlusal
19 plane. The head-corrected data was rotated so that the origin of the spatial coordinates
20 corresponds to the occlusal plane at the front teeth. All articulatory signals were smoothed using
21 Garcia’s robust smoothing algorithm (Garcia, 2010).

22

1 **Analysis**

2 ***Presence of voicing***

3 One of the authors and a research assistant each went through the spectrograms and waveforms
4 of all the tokens, and confirmed that /u/ in the devoicing environments are all devoiced (Figure
5 5 below provides a sample spectrogram), whereas /u/ in the voicing environments was voiced.
6 This is an unsurprising result, given that vowel devoicing is reported to be obligatory in the
7 normal speech style of Tokyo Japanese speakers (Fujimoto 2015).⁹

8

9 ***Lingual targets***

10 All stimulus items were selected so that the vowels preceding and following the target /u/ were
11 non-high. In order to progress from a non-high vowel to /u/, the tongue body must rise. For
12 some stimulus items, e.g., /ϕusoku/, /ϕuzoku/, /ʃutaise:/, /ʃudaika/, the tongue body may also
13 retract from the front position required for /e/ in the carrier phrase (*okee _____ to itte*) to the
14 more posterior position required for /u/. The degree to which the tongue body retracts for /u/,
15 i.e., the degree to which /u/ is a back vowel, has been called into question, with some data
16 suggesting that /u/ in Japanese is central (Nogita, Yamane, & Bird, 2013), similar to “fronted”
17 variants of /u/ in some dialects of English (e.g., Blackwood-Ximenes et al., 2017; Harrington,
18 Kleber, & Reubold, 2008). We therefore focus on the height dimension, in which /u/ is

⁹ Though see Maekawa & Kikuchi (2005) who show that devoicing may not be entirely obligatory in spontaneous speech—in their study, overall, /u/ is devoiced about 84% of the time in the devoicing environment. However, as Hirayama (2009) points out, their study is likely to contain environments where there are two consecutive high vowels, which sometimes resist devoicing (Kondo 2001; Nielsen 2015).

1 uncontroversially distinct from /o/. As an index of tongue body height, we used the TD sensor,
2 the most posterior sensor of the three sensors on the tongue, which provides comparable data
3 to past work using fleshpoint tracking to examine vowel articulation (Browman & Goldstein,
4 1992a; Johnson, Ladefoged, & Lindau, 1993).

5

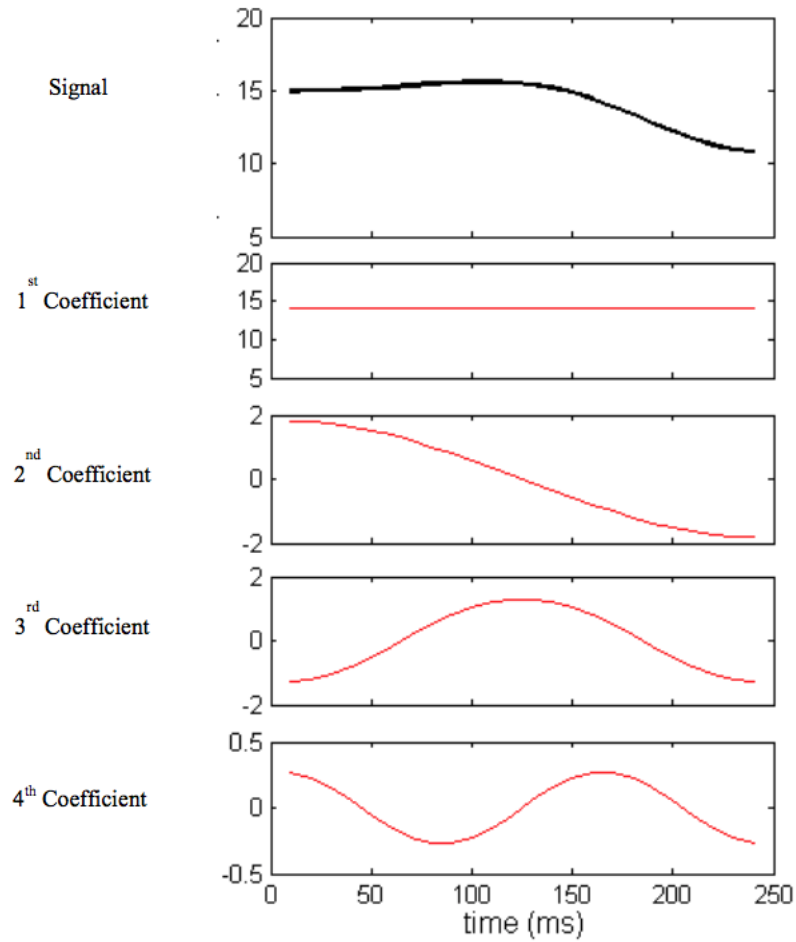
6 Our analytical framework makes use of the computational toolkit for assessing phonological
7 specification proposed by Shaw and Kawahara (submitted-a). This framework evaluates
8 presence/absence of an articulatory target based upon analysis of continuous movement of the
9 tongue body across $V_1C_1\underline{u}C_2V_3$ sequences (Figure 3). Analysis involves four steps: (1) fit
10 Discrete Cosine Transform (DCT) components to the trajectories of interest; (2) define the
11 targetless hypothesis based on linear movement trajectories between the vowels flanking /u/;
12 (3) simulate “noisy” linear trajectories using variability observed in the data, in which the means
13 are taken from the DCT coefficients of the targetless trajectory, and the standard deviations are
14 taken from the observed data; (4) classify voiceless tokens as either “vowel present” or “vowel
15 absent” based on comparison to the “vowel present” training data (voiced vowels) and the
16 “vowel absent” training data (based on linear interpolation), using a Bayesian classifier.

17

| 18

1 **Figure 3: Steps illustrating the computation analysis (based on Shaw & Kawahara,**
2 **submitted-a)**

3 **(A) Step 1: Fitting four DCT components to a trajectory. The top panel is the raw signal**
4 **of /V1C1uC2V3/. The rest of the panels shows each DCT components.**

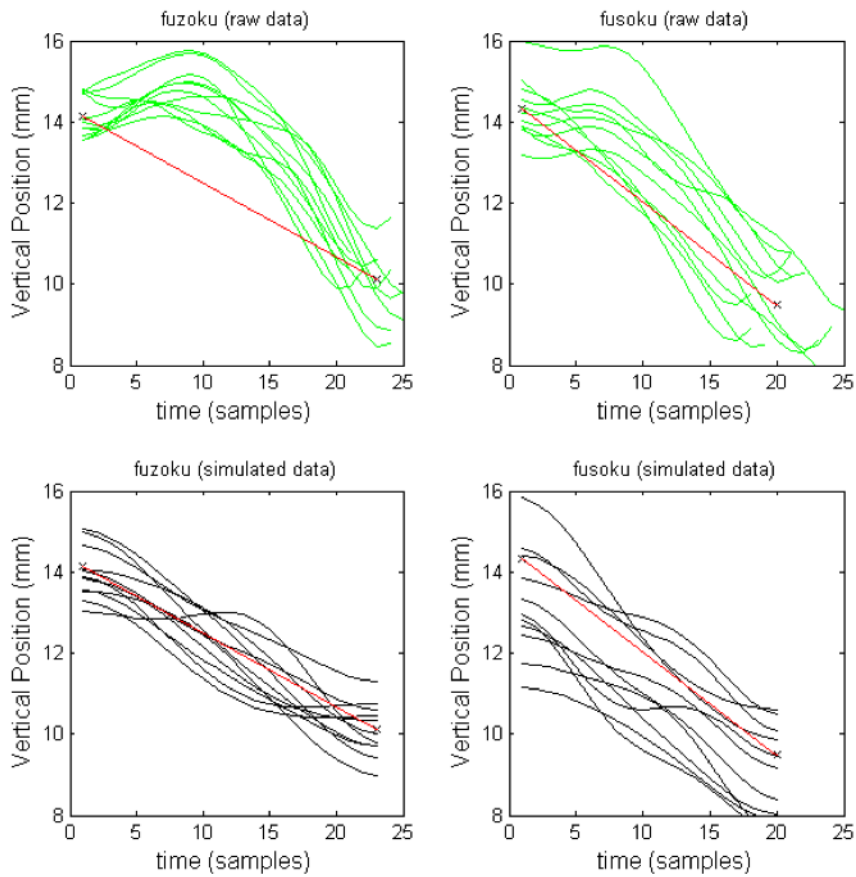


5

6

| 7

1 (B) Steps 2 & 3: Define a linear interpolation (shown as red lines) and generate a noisy
2 null trajectory (the bottom two panels), based on the variability found in the raw data
3 (the top two panels). Left=trajectories for / $\epsilon\phi$ uzo(ku)/; right=trajectories for
4 / $\epsilon\phi$ uso(ku)/.



5

| 6

1 **(C) Step 4: Train a Bayesian classifier in terms of DCT coefficients. “Vowel present”**
2 **defined in terms of voiced vowels. “Vowel absent” defined in terms of the linear trajectory.**

3

$$4 \quad p(T|Co_1, Co_2, Co_3, Co_4) = \frac{p(T) p(Co_1, Co_2, Co_3, Co_4|T)}{p(Co_1, Co_2, Co_3, Co_4)}$$

5

6 where

7

8 $T =$ targetless (linear interpolation), target present (voiced vowel)

9 $Co_1 =$ 1st DCT Coefficient

10 $Co_2 =$ 2nd DCT Coefficient

11 $Co_3 =$ 3rd DCT Coefficient

12 $Co_4 =$ 4th DCT Coefficient

13

14 Shaw and Kawahara (submitted-a) demonstrate that four DCT coefficients are sufficient for
15 representing TD height trajectories over VC_uCV intervals with an extremely high degree of
16 precision ($R^2 > .99$). They show moreover that each DCT coefficient has a plausible linguistic
17 interpretation: the 1st coefficient generally picks out general TD height across the analysis
18 interval, the 2nd coefficient captures V1-V3 movement, the 3rd coefficient picks out movement
19 associated with /u/, if any, and the 4th coefficient captures additional coarticulatory effects from
20 surrounding consonants.

21

22 We view the potential for interpreting DCT components as signal modulations associated with
23 linguistically relevant units, i.e., gestures, as an interesting difference from other times series
24 analyses (e.g., GAMMs, functional data analysis, SSANOVA). Nevertheless, we would like to
25 stress that, although it seems clear that the magnitude of the /u/ gesture is related to the 3rd DCT
26 component, we cannot demonstrate that the 3rd DCT coefficient is picking out all and only the
27 rising movement associated with /u/. For example, an increased magnitude of /u/ may also lead

1 to an increase in the 1st DCT coefficient, which is related to the average trajectory height. For
2 this reason, we took what we believe to be the most conservative approach and based our
3 classification of trajectories on all four DCT coefficients. In this way, the properties of DCT
4 that are most pertinent to our analysis are the compression property—four DCT coefficients
5 provide a very close approximation to the raw trajectories—and the statistical independence of
6 the coefficients, an assumption of the naïve Bayes classifier that is met by DCT coefficients.
7 We report the classification results for each devoiced token in terms of the posterior probability
8 of targetlessness, i.e. the likelihood that the trajectory follows a linear interpolation between V1
9 and V3 instead of rising like the tongue body does in voiced tokens of /u/.

10

11 We acknowledge that our choice to model the entire VC₁uC₂V trajectory has the consequence
12 that any differences in TD height associated with the voicing specification of C₂ will also factor
13 into the results. Consonants that contrast minimally in voicing are known to have some
14 differences in lingual articulation, owing in part to aerodynamic factors associated with voicing.
15 In Japanese, voiceless stops and affricates tend to have tighter constrictions than voiced stops,
16 as indicated by greater lingual-palatal contact (Kochetov & Kang, 2017) while the pattern is
17 reversed for fricatives—voiced fricatives tend to have more contact than voiceless fricatives, at
18 least in the anterior portion of the palate (Nakamura, 2003). In our items, C₂ was always a
19 coronal consonant. In this precise environment (C₂ following a devoiced consonant), Nakamura
20 (2003) reports EPG data showing a difference in palate contact between Japanese /z/ and /s/
21 (more contact for /z/) but only in the anterior portion of the palate. There were no differences
22 at the tongue dorsum, the focus of our analysis.

23

24 The four hypotheses introduced in (1) can each be evaluated by examining the distribution of
25 posterior probabilities across tokens. Figure 4 presents hypothetical distributions corresponding

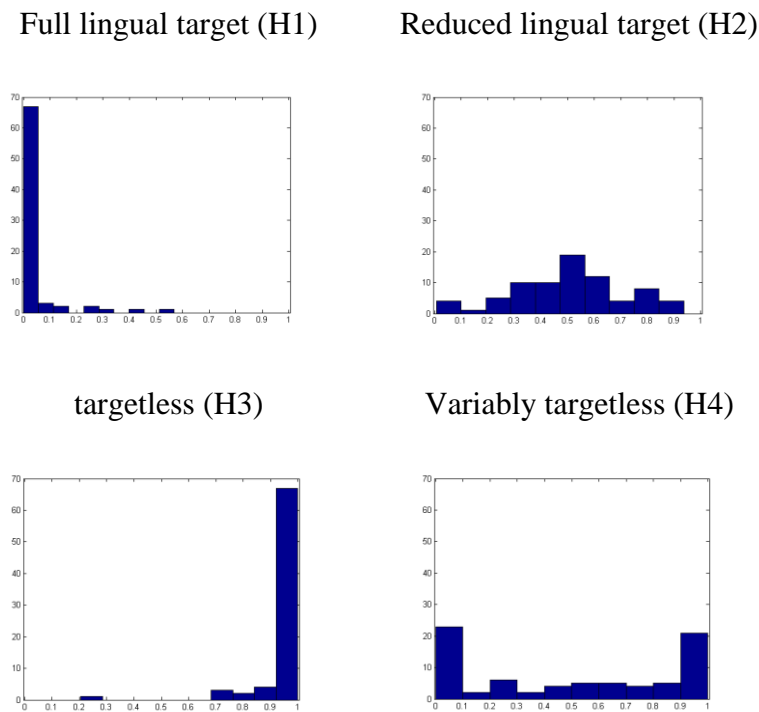
Lingual articulation of devoiced /u/

1 to each of the hypotheses. If devoiced vowels have full lingual targets (H1), then we expect the
2 probability of targetless to be low, as is shown in the top left panel of Figure 4. This figure was
3 made by submitting items with voiced /u/ to the classifier. If H2 is correct, and devoiced vowels
4 have reduced lingual targets, then the distribution of posterior probabilities should be centered
5 around .5, as in the top right panel. This figure was made by submitting DCT components half
6 the magnitude of voiced vowels to the classifier. If devoiced vowels lack lingual targets
7 altogether (H3), then the distribution of posterior probabilities should be near 1.0, as in the
8 bottom left panel of the figure. This figure was made by submitting DCT coefficients from
9 noisy linear trajectories (simulated) to the classifier. Lastly, if lingual targets are variably
10 present, as in H4, then we expect to see a bimodal distribution, with one mode near 0 and the
11 other near 1.0, as shown in the bottom right panel. This figure was made by submitting a mix
12 of tokens from sampled from linear trajectories and from voiced vowels to the classifier.

13

| 14

1 **Figure 4: Four hypothetical posterior probability patterns. The vertical axis of each**
2 **histogram shows posterior probabilities generated by the Bayesian classifier**
3 **summarized in Figure 3(C). The histogram in the top left panel was obtained by**
4 **submitting the /ϕuzoku/ (voiced vowel) tokens to the Bayesian classifier. The histogram**
5 **in the bottom left was obtained by submitting the same number of simulated linear**
6 **interpolation “vowel absent” trajectories to the classifier. The top right panel was**
7 **generated by stochastic sampling of DCT coefficients that were averaged between**
8 **“target present” (H1) and “target absent” (H3) values. The right bottom panel was**
9 **created by sampling over both targetless and full vowel target tokens.**
10



Posterior probabilities of targetlessness

11
12 An anonymous reviewer points out that the velar stop preceding /u/ in /hakusai/ complicates
13 the interpretation of our classification results for this item. In our other target items, the

1 devoiced /u/ is immediately preceded by coronal, /f/, /ts/, /s/, or labial, /p/, consonants, which
2 do not dictate a large rise in tongue dorsum (TD) height. For these items, we interpret a rise in
3 TD height as progress towards the goal of /u/ production. More precisely, our analytical
4 approach assesses whether the height of the devoiced /u/ trajectory is closer to a linear trajectory
5 between flanking vowels, i.e., no rise at all, or the TD rise observed in a voiced vowel in the
6 same consonantal environment. In /hakusai/, the devoiced /u/ is preceded by a velar stop, /k/.
7 The TD rises from the first /a/ towards the target of /k/ before falling again for the following /a/
8 (see Figure 6, fourth row from the top). The large TD rise for the /k/ immediately preceding the
9 target /u/ may obscure differences in vowel height specification across voiced and devoiced /u/.
10 Specifically, the TD trajectory in the /akusa/ portion of /hakusai/ may be more similar to the
11 TD trajectory of the /akuza/ portion of /jakuzai/ than to a linear interpolation between /a/ and
12 /a/ because of the influence that /k/ exerts over the TD in both /akusa/ and /akuza/. For this
13 reason, we have excluded /hakusai/ from the classification analysis. The raw TD trajectories
14 from /hakusai/ and /jakuzai/ are reported and these data are included in the analyses of inter-
15 consonantal timing described below, as this analysis does not require that we assess the
16 influence of /u/ on the TD signal.

17

18 ***Consonantal timing across devoiced vowels***

19 In addition to examining the continuous trajectory of the tongue dorsum, we also investigated
20 the timing of consonants preceding and following /u/. If the consonants flanking /u/ are
21 coordinated in time with the lingual gesture for the vowel (e.g., see Smith, 1995 for a concrete
22 proposal), reduction or deletion of that vowel gesture may perturb consonant timing. In this
23 way, consonant timing offers another angle on how devoicing influences lingual articulation.

24

1 For this analysis, we identified articulatory landmarks from consonants flanking /u/ based on
2 the primary oral articulator, e.g., tongue tip for /t/, /s/, tongue blade for /ʃ/, tongue dorsum for
3 /k/, lips for /ɸ/, etc, for each gesture. The affricate /ts/ in /katsutoki/ and /katsudou/ was treated
4 as a single segment. We determined the start and end of consonantal constrictions with
5 reference to the velocity signal in the movements toward and away from constriction. Data were
6 displayed in Mview and articulatory landmarks were extracted using *findgest*, an Mview
7 labeling procedure (Tiede, 2005). Both the start of the constriction, a.k.a. the achievement of
8 target of the consonant, and the end of the constriction, a.k.a. the release landmark, were
9 extracted at the timepoint corresponding to 20% of peak velocity in the movement
10 towards/away from consonantal constrictions, a heuristic applied extensively in other recent
11 (Bombien, Mooshammer, & Hoole, 2013; Gafos, Hoole, Roon, & Zeroual, 2010; Marin, 2013;
12 Marin & Pouplier, 2010; Shaw, Chen, Proctor, & Derrick, 2016; Shaw, Gafos, Hoole, &
13 Zeroual, 2011; Shaw, Gafos, Hoole, & Zeroual, 2009). As a first pass, we used all three spatial
14 dimensions of the positional signal and corresponding tangential velocities to parse consonantal
15 landmarks. This approach was practicable for a majority of the consonant tokens in our corpus.
16 However, for some tokens, parsing landmarks based on the tangential velocity was problematic.
17 One issue was that the amount of spatial displacement associated with a consonant gesture
18 was sometimes too small to produce a prominent velocity peak. This occurred most often for
19 the tongue tip movement from /s/ to /d/ in /masuda/. When the velocity peak is very small it
20 cannot reliably delineate controlled movement. Consonant tokens that could not be parsed
21 from the velocity signal were excluded from analysis (25 tokens, 3.6% of the data). Another
22 issue was that, in some tokens, there were not distinct tangential velocity peaks associated with

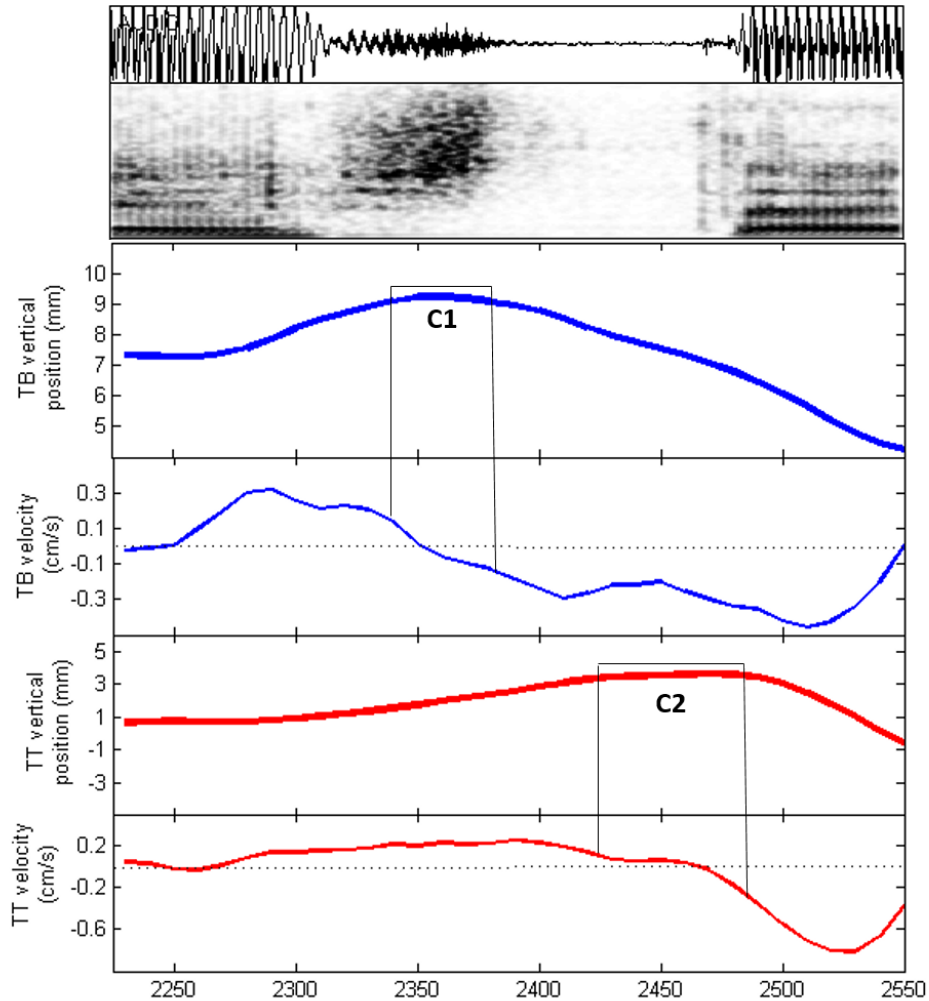
1 the release of C1 and the movement towards C2. This is because movement towards C2 in one
2 dimension, such as advancement of the tongue for /t/ in /ʃutasise:/ or /d/ in /ʃudaika/ overlapped
3 in time with movement in another dimension associated with C1, such as lowering of tongue
4 for /ʃ/. For many of these cases, we were able to isolate distinct velocity peaks for C1 and C2
5 by focusing on the primary spatial dimension of movement for the gesture, such as lowering
6 for the release of /ʃ/ and raising toward the target for /t/. This approach, suggested in
7 *Guidelines for using Mview* (Gafos, Kirov, & Shaw, 2010) allowed us to consider a greater
8 number of tokens for analysis. The parseability of consonants based on tangential velocity (as
9 opposed to component velocity) was unrelated to whether C2 was voiced or voiceless. The total
10 number of tokens parsed by tangential vs. component velocities is provided by item in the
11 Appendix.

12
13 Figure 5 provides an example of consonantal landmarks parsed for for /ʃ/ and /t/, the consonants
14 flanking target /u/ in /ʃutaise:/ based on movement in the vertical dimension only. The vertical
15 black lines indicate the timestamp of the consonantal landmarks. They extend from the
16 threshold of the velocity peak used as a heuristic for parsing the consonant up to the
17 corresponding positional signal. The dotted line in the velocity panels extends from 0 (cm/s),
18 or minimum velocity.

19
20 The interval between landmarks, labeled C₁ for /ʃ/ and C₂ for /t/, are the consonant constriction
21 durations. The interval between the consonants, or inter-consonantal interval (ICI), defined as
22 the achievement of target of C₂ minus the release of C₁ (see also, inter-plateau interval in Shaw

1 & Gafos, 2015) was also analyzed. At issue is whether this interval varies with properties of
2 the lingual gesture for the intervening /u/.

3



4

5 **Figure 5: An illustration of how consonantal landmarks, /ʃ/ (C1) and /t/ (C2), were**
6 **parsed from a token of /ʃutaise:/.** The thick blue line shows the vertical position of the
7 **tongue blade (TB); the thin blue line shows the corresponding velocity signal. The thick**
8 **red line shows the vertical position of the TT; the thin red line shows the corresponding**
9 **velocity signal. Consonant onsets and offsets were based on a threshold of peak velocity,**
10 **vertical black lines, in the movements toward and away from target.**

11

1 ***Data exclusion***

2 The tongue tip sensor became unresponsive for S04 on the sixth trial. We think that this was
3 due to wire malfunction, possibly due to the participant biting on the wire. The tongue tip
4 trajectory is relevant only for the analysis of consonant timing. Due to missing data, the
5 consonant timing analysis for this speaker is based on only 5 trials. Analyses involving other
6 trajectories are based on 11 trials for this participant.

7

8 **Results**

9 ***Presence/absence of articulatory height targets***

10 Figure 6 summarizes the data on TD height across speakers and words. Each panel shows TD
11 height (y-axis) over time (x-axis) for tokens of a target word with a voiced vowel (blue lines)
12 and devoiced counterpart (red lines). The columns show data from different speakers and the
13 rows show the different voiced-devoiced dyads. An interval of 350 ms (35 samples of data),
14 beginning with the vowel preceding the target, is shown in each panel. Despite the variation in
15 speech rate (note, for example, that the rise of the TD for /k/ at the right of panels displaying
16 /ϕusoku/~/ϕuzoku/ in the top row is present to different degrees across speakers), a 350 ms
17 window is sufficient to capture the three-vowel sequence including the target /u/ and preceding
18 and following vowels for all tokens of all words across all speakers.

19

20 To facilitate interpretation of the tongue height trajectories, annotations are provided for the
21 first speaker (leftmost column of panels). The trajectories begin with the vowel preceding the
22 target vowel, e.g., /e/ from the carrier phrase in the case of /ϕusoku/ and /ʃutaise:/, /a/ in
23 /katsutoki/, etc. Movements corresponding to the vowel following /u/ are easy to identify—
24 since the vowel following /u/ is always non-high, it corresponds to a lowering of the tongue.

Lingual articulation of devoiced /u/

1 The label for the target vowel, /u/, has been placed in slashes between the flanking vowels. We
2 note also that the vertical scales have been optimized to display the data on a panel-by-panel
3 basis and are therefore not identical across all panels. In particular, across speakers the TD is
4 lower in the production of *masuda~masutaa* than for many of the other dyads and that this
5 influences the scale for most speakers (S02-S06).

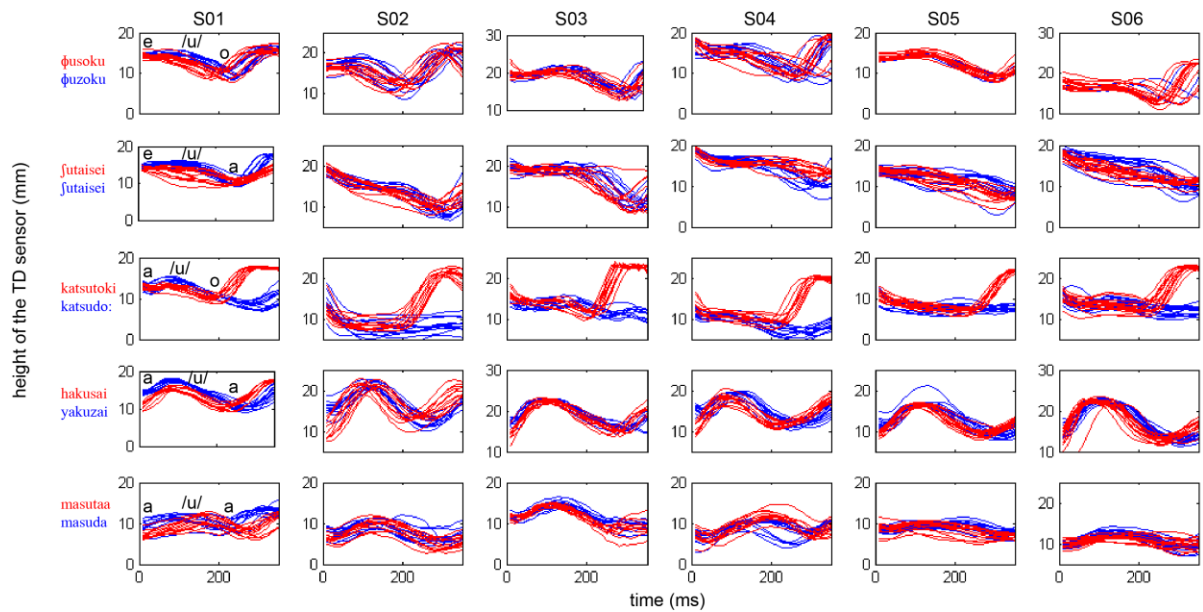
6

7 Differences between voiced and devoiced dyads (red and blue lines, respectively, in the figure)
8 include cases in which the tongue tends to be higher in the neighborhood of /u/ for the voiced
9 than for the devoiced member of the dyad. The top left panel, /ϕsoku/ for S01, exemplifies
10 this pattern (a zoom-in plot is provided in Figure 7). In this panel, the blue lines rise from /e/ to
11 /u/ while the devoiced vowel trajectory is a roughly linear trajectory between /e/ and /o/. Our
12 analysis assesses this possibility specifically by setting up stochastic generators of the
13 competing hypothesis, lingual target present (based on the voiced vowel) vs. lingual target
14 absent (based on linear interpolation), and evaluates the degree to which the voiceless vowel
15 trajectory is consistent with these options.

16

1

Figure 6: vertical TD trajectories of all speakers, all items. The



2

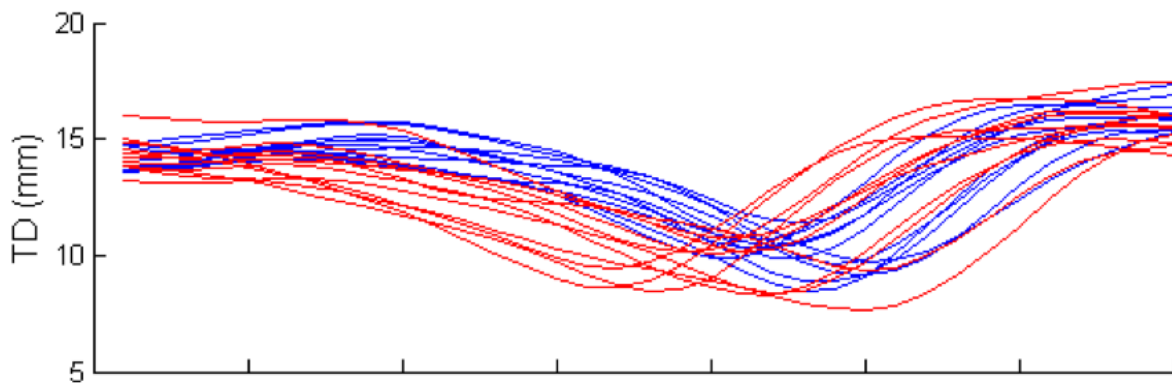
3 **Figure 7: TD trajectories of /ϕ_usoku/ (red) and /ϕ_uzoku/ (blue) for S01. The blue lines**
 4 **generally show a rise from the /e/ in the carrier phrase to the /u/ before lowering again to**

5 **the target for /o/. Some red lines show a similar rise from /e/ to /u/; others follow a**

6 **downward cline from /e/ to /o/ without rising for /u/. The computational analysis**

7 **described above classifies each red line as either belonging to the voiced category (blue**

8 **lines) or to a linear trajectory from /e/ to /o/.**



9

10

1 Table 3 provides the average posterior probabilities across tokens by speaker and by word.
 2 Since our analysis involves a stochastic component—the linear trajectory corresponding to the
 3 vowel absent scenario is stochastically sampled—we repeated the simulation and classification
 4 steps multiple times. The standard deviations of the posterior probabilities across 1000
 5 repetitions are given in parenthesis. The probability of targetlessness varies across speakers
 6 rather widely, from speakers that have a high probability of targetless vowels, e.g. 0.70 for S01,
 7 to speakers with a much lower probability of targetlessness, e.g., 0.23 for S05. There are also
 8 differences across items. The targetlessness probability is highest for /futaise:/ (0.69) and
 9 /katsutoki/ (0.60) with lower probabilities for /ϕusoku/ (0.43) and /masutaa/ (0.36).

10

11 **Table 3: Posterior probability of lingual targetlessness (vowel height)**

	S01	S02	S03	S04	S05	S06	average
<u>ϕ</u> usoku	0.47(.11)	0.38(.14)	0.75(.07)	0.85(.08)	0.01(.01)	0.11(.06)	0.43
<u>f</u> utaise:	0.92(.03)	0.66(.17)	0.81(.06)	0.92(.07)	0.05(.05)	0.80(.15)	0.69
kats <u>u</u> toki	0.70(.22)	0.23(.11)	0.81(.13)	0.93(.03)	0.13(.11)	0.78(.13)	0.60
mas <u>u</u> taa	0.73(.20)	0.11(.08)	0.04(.05)	0.02(.02)	0.74(.20)	0.52(.14)	0.36
average	0.70	0.35	0.60	0.68	0.23	0.55	

12

13 Figure 8 shows histograms of posterior probabilities both across speakers (left panels) and
 14 within speakers (right panels). Probabilities close to 1 indicate a linear trajectory, i.e., no rise in
 15 TD height for /u/, while probabilities near 0 indicate that the trajectory for the devoiced vowels
 16 resembles the trajectory for voiced vowels. As illustrated earlier (Figure 4), each of the
 17 hypotheses about lingual articulation of devoiced vowels motivated from the literature makes
 18 distinct predictions about the shape of these histograms. Figure 8 shows that, for all words, the
 19 distribution of probabilities is distinctly bimodal. This indicates that many of the tokens of

1 devoiced vowels in our corpus were either produced with a full lingual target or produced as
2 linear interpolation between flanking vowels. Only a small number of tokens are intermediate,
3 i.e., posterior probabilities in the range of .5 indicating reduction relative to voiced vowels but
4 not to the degree that would result in linear interpolation. The bimodal distributions support the
5 “optional targetless hypothesis” (=H4) (see also Figure 5), at least for the population of tokens
6 drawing from the six speakers in this study.

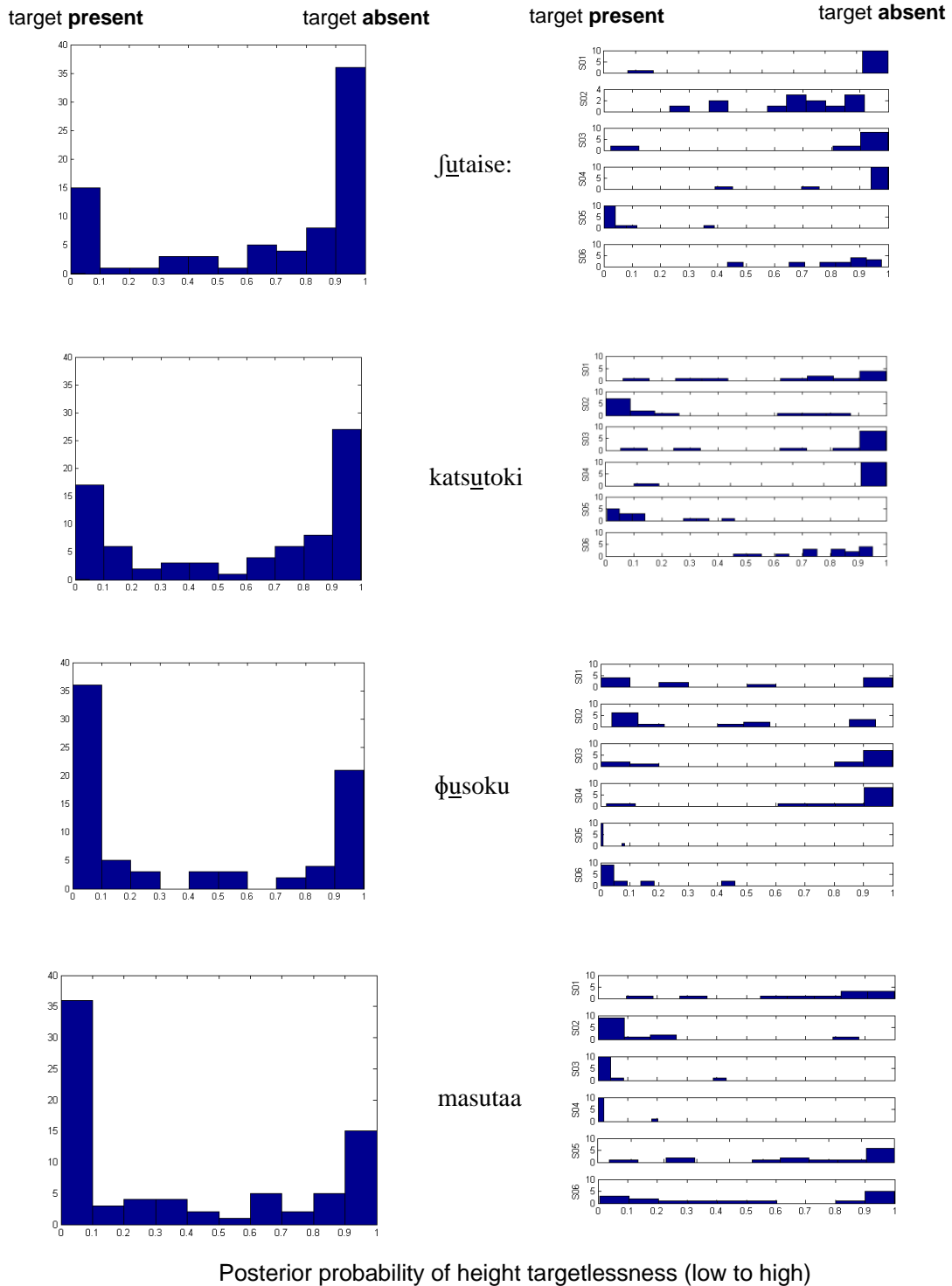
7

8 The right panels of Figure 8 present posterior probabilities by speaker. These figures illustrate
9 that both inter- and intra-speaker variation contribute to the bimodality of posterior
10 probabilities. Consider /futaise:/. Of the six speakers in the study, three of them, S01, S03, S04,
11 show strong tendencies towards a linear trajectory, i.e., an /u/ lacking a height target. For these
12 three speakers, the vast majority of tokens have targetless probabilities greater than .9. Two
13 others, S02 and S06, show slightly weaker tendencies towards targetlessness. The majority of
14 the tokens for these two speakers are between .7 and .9, values that indicate these tokens are
15 still much closer to the linear trajectory than to the voiced vowel. Only one speaker, S05, shows
16 a strong tendency towards producing a full vowel height target. This speaker’s data contributes
17 most heavily to the mode near 0 in the group data. Of the 15 tokens of /futaise:/ with targetless
18 probabilities of less than .1, 10 tokens are from S05. Thus, part of the bimodality in the group
19 data derives from inter-speaker differences: speaker S05 tends to produce full vowel height
20 targets in devoiced vowels; the other five speakers do not. However, we can see that optionality
21 within speakers, i.e., intra-speaker variation, also contributes somewhat to bimodality in the
22 group data. Speakers S01, S03, and S04, those who showed the strongest tendency towards
23 targetlessness, indicated by the peaks in the histograms between .9 and 1, all produced one
24 token that is at the other end of the histogram, indicating that it resembles the height trajectory
25 of the voiced vowel. The other items tell similar stories vis the contribution of both inter- and

Lingual articulation of devoiced /u/

1 intra-speaker variation to bimodality. Overall, it seems that bimodality in the group data derives
2 both from across speaker variation in the probability with which /u/ is produced with a height
3 target as well as intra-speaker variation. All six speakers demonstrated that the linear trajectory
4 for /u/ is within their production repertoire, producing at least one item with high probability of
5 targetlessness.
6

Lingual articulation of devoiced /u/



1 **Figure 8: Posterior probability of targetlessness for each token organized by item. The**
 2 **left panels aggregate across speakers; the right panels show probabilities for each**
 3 **speaker from S01 (top) to S06 (bottom). All items show a roughly bimodal pattern**
 4 **anchored by tokens with high probabilities of targetlessness (numbers close to 1) on the**
 5 **right side of the figures and tokens with low probabilities of targetlessness (numbers**
 6 **close to 0) on the left side of the figures.**
 7

1 ***The inter-consonantal interval***

2 We next turn to the timing of the consonants flanking /u/, asking whether devoicing influences
3 relative timing between the preceding and following consonants. Past work on laryngeal control
4 of devoiced vowels indicates that the laryngeal gestures associated with flanking consonants
5 aggregate to form one large laryngeal gesture near the center of the vowel (Fujimoto et al. 2002;
6 Fujimoto 2015; see Figure 1). Here, we investigate the consequences of this laryngeal
7 reorganization for the oral gestures associated with the consonants. A decrease in the inter-
8 consonantal interval, defined as the interval spanning from the release of C₁ to the achievement
9 of target of C₂, across devoiced vowels relative to voiced vowels would indicate that laryngeal
10 reorganization “pulls” the oral gestures of the consonants closer together in time. Alternatively,
11 a consistent ICI across voiced and devoiced vowels would indicate that articulatory binding
12 (consonant-internal temporal organization) of oral and laryngeal gestures is perturbed to
13 achieve devoicing.

14

15 Alongside our interest in the ICI interval as an indication of how oral gestures respond to
16 laryngeal reorganization, given the results above on vowel height targets, we can also ask
17 whether ICI is impacted by the presence/absence of a lingual height target for the intervening
18 vowel. We opted to analyze the data in terms of this categorical difference (instead of using the
19 raw probabilities as predictors) because, as shown in the histograms in Figure 8, the data are
20 largely categorical in nature. We applied the Bayesian decision rule, interpreting (targetless)
21 probabilities greater than .5 as indicating that the vowel height target was absent and
22 probabilities less than .5 indicate that the target was present.

23

24 Figure 9 summarizes the inter-consonantal interval (ICI) across words containing a voiced
25 vowel and words containing a devoiced vowel. Since ICI is a difference, it is possible for it to

1 be negative. This can happen when C2 achieves its target before the release of C1, as commonly
2 observed in English consonant clusters (e.g., Byrd, 1996). The average ICI is around zero for
3 $\phi_{\underline{u}}soku \sim \phi_{\underline{u}z}oku$ indicating that the lips remain approximated until the tongue blade achieves its
4 target, a temporal configuration which does not at all interfere with the achievement of a height
5 target for the vowel. The / ϕs / sequence of consonants is a front-to-back sequence, in that the
6 place of articulation for the labial fricative, the first consonant, is anterior to the place of
7 articulation of the alveolar fricative, the second consonant. Consonant clusters with a front-to-
8 back order of place tend to have greater overlap than back-to-front clusters (Chitoran,
9 Goldstein, & Byrd, 2002; Gafos, Hoole, et al., 2010; Wright, 1996; Yip, 2013). For other dyads,
10 the consonantal context requires longer average ICI's, with median values ranging from ~50
11 ms to ~170 ms. The variation reflects the broader fact about Japanese that consonantal context
12 has a substantial influence on the duration of the following vowel (see Shaw & Kawahara,
13 submitted-b for a large scale corpus study).

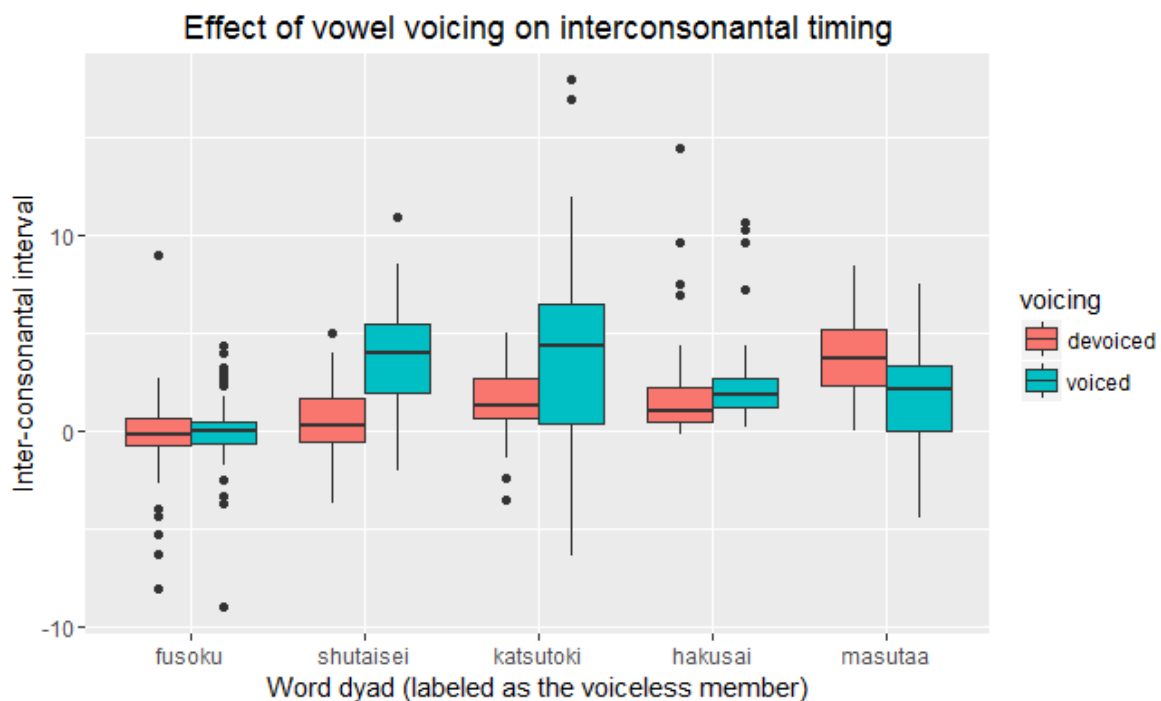
14

15 With respect to how devoicing influences ICI, Figure 9 shows that the effect of vowel devoicing
16 on ICI varies across consonantal environments. For *futaise:~fudaika* and *katsutoki~katsudo*”,
17 and to a lesser degree *hakusai~jakuzai*, ICI is longer when the vowel is voiced. For
18 $\phi_{\underline{u}}soku \sim \phi_{\underline{u}z}oku$, ICI is quite similar across words and for *masutaa~masuda*, the presence of
19 vowel voicing actually results in shorter ICI.

20

21 To assess the statistical significance of the trends shown in Figure 9, we fitted a series of three
22 nested linear mixed effects models to the ICI data using the *lme4* package (Bates, Maechler,
23 Bolker, & Walker, 2014) in R. The baseline model included word DYAD as a fixed factor, since
24 it is clear that ICI depends in part on the identity of the particular consonants, and random
25 intercepts for speaker, to account for speaker-specific influences, e.g., speech rate, on ICI. The

1 second model added VOWEL VOICING to the baseline model as a fixed factor. The third model
2 included the interaction between DYAD and VOWEL VOICING.
3
4 Table 4 summarizes the model comparison. Adding VOWEL VOICING to the model leads to
5 significant improvement over the baseline, indicating a general trend for shorter ICI across
6 devoiced vowels. However, as we observed, VOWEL VOICING affects ICI for some dyads
7 positively and for other dyads negatively. Adding the interaction between DYAD and VOWEL
8 VOICING leads to further improvement over the model with DYAD and VOWEL VOICING as non-
9 interacting fixed factors. These results indicate that the trends observed in Figure 9 are
10 statistically reliable. ICI varies depending on whether the vowel is voiced or devoiced; how it
11 varies depends also on the word or possibly on the identity of the flanking consonants within
12 the word.



13
14 **Figure 9: Inter-consonantal interval duration for each item, classified by target**
15 **absence/presence.**

1

2

Table 4: Model comparison showing effect of (de)voicing on ICI duration.

Model of ICI	Df	AIC	BIC	logLik	Chisq	Pr(>Chisq)
DYAD + (1 speaker)	7	7869	7901	-3928	-	-
DYAD + VOWEL_VOICING + (1 speaker)	8	7851	7887	-3917	20.39	0.000006***
DYAD* VOWEL_VOICING + (1 speaker)	12	7802	7856	-3889	56.99	1.2e-11***

3

4

In the light of our analysis of vowel height targets in the preceding section, it is notable that the

5

two dyads that showed the greatest effect of devoicing on ICI—*futaise:~fudaika* and

6

katsutoki~katsudo—are also those that have the highest probability of vowel height

7

targetlessness. Given this pattern of results, it is possible that the effect of vowel voicing on ICI

8

may be attributable to the occasional (height) targetlessness of flanking vowels. The intuitive

9

idea is that consonants move closer together when the intervening vowel is reduced or absent.

10

We assessed this possibility through further model comparison. In one comparison, we replaced

11

the VOWEL_VOICING factor with another factor, TARGET PRESENCE, which we determined

12

according to the Bayesian classification of tongue dorsum trajectories. We also fit a model

13

including both factors, VOWEL_VOICING and TARGET PRESENCE, to see if TARGET PRESENCE

14

would explain variance in ICI above and beyond the VOWEL_VOICING factor. We note here that

15

these model comparisons required that we exclude the *hakusai~jakuzai* dyad (see discussion in

16

previous section) because of the ambiguity that the velar consonant introduces into the

17

classification results. We therefore refit the baseline models to the subset of data for which we

18

were able to report classification results (525 tokens). As with the models of the full data set

19

reported in Table 4 (665 tokens), including an interaction between DYAD and VOWEL_VOICING

20

showed significant improvement over baseline. However, as shown in Table 5, adding TARGET

21

PRESENCE and the interaction between TARGET PRESENCE and DYAD as fixed factors, despite the

1 added complexity, resulted in only marginal improvement ($p = .07$). Moreover, the direct
 2 comparison of VOWEL VOICING and TARGET PRESENCE rendered negligible differences. From
 3 these results, we conclude that the observed effect of devoicing on ICI is not due to vowel
 4 height targetlessness alone; rather, it appears to be a genuine effect of devoicing. This is true
 5 despite its largest effects coming in words that also happen to frequently lack vowel height
 6 targets.

7

8 **Table 5: Model comparison showing effect of vowel height target on ICI duration.**

Model of ICI	Df	AIC	BIC	logLik	Chisq	Pr(>Chisq)
DYAD*VOWEL_VOICING + (1 speaker)	10	6194	6236	-3086		
DYAD*VOWEL_VOICING + DYAD* TARGET_PRESENSE + (1 speaker)	14	6193	6253	-3083	8.67	0.0698

9

10 To summarize the results on the ICI interval, we found that the effects of vowel devoicing show
 11 substantial variation across words. Some combinations of oral gestures are pulled closer
 12 together in time when they flank a devoiced vowel than when they flank a voiced vowel.
 13 However, shortening of ICI does not seem to be due to the occasional absence of a vowel height
 14 target. Rather, the chain of causation may run the other direction. Vowel reduction, including
 15 the absence of a height target, may be driven in part by the proximity of the flanking consonants.
 16 Supra-laryngeal gestures brought closer together in time due to laryngeal reorganization may
 17 encourage non-specification of vowel gestures on one or many dimensions.

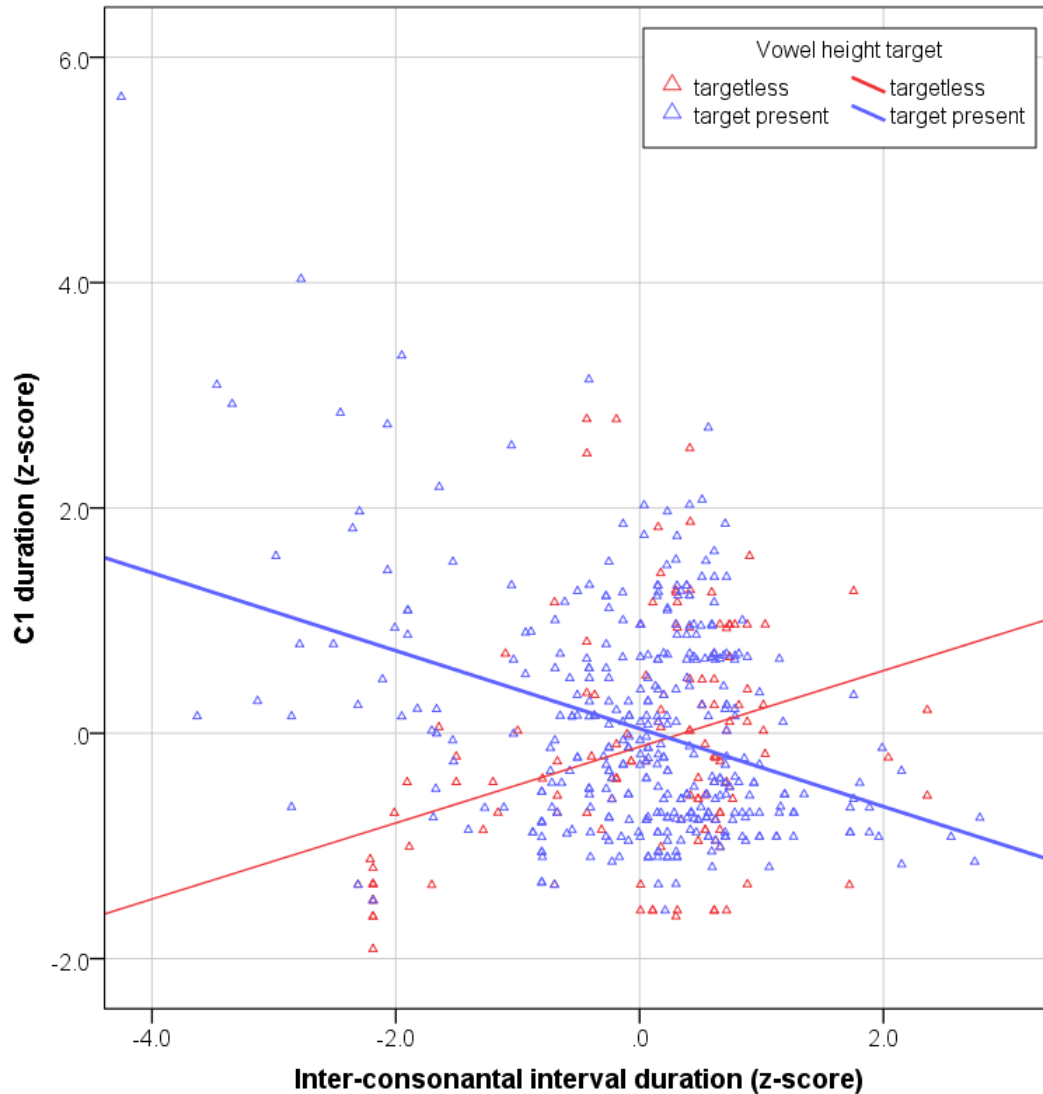
18

19 We conclude the presentation of the results by evaluating the correlation between ICI duration
 20 and C₁ duration. The predicted relation between these intervals depends on coordination
 21 topology (Table 1). As long as the /u/ is present (phonetically specified), we assume that it will

1 be coordinated with the preceding consonant, i.e., C-V coordination (Browman & Goldstein,
2 2000; Gafos, 2002; Smith, 1995). If /u/ is absent, then the consonants flanking the vowel may
3 be timed to each other, i.e., C-C coordination (Browman & Goldstein, 2000; Gafos, 2002; Shaw
4 & Gafos, 2015). As summarized in Table 1, C-V coordination predicts a trade-off (negative
5 correlation) between C_1 duration and ICI duration in our data, because as C_1 shortens, it will
6 expose more of the vowel and ICI will lengthen. C-C coordination on the other hand, predicts
7 no such relation between C_1 and ICI. However, all else equal, we expect a positive correlation
8 between adjacent intervals, such as C_1 and ICI, as both would be influenced by common factors,
9 such as speech rate and ease of lexical access. Investigating the correlation between C_1 and ICI
10 thus affords the opportunity for an assessment of vowel presence independent from our analysis
11 of TD trajectories, which focused only on the height dimension. A negative correlation between
12 C_1 and ICI, a characteristic of C-V coordination, provides evidence, albeit indirect, that a vowel
13 gesture is present in the signal. Figure 10 shows a scatterplot of C_1 duration and ICI. Since we
14 are comparing across items that have different average ICI, we have z-scored C_1 duration and
15 ICI within item to avoid obtaining spurious correlation. The blue triangles are tokens that
16 contain a vowel height target, according to our analysis of TD trajectories. The red triangles
17 represent tokens that lack a height target. The blue and red lines are linear fits to the tokens with
18 and without vowel height targets, respectively.

19

1



2
3

4 **Figure 10: The correlation between C1 duration (y-axis) and ICI (x-axis) across all items**
5 **in the corpus. Blue triangles represent voiced vowels and devoiced vowels classified as**
6 **having a vowel height target. Red triangles represent devoiced vowels classified as**
7 **lacking a vowel height target. The linear regression line fit to the blue triangles (target**
8 **present cases) shows a significant negative trend; the red triangles (target absent cases)**
9 **show the opposite direction of correlation.**

10

1 The distributions of C1 and ICI values both showed significant deviations from normality (C1
2 was right skewed; ICI was left skewed) according to a Shapiro-Wilkes test ($C1_{\text{target_present}}$:
3 $w(386) = .894, p < .001$; $C2_{\text{targetless}}$: $w(139) = .895, p < .001$; $ICI_{\text{target_present}}$: $w(386) = .975, p <$
4 $.001$; $ICI_{\text{targetless}}$: $w(139) = .969, p = .003$). We therefore conducted nonparametric Spearman
5 correlation analyses. There was a significant negative correlation between C1 and ICI when the
6 vowel height target is present ($\rho(386) = -0.19; p < .001$). When absent, the relation between C1
7 and ICI is positive and also statistically significant ($\rho(139) = 0.289, p < .001$). Thus, the
8 negative tradeoff between consonant duration and ICI duration is only maintained when the
9 lingual vowel height target is present. When the vowel is absent, C1 is positively correlated
10 with ICI. This provides a converging argument for the variation in vowel specification found
11 across tokens in our analysis of TD height trajectories. There are clearly systematic differences
12 in temporal organization between those tokens classified as containing a vowel height target
13 and those tokens classified as lacking a vowel height target. Moreover, the direction of the
14 differences is as expected is C-V coordination were only available for tokens that contain a
15 vowel height target. The absence of the negative correlation between C1 and ICI for tokens that
16 lack a vowel height target suggests that at least some of these tokens lack sufficient specification
17 to enter into C-V coordination.

18

19 The analyses of tongue dorsum trajectories and of C-V timing provide converging support for
20 H4, the hypothesis that devoiced vowels in Japanese are optionally targetless. Change over time
21 in tongue dorsum height, the most salient dimension of Japanese /u/, approximated a linear
22 trajectory between flanking non-high vowels in some devoiced tokens. In these tokens, the
23 tongue dorsum did not rise from its position for the preceding non-high vowel to /u/. Linear
24 interpolation of TD height across VCuCV sequences alternated with another pattern. Some
25 tokens containing devoiced /u/ were indistinguishable from voiced /u/—that is, the tongue

1 dorsum rose from its position for the preceding non-high vowel to achieve a vowel height target
2 for /u/ before lowering again for the following non-high vowel.

3

4 Speakers differ in the degree to which they produced words without a vowel height target for
5 devoiced /u/, but all speakers produced /u/ without a lingual target in some words some of the
6 time. When the vowel height target was present (but not when it was absent), we observed a
7 negative correlation between C₁ duration and ICI, as predicted by C-V coordination (Table 1).
8 The absence of this correlation across tokens lacking a vowel height target suggests a different
9 coordination regime, i.e., C-C coordination. The tendency to produce a targetless vowel also
10 varies across words, possibly due to the consonantal environment in which /u/ occurs.
11 Devoicing also impacted the timing between flanking consonants. The interval between
12 consonants flanking /u/ (ICI) decreased with devoicing, a pattern which also varied significantly
13 across words. Notably, the effect of devoicing on ICI was greatest for /ʃutais:/ and /katsutoki/,
14 the words which, on average, showed the highest degree of height targetlessness.

15

16

Discussion

17 The primary purpose of the study was to examine the lingual articulation of devoiced vowels
18 in Japanese. One of our foci was the height dimension of /u/. We focused on /u/, as opposed to
19 /i/ (which also devoices in Tokyo Japanese), since there is no previous lingual articulatory data
20 on this vowel. We focused on the height dimension because height is the most salient
21 characteristic of Japanese /u/. Even for voiced variants, the degree of backness and rounding in
22 Tokyo Japanese /u/ is often reduced relative to how /u/ is defined by the IPA. At the outset, we
23 formulated four specific hypotheses based on previous research, which we tested in an EMA
24 experiment. The data were analyzed via Bayesian classification of DCT components fit to the
25 TD trajectories. One innovative aspect of our analysis is that we defined a category lacking a

1 vowel height target in terms of a linear trajectory between flanking vowels. This approach
2 allowed us to consider on a token-by-token basis whether the vowel was specified for a height
3 target. Results support the hypothesis that vowel height targets are optional in devoiced vowels.
4 When the devoiced vowels in our study were produced without a height target, the TD height
5 trajectory for /u/ followed a path of linear interpolation between flanking vowels. All speakers
6 produced at least some tokens that were classified as linear interpolations, although the
7 probability of such tokens varied across speakers and across words.

8
9 Recall that there are competing claims in past work about the environments that condition
10 devoicing vs. deletion (Kawakami 1971; Whang 2014). Comparing Table 2, which shows the
11 claims of past work, and Table 3, which shows the actual deletion probabilities, our results do
12 not match with either. In particular, both Kawakami (1971) and Whang (2014) argue for
13 devoicing, not deletion, in [ʃutaise:], but we found the highest probability of height
14 targetlessness in that environment. Across speakers, vowel height targets were identified most
15 often in /masutaa/ followed by /ϕusoku/ and then /katsutoki/ and /ʃutaise:/. These four items
16 differ in many ways which may contribute to this result, including lexical statistics, the
17 presence/absence of morphological boundaries, and aspects of the vowel and consonantal
18 environments of the target /u/. For example, although the target /u/ was always flanked by non-
19 high vowels, the precise vowel contexts were different across target items: /a-/u-/a/ in
20 /masuta/, /e-/u-/o/ in /ϕusoku/, /a-/u-/o/ in /katsutoki/, and /e-/u-/a/ in /ʃutaise:/. We note
21 that the /a-/u-/a/ in /masutaa/ requires the largest articulatory movement for /u/, and it is here
22 where we observe vowel height targets most frequently. It could be that the lingual articulatory
23 difference between voiced and voiceless vowels are conditioned by coarticulatory context such
24 that differences between voiced and devoiced vowels are minimized for more extreme

1 articulatory trajectories.¹⁰ It is also the case that /masutaa/ is the only recent loanword in the
2 study and that it derives from a word in English with a /st/ cluster, although this etymology did
3 not seem to encourage vowel target absence. We also note that the /e/ preceding target /u/ in
4 /e#ϕsoku/ and /e#ʃutaise:/ comes from the carrier phrase, while the preceding vowels in the
5 other items are internal to the target word. While we took steps to discourage insertion of a
6 phonological phrase boundary between the carrier phrase and the target word (see methods), as
7 this is known to influence vowel-to-vowel coarticulation (e.g., Cho, 2006), the presence of a
8 morpheme boundary may also affect timing across gestures (Cho, 2001; Lee-Kim, Davidson,
9 & Hwang, 2013)), leading to differences across items. Specifically, increased vowel overlap
10 across morphological boundaries could reduce the rise in TD trajectory for both voiced and
11 voiceless /u/, pushing both towards the linear interpolation trajectory that served as the basis
12 for our “target absent” classification. We also note that, as there is also a morpheme boundary in
13 /katsu#toki/, the three items with the greatest incidence of height targetlessness contained a
14 morpheme boundary within the target VCVCV trajectory. The flanking consonants may also
15 play a role in conditioning item specific differences. The degree to which vowels coarticulate
16 across consonants is known to be influenced by the degree of palatal contact such that a
17 consonant like /ʃ/ with a high degree of palate contact has a greater degree of coarticulatory
18 resistance than /t/ which has less contact with the palate (Recasens, 1989). Although
19 consonantal context was controlled across voiced-voiceless dyads, this could potentially be a
20 source of difference across dyads. Our analysis of the inter-consonantal interval (ICI) showed
21 that some combinations of consonants, /ʃ_t/ and /ts_t/ in particular, were pulled closer together
22 when intervening vowels are devoiced, regardless of targetlessness. It was in these contexts that

¹⁰ We would like to thank the editor, Taehong Cho, for this suggestion.

1 we also observed increased frequency of targetlessness. This suggests a possible chain of
2 causation whereby laryngeal reorganization dictated by devoicing pulls consonants closer
3 together which obliterates the intervening vowel. Although teasing apart the various factors that
4 may be conditioning item difference requires future studies, we find the role of flanking
5 consonants to be one of the most intriguing possibilities.

6

7 We acknowledge that our focus on height in the analysis of TD trajectories prevents us from
8 drawing any firm conclusions about complete vowel absence from the classification analysis
9 alone. As mentioned above, height is the most salient feature of /u/ in Japanese. Nevertheless,
10 linear interpolation on the height dimension, indicating a lack of vowel height target, does not
11 preclude phonetic specification in other dimensions. Past articulatory data from ultrasound and
12 MRI indicate that the highest position of the tongue for /u/ is rather central, c.f., the substantially
13 more posterior position observed for /o/; the labial component of /u/ has long been recognized
14 as different from similar instances of this phone in other languages (Vance, 2008). Neither the
15 labial nor backing components of /u/ involve particularly large movements, making it difficult
16 to discern differences between voiced and devoiced /u/.

17

18 Besides the possibility of lip compression or backness targets for /u/, there is, as mentioned in
19 the introduction, evidence that laryngeal gestures originating with voiceless consonants are
20 controlled, c.f., passive coarticulation, to yield vowel devoicing. This conclusion goes back at
21 least to Sawashima (1971),¹¹ who writes (pg 13): “The opening of the glottis for the medial [kt]
22 and [kk] were significantly larger than for [tt] and [kk] which lasted for approximately the same

¹¹ We would like to thank an anonymous reviewer for pointing us to this literature.

1 durations, and this fact shows that the glottal adjustments for devoicing of the vowel are not a
2 mere skipping of the phonatory adjustments for the vowel but a positive effort of widening of
3 the glottis for the devoiced vowel segment, even though there is no phonemic distinction
4 between the voiced and devoiced vowels.” Moreover, there is evidence from EPG that the
5 tongue maintains fricative-like contact with the palate during devoiced vowels, which may be
6 controlled to sustain devoicing. These aspects of devoiced vowel articulation likely persist (we
7 have no evidence to indicate that they don’t) even when the tongue does not move towards a
8 height target for /u/. For these reasons, we cannot equate the absence of a height target with the
9 absence of vowel. Nevertheless, the data strongly suggest that the height target is categorically
10 present or absent across tokens, even if other aspects of the devoiced vowel remain under
11 speaker control.

12

13 One theoretical consequence of our results is that they constitute a categorical but optional
14 pattern. As recently pointed out by Bayles, Kaplan & Kaplan (2016), “optionality” in many
15 studies can come from averaging over inter-speaker differences, and it is important to examine
16 whether a categorical pattern can show true optionality *within a speaker*. In the case of devoiced
17 vowels in Japanese, the average trajectory of devoiced /u/ would point to the misleading
18 conclusion that /u/ is reduced, because it averages over “full target” and “no target” tokens (see
19 Shaw & Kawahara, submitted-a for further details). Although this was not the main purpose of
20 this experiment, our finding supports the view expressed by Bayles et al. (2016) that indeed, a
21 categorical pattern, like French schwa deletion, can be optional within a speaker. Other
22 articulatory research has identified similarly optional patterns, including nasal place
23 assimilation in English (Ellis & Hardcastle, 2002) and place assimilation in Korean (Kochetov
24 & Pouplier, 2008). Revealing such patterns requires that the phonological status of each token
25 is assessed individually. The data from Bayles et al. (2016) draws from an already segmented

1 corpus, essentially trusting the judgements about when a vowel appears or does not appear.
2 Ellis & Hardcastle (2002) identified optional patterns of nasal place assimilation in English
3 through visual inspection of EMA and EPG data. They collected baseline data of a (lexically)
4 velar nasal followed by a velar stop to evaluate possible assimilation of a coronal nasal to a
5 following velar stop. Kochetov & Pouplier (2008) went beyond visual inspection but found
6 similar patterns. For their study of place assimilation in Korean, they established a criterion—
7 two standard deviations from the mean of the canonical category—which they used to classify
8 tokens as either fully articulated (within two standard deviations) or reduced. Our approach
9 goes one step further, as we explicitly define categories based both on the full target and on
10 phonetic interpolation (target absent scenario). Despite some differences in analytical method,
11 we found a similar type of variation as Bayles et al. (2016), i.e., within-item and, in many cases,
12 within speaker variation of a largely categorical nature. To these results we can add vowel
13 height specification (presence/absence) to the list of categorical but optional processes.

14

15 There are numerous theoretical frameworks capable of handling the optional presence/absence
16 of a vowel or even a particular vowel feature. These include varbrul (Guy, 1988) and exemplar
17 theory (Pierrehumbert, 2006) as well as generative models developed to handle categorical
18 variation, including Stochastic OT (Hayes & Londe, 2006) and (Noisy) Harmonic Grammar
19 (Coetzee & Kawahara, 2013; McPherson & Hayes, 2016). In contrast to other optional
20 phonological patterns to which these models have been applied, the Japanese vowel devoicing
21 case is of particular theoretical interest because of learnability issues. Most theoretical
22 approaches to variability proceed by matching frequencies in the input data. Devoicing
23 impoverishes the acoustic signature of tongue dorsum height, disrupting as well the auditory
24 feedback that the learner may receive from variation in their own production of /u/. Hence, the
25 degree to which a learner can match productions in their own speech to frequencies in the

1 ambient environment is limited in this case. Learners may be restricted to somatosensory
2 feedback from their own productions (Tremblay, Shiller, & Ostry, 2003). Relevant to the
3 learnability issue is the observation that infant-directed speech in Japanese contains vowel
4 devoicing at approximately the same rates as adult-directed speech (Fais, Kajikawa, Amano, &
5 Werker, 2010). To the extent that systematic patterns, e.g., conditioning environments for vowel
6 height targetlessness, are found in the population, they may emerge from analytical bias as
7 opposed to pattern matching (Moreton, 2008).

8

9 The continuity of phonetic measurements makes it natural to consider the possibility that a
10 gesture is categorically present but reduced in magnitude or coordinated in time with other
11 gestures such that the acoustic consequences of the gesture are attenuated (Browman &
12 Goldstein, 1992b; Iskarous, McDonough, & Whalen, 2012; Jun, 1996; Jun & Beckman, 1993;
13 Parrell & Narayanan, 2014). Particularly with vowels, gradient and gradual phonetic shift is
14 well-documented and is often treated as the primary mechanism of variation (e.g., Labov, Ash,
15 & Boberg, 2005; Wells, 1982 for comprehensive overviews). This underscores the importance
16 of deploying rigorous methods to support claims about the presence/absence of a gesture in the
17 phonetic signal of which the Japanese data offer a clear case. We emphasize at this point also,
18 that our method is useful in testing the general “phonetic underspecification analysis” (Keating,
19 1988). Several studies have argued that certain segments lack phonetic targets in some
20 dimensions (Cohn, 1993; Keating, 1988; Pierrehumbert & Beckman, 1988). By rigorously
21 relating the signal to both a full vowel and the linear interpolation hypothesis, the current
22 computational analysis offers a general approach that can be used to evaluate phonetic
23 underspecification.

24

1 The variable targetlessness of devoiced /u/ raises several new research questions. Firstly, we
2 have observed that targetless probabilities differ across words, but we have left the precise
3 conditioning environments to future study. The identity of the flanking consonants, the lexical
4 statistics of the target words, the informativity of the vowel, etc., are all possibilities. Due to the
5 small number of words recorded in this experiment we hesitate to speculate on which of these
6 factors (and to what extent) may influence targetlessness, but we plan to follow up with another
7 study that expands the number of words instantiating the consonantal environments reported
8 here. A second question is the syllabic status of consonant clusters that result from targetless
9 /u/. Matsui (2014) speculates that the preceding consonant forms an independent syllable, a
10 syllabic consonant, while Kondo (1997) argues that it is resyllabified into the following
11 syllable, forming a complex onset. A third question is the status of /i/, the other of the two high
12 vowels that are categorically devoiced in Tokyo Japanese. Our conclusion thus far is solely
13 about /u/, which is independently known to be variable in its duration (Kawahara & Shaw,
14 submitted-b), and may be more susceptible to coarticulation than /i/ (c.f., Recasens & Espinosa,
15 2009). At this point we have nothing to say about whether devoiced /i/ may also be targetless
16 in some contexts but hold this to be an interesting question for future research. Finally, the
17 observed shift from C-V to C-C coordination may bear on the broader theoretical issue of how
18 higher level structure, i.e., the mora in Japanese, relates to the temporal organization of
19 consonant and vowel gestures. The CV mora may be less determinate of rhythmic patterns of
20 Japanese than is sometimes assumed (Beckman 1982; Warner and Arai, 2001).

21

22

Conclusion

23 The current experiment was designed to address the question of whether devoiced /u/ in
24 Japanese is simply devoiced (Jun & Beckman 1993) or whether it is also targetless (Kondo

1 2001). Since previous studies on this topic have used acoustic data (Whang 2014) or
2 impressionistic observations (Kawakami 1977), we approached this issue by collecting
3 articulatory data. Using EMA, we recorded tongue dorsum trajectories in words containing
4 voiced and devoiced /u/. Some devoiced tokens showed a linear trajectory between flanking
5 vowels, indicating that there is no height target for the vowel. Other devoiced vowel tokens had
6 trajectories like voiced vowels. Tokens that were reduced, showing trajectories intermediate
7 between the voiced vowels and linear interpolation between flanking vowels, were less
8 common. We conclude that /u/ is optionally targetless; i.e., there is token-by-token variability
9 in whether the lingual gesture is present or absent. The patterns of covariation between C1
10 duration and ICI provided further support for this conclusion. For tokens classified as
11 containing a lingual target, there is a significant negative correlation between C1 duration and
12 ICI, a prediction of C-V timing. The correlation is in the opposite direction for tokens that lacks
13 a lingual height target for the vowel.

14

15 Achieving the above descriptive generalization—that devoiced vowels are optionally
16 targetless—required some methodological advancements, including analytical tools for
17 assessing phonological status on a token-by-token basis. We analyzed the data using Bayesian
18 classification of a compressed representation of the signal based on Discrete Cosine Transform
19 (following Shaw & Kawahara, submitted-a). The posterior probabilities of the classification
20 showed a bimodal distribution, supporting the conclusion that devoiced /u/ in Tokyo Japanese
21 is variably targetlessness.

22

23 Overall, establishing that devoiced /u/ tokens are sometimes targetless, in that they do not differ
24 from the linear interpolation of flanking gestures, answers a long-standing question in Japanese
25 phonetics/phonology while raising several new questions to pursue in future research, including

1 the syllabic status of consonant clusters flanking a targetless vowel, the role of the mora (or
2 lack thereof) in Japanese timing, the phonological contexts that favor targetless /u/, and whether
3 other devoiced vowels, particularly /i/, may also be variably targetless.

4

5 **Acknowledgements**

6 This project is supported by JSPS grant #15F15715 to the first and second authors, #26770147
7 and #26284059 to the second author. Thanks to audience at Yale, ICU, Keio, RIKEN, and
8 Phonological Association in Kansai and “Syllables and Prosody” workshop at NINJAL, in
9 particular Mary Beckman, Lisa Davidson, Junko Ito, Michinao Matsui, Reiko Mazuka, Armin
10 Mester, Haruo Kubozono, and three anonymous reviewers. All remaining errors are ours.

11

12 **References**

- 13 Aylett, M., & Turk, A. (2004). The smooth signal redundancy hypothesis: A functional
14 explanation for relationships between redundancy, prosodic prominence, and duration
15 in spontaneous speech. *Language and Speech*, 47(1), 31-56.
- 16 Aylett, M., & Turk, A. (2006). Language redundancy predicts syllabic duration and the spectral
17 characteristics of vocalic syllable nuclei. *Journal of Acoustical Society of America*,
18 119(5), 3048-3059.
- 19 Bates, D., Maechler, M., Bolker, B., & Walker, S. (2014). lme4: Linear mixed-effects models
20 using Eigen and S4. *R Package Version*, 1(7).
- 21 Bayles, A., Kaplan, A., & Kaplan, A. (2016). Inter-and intra-speaker variation in French schwa.
22 *Glossa: a journal of general linguistics*, 1(1).
- 23 Beckman, M. (1982). Segmental duration and the 'mora' in Japanese. *Phonetica*, 39, 113-135.
- 24 Beckman, M. (1986). *Stress and Non-Stress Accent*. Dordrecht: Foris.
- 25 Beckman, M., & Shoji, A. (1984). Spectral and Perceptual Evidence for CV Coarticulation in
26 Devoiced/si/and/syu/in Japanese. *Phonetica*, 41(2), 61-71.
- 27 Bell, A., Brenier, J. M., Gregory, M., Girand, C., & Jurafsky, D. (2009). Predictability effects
28 on durations of content and function words in conversational English. *Journal of*
29 *Memory and Language*, 60(1), 92-111.
- 30 Berry, J. J. (2011). Accuracy of the NDI wave speech research system. *Journal of Speech,*
31 *Language, and Hearing Research*, 54(5), 1295-1301.
- 32 Blackwood-Ximenes, A., Shaw, J., & Carignan, C. (2017). A comparison of acoustic and
33 articulatory methods for analyzing vowel variation across American and Australian
34 dialects of English. *The Journal of Acoustical Society of America*, 142(2), 363-377.
- 35 Bombien, L., Mooshammer, C., & Hoole, P. (2013). Articulatory coordination in word-initial
36 clusters of German. *Journal of Phonetics*, 41(6), 546-561.

- 1 Browman, & Goldstein, L. (1992a). 'Targetless' schwa: An articulatory analysis. In G. Docherty
2 & R. Ladd (Eds.), *Papers in Laboratory Phonology II: Gesture, Segment, Prosody* (pp.
3 26-56). Cambridge: Cambridge University Press.
- 4 Browman, C., & Goldstein, L. (1992b). Articulatory phonology: An overview. *Phonetica*, 49,
5 155-180.
- 6 Browman, C. P., & Goldstein, L. M. (2000). Competing Constraints on Intergestural
7 Coordination and Self-Organization of Phonological Structures. *Les cahiers de l'ICP,*
8 *Bulletin de la Communication Parlee*, 5, 25-34.
- 9 Byrd, D. (1996). Influences on Articulatory Timing in Consonant Sequences. *Journal of*
10 *Phonetics*, 24(2), 209-244.
- 11 Cedergren, H. J., & Simoneau, L. (1985). La chute des voyelles hautes en français de
12 Montréal: 'As-tu entendu la belle syncope?'. *Les tendances dynamiques du français*
13 *parlé à Montréal*, 1, 57-145.
- 14 Chitoran, I., Goldstein, L., & Byrd, D. (2002). Gestural overlap and recoverability: articulatory
15 evidence from Georgian. In C. Gussenhoven & N. Warner (Eds.), *Laboratory*
16 *Phonology* 7 (pp. 419-447). Berlin, New York: Mouton de Gruyter.
- 17 Cho, T. (2001). Effects of morpheme boundaries on intergestural timing: Evidence from
18 Korean. *Phonetica*, 58(3), 129-162.
- 19 Cho, T. (2004). Prosodically conditioned strengthening and vowel-to-vowel coarticulation in
20 English. *Journal of Phonetics*, 32(2), 141-176.
- 21 Coetzee, A. W., & Kawahara, S. (2013). Frequency biases in phonological variation. *Natural*
22 *Language & Linguistic Theory*, 31(1), 47-89.
- 23 Coetzee, A. W., & Pater, J. (2011). The place of variation in phonological theory. *The*
24 *Handbook of Phonological Theory, Second Edition*, 401-434.
- 25 Cohn, A. C. (1993). Nasalisation in English: phonology or phonetics. *Phonology*, 10(01), 43-
26 81.
- 27 Dauer, R. M. (1980). The reduction of unstressed high vowels in Modern Greek. *Journal of the*
28 *International Phonetic Association*, 10(1-2), 17-27.
- 29 Davis, C., Shaw, J., Proctor, M., Derrick, D., Sherwood, S., & Kim, J. (2015). Examining
30 speech production using masked priming. *18th International Congress of Phonetic*
31 *Sciences*.
- 32 Dupoux, E., Kakehi, K., Hirose, Y., Pallier, C., & Mehler, J. (1999). Epenthetic vowels in
33 Japanese: A perceptual illusion? *Journal of Experimental Psychology: Human*
34 *Perception and Performance*, 25, 1568-1578.
- 35 Eftychiou, E. (2010). Routes to lenition: an acoustic study. *PLoS One*, 5(3), e9828.
- 36 Ellis, L., & Hardcastle, W. J. (2002). Categorical and gradient properties of assimilation in
37 alveolar to velar sequences: evidence from EPG and EMA data. *Journal of Phonetics*,
38 30(3), 373-396.
- 39 Faber, A., & Vance, T. J. (2000). More acoustic traces of 'deleted' vowels in Japanese.
40 *Japanese/Korean Linguistics*, 9, 100-113.
- 41 Fais, L., Kajikawa, S., Amano, S., & Werker, J. F. (2010). Now you hear it, now you don't:
42 vowel devoicing in Japanese infant-directed speech. *Journal of Child Language*, 37(2),
43 319-340.
- 44 Fujimoto, M. (2015). Chapter 4: Vowel devoicing. In H. Kubozono (Ed.), *The handbook of*
45 *Japanese phonetics and phonology*. Berlin: Mouton de Gruyter.
- 46 Fujimoto, M., Murano, E., Niimi, S., & Kiritani, S. (2002). Differences in glottal opening
47 pattern between Tokyo and Osaka dialect speakers: factors contributing to vowel
48 devoicing. *Folia phoniatrica et logopaedica*, 54(3), 133-143.

- 1 Gafos, A. (2002). A grammar of gestural coordination. *Natural Language and Linguistic Theory*, 20, 269-337.
- 2
- 3 Gafos, A., & Goldstein, L. (2012). Articulatory representation and organization. *The handbook of laboratory phonology*, 220-231.
- 4
- 5 Gafos, A., Hoole, P., Roon, K., & Zeroual, C. (2010). Variation in timing and phonological grammar in Moroccan Arabic clusters. In C. Fougeron, B. Kühnert, M. d'Imperio, & N. Vallée (Eds.), *Laboratory Phonology* (Vol. 10, pp. 657-698). Berlin, New York: Mouton de Gruyter.
- 6
- 7
- 8
- 9 Gafos, A., Kirov, C., & Shaw, J. (2010). *Guidelines for using Mview*.
- 10 Garcia, D. (2010). Robust smoothing of gridded data in one and higher dimensions with missing values. *Computational statistics & data analysis*, 54(4), 1167-1178.
- 11
- 12 Goldstein, L., Nam, H., Saltzman, E., & Chitoran, I. (2009). Coupled oscillator planning model of speech timing and syllable structure. *Frontiers in phonetics and speech science*, 239-250.
- 13
- 14
- 15 Guy, G. (1988). Advanced VARBRUL analysis. *Linguistic change and contact*, 124-136.
- 16 Hall, K. C., Hume, E., Jaeger, F., & Wedel, A. (2016). *The Message Shapes Phonology*.
- 17 Harrington, J., Kleber, F., & Reubold, U. (2008). Compensation for coarticulation, /u/-fronting, and sound change in standard southern British: An acoustic and perceptual study. *The Journal of the acoustical society of America*, 123(5), 2825-2835.
- 18
- 19
- 20 Hayes, B., & Londe, Z. C. (2006). Stochastic phonological knowledge: The case of Hungarian vowel harmony. *Phonology*, 23(1), 59-104.
- 21
- 22 Hirayama, M. (2009). Postlexical Prosodic Structure and Vowel Devoicing in Japanese (2009). *Toronto Working Papers in Linguistics*.
- 23
- 24 Hirose, H. (1971). The activity of the adductor laryngeal muscles in respect to vowel devoicing in Japanese. *Phonetica*, 23(3), 156-170.
- 25
- 26 Imaizumi, S., & Hayashi, A. (1995). Listener-adaptive adjustments in speech production: evidence from vowel devoicing. *Ann. Bull. RIPL*, 29, 43-48.
- 27
- 28 Iskarous, K., McDonough, J., & Whalen, D. (2012). A gestural account of the velar fricative in Navajo.
- 29
- 30 Isomura, K. (2009). *Nihongo-wo Oshieru [Teaching Japanese]*. Tokyo Hitsuji.
- 31 Johnson, K., Ladefoged, P., & Lindau, M. (1993). Individual differences in vowel production. *Journal of Acoustical Society of America*, 94(2), 701-714.
- 32
- 33 Jun, J. (1996). Place assimilation is not the result of gestural overlap: evidence from Korean and English. *Phonology*, 13(03), 377-407.
- 34
- 35 Jun, S.-A., & Beckman, M. (1993). *A gestural-overlap analysis of vowel devoicing in Japanese and Korean*. Paper presented at the 67th annual meeting of the Linguistic Society of America, Los Angeles.
- 36
- 37
- 38 Jun, S.-A., Beckman, M. E., & Lee, H.-J. (1998). Fiberscopic evidence for the influence on vowel devoicing of the glottal configurations for Korean obstruents. *UCLA Working Papers in Phonetics*, 43-68.
- 39
- 40
- 41 Jurafsky, D., Bell, A., Gregory, M., & Raymond, W. D. (2001). Probabilistic relations between words: Evidence from reduction in lexical production. *Typological studies in language*, 45, 229-254.
- 42
- 43
- 44 Kawahara, S. (2015). A catalogue of phonological opacity in Japanese: Version 1.2. *慶応義塾大学言語文化研究所紀要*(46), 145-174.
- 45
- 46 Kawakami, S. (1977). Outline of Japanese Phonetics [written in Japanese as "Nihongo Onsei Gaisetsu"]: Tokyo: Oofuu-sha.
- 47
- 48 Keating, P. (1988). Underspecification in phonetics. *Phonology*, 5, 275-292.

- 1 Kilbourn-Ceron, O., & Sonderegger, M. (2017). Boundary phenomena and variability in
2 Japanese high vowel devoicing. *Natural Language & Linguistic Theory*, 1-43.
- 3 Kingston, J. (1990). Articulatory binding. In J. Kingston & M. Beckman (Eds.), *Papers in*
4 *Laboratory Phonology I* (pp. 406-434). Cambridge: Cambridge University Press.
- 5 Kochetov, A., & Kang, Y. (2017). Supralaryngeal implementation of length and laryngeal
6 contrasts in Japanese and Korean. *Canadian Journal of Linguistics/Revue canadienne*
7 *de linguistique*, 62(1), 18-55.
- 8 Kochetov, A., & Pouplier, M. (2008). Phonetic variability and grammatical knowledge: An
9 articulatory study of Korean place assimilation. *Phonology*, 25(3), 399-431.
- 10 Kondo, M. (1997). Mechanisms of vowel devoicing in Japanese.
- 11 Kondo, M. (2001). Vowel Devoicing and Syllable Structure in Japanese. In M. Nakayama &
12 C. J. Quinn (Eds.), *Japanese/Korean Linguistics* (Vol. 9). Stanford: CSLI.
- 13 Kondo, M. (2005). Syllable structure and its acoustic effects on vowels in devoicing
14 environments. *Voicing in Japanese*, 84, 229.
- 15 Kuriyagawa, F., & Sawashima, M. (1989). Word accent, devoicing and duration of vowels in
16 Japanese. *Annual Bulletin of the Research Institute of Language Processing*, 23, 85-
17 108.
- 18 Labov, W., Ash, S., & Boberg, C. (2005). *The atlas of North American English: Phonetics,*
19 *phonology and sound change*: Walter de Gruyter.
- 20 Lee-Kim, S.-I., Davidson, L., & Hwang, S. (2013). Morphological effects on the darkness of
21 English intervocalic/l. *Laboratory Phonology*, 4(2), 475-511.
- 22 Lindblom, B. (1990). Explaining phonetic variation: A sketch of the H&H theory *Speech*
23 *production and speech modelling* (pp. 403-439): Springer.
- 24 Maekawa, K. (1990). *Production and perception of the accent in the consecutively devoiced*
25 *syllables in Tokyo Japanese*. Paper presented at the ICSLP.
- 26 Maekawa, K., & Kikuchi, H. (2005). Corpus-based analysis of vowel devoicing in spontaneous
27 Japanese: an interim report. In J. v. d. Weijer, K. Nanjo, & T. Nishihara (Eds.), *Voicing*
28 *in Japanese* (pp. 205-228). Berlin, New York: Mouton de Gruyter.
- 29 Maekawa, K., Yamazaki, M., Ogiso, T., Maruyama, T., Ogura, H., Kashino, W., . . . Den, Y.
30 (2014). Balanced corpus of contemporary written Japanese. *Language resources and*
31 *evaluation*, 48(2), 345.
- 32 Marin, S. (2013). The temporal organization of complex onsets and codas in Romanian: A
33 gestural approach. *Journal of Phonetics*, 41(3), 211-227.
- 34 Marin, S., & Pouplier, M. (2010). Temporal organization of complex onsets and codas in
35 American English: Testing the predictions of a gesture coupling model. *Motor Control*,
36 14, 380-407.
- 37 Matsui, M. (2014). Vowel devoicing, VOT Distribution and Geminate Insertion of Sibilants 歯
38 擦音の母音無声化・VOT 分布・促音挿入. *Theoretical and applied linguistics at*
39 *Kobe Shoin: トークス*, 17, 67-106.
- 40 McPherson, L., & Hayes, B. (2016). Relating application frequency to morphological structure:
41 the case of Tommo So vowel harmony. *Phonology*, 33(1), 125-167.
- 42 Moreton, E. (2008). Analytic bias and phonological typology. *Phonology*, 25(01), 83-127.
- 43 Munhall, K., & Lofqvist, A. (1992). Gestural aggregation in speech: laryngeal gestures. *Journal*
44 *of Phonetics*, 20(1), 111-126.
- 45 Nakamura, M. (2003, August 3-9). *The Spatio-temporal Effects of Vowel Devoicing on Gestural*
46 *Coordination: an EPG study*. Paper presented at the 15th International Congress of
47 Phonetics Sciences, Barcelona.

- 1 Nielsen, K. Y. (2015). Continuous versus categorical aspects of Japanese consecutive
2 devoicing. *Journal of Phonetics*, 52, 70-88.
- 3 Nogita, A., Yamane, N., & Bird, S. (2013). *The Japanese unrounded back vowel/w/ is in fact*
4 *rounded central/front [ɰ-ɾ]*. Paper presented at the Ultrafest VI, Edinburgh.
- 5 Parrell, B., & Narayanan, S. (2014). *Interaction between general prosodic factors and*
6 *languagespecific articulatory patterns underlies divergent outcomes of coronal stop*
7 *reduction*. Paper presented at the International Seminar on Speech Production (ISSP)
8 Cologne, Germany.
- 9 Pierrehumbert, J. (2006). The next toolkit. *Journal of Phonetics*, 34(6), 516-530.
- 10 Pierrehumbert, J., & Beckman, M. (1988). *Japanese Tone Structure*. Cambridge, Mass.: MIT
11 Press.
- 12 Pierrehumbert, J. B. (2001). Stochastic phonology. *Glott International*, 5(6), 195-207.
- 13 Poser, W. J. (1990). Evidence for foot structure in Japanese. *Language*, 66, 78-105.
- 14 Recasens, D. (1989). Long range coarticulation effects for tongue dorsum contact in VCVCV
15 sequences. *Speech Communication*, 8(4), 293-307.
- 16 Recasens, D., & Espinosa, A. (2009). An articulatory investigation of lingual coarticulatory
17 resistance and aggressiveness for consonants and vowels in Catalan. *The Journal of the*
18 *acoustical society of America*, 125(4), 2288-2298.
- 19 Sawashima, M. (1971). Devoicing of vowels. *Annual Bulletin of the Research Institute of*
20 *Logopedics and Phoniatics*, 5, 7-13.
- 21 Shaw, J. A., Chen, W.-r., Proctor, M. I., & Derrick, D. (2016). Influences of tone on vowel
22 articulation in Mandarin Chinese. *Journal of Speech, Language, and Hearing Research*,
23 59(6), S1566-S1574.
- 24 Shaw, J. A., Gafos, A., Hoole, P., & Zeroual, C. (2011). Dynamic invariance in the phonetic
25 expression of syllable structure: a case study of Moroccan Arabic consonant clusters.
26 *Phonology*, 28(3), 455-490.
- 27 Shaw, J. A., & Gafos, A. I. (2015). Stochastic Time Models of Syllable Structure. *PLoS One*,
28 10(5), e0124714.
- 29 Shaw, J. A., Gafos, A. I., Hoole, P., & Zeroual, C. (2009). Syllabification in Moroccan Arabic:
30 evidence from patterns of temporal stability in articulation. *Phonology*, 26, 187-215.
- 31 Shaw, J. A., & Kawahara, S. (submitted-a). A computational toolkit for assessing phonological
32 specification in phonetic data: Discrete Cosine Transform, Micro-Prosodic Sampling,
33 Bayesian Classification. *Phonology*, 34.
- 34 Shaw, J. A., & Kawahara, S. (submitted-b). Effects of entropy and surprisal on vowel duration
35 in Japanese. *Language and Speech*.
- 36 Sjoberg, A. F. (1963). *Uzbek structural grammar* (Vol. 18): Indiana university.
- 37 Smith, C. L. (1995). Prosodic patterns in the coordination of consonant and vowel gestures. In
38 B. Connell & A. Arvaniti (Eds.), *Papers in laboratory phonology IV: phonology and*
39 *phonetic evidence* (pp. 205-222). Cambridge: Cambridge University Press.
- 40 Smith, C. L. (2003). Vowel devoicing in contemporary French. *Journal of French Language*
41 *Studies*, 13(02), 177-194.
- 42 Sugito, M., & Hirose, H. (1988). Production and perception of accented devoiced vowels in
43 Japanese. *Annual Bulletin of Research Institute of Logopedics and Phoniatics*, 22, 19-
44 36.
- 45 Tiede, M. (2005). MVIEW: software for visualization and analysis of concurrently recorded
46 movement data. New Haven, CT: Haskins Laboratories.
- 47 Tremblay, S., Shiller, D. M., & Ostry, D. J. (2003). Somatosensory basis of speech production.
48 *Nature*, 423(6942), 866.

- 1 Tsuchida, A. (1997). *Phonetics and phonology of Japanese vowel devoicing*. Ph. D.
 2 Dissertation: University of Cornell.
 3 Vance, T. (1987). *An Introduction to Japanese Phonology*. New York: SUNY Press.
 4 Vance, T. J. (2008). *The sounds of Japanese with audio CD*: Cambridge University Press.
 5 Warner, N., & Arai, T. (2001). Japanese mora-timing: A review. *Phonetica*, 58(1-2), 1-25.
 6 Wells, J. C. (1982). *Accents of English, vol.2: The British Isles*. Cambridge: Cambridge
 7 University Press.
 8 Whang, J. (2014). Effects of predictability on vowel reduction. *Journal of Acoustical Society*
 9 *of America*, 135(4), 2293.
 10 Wright, R. (1996). *Consonant Clusters and Cue Preservation in Tsou*. (PhD dissertation),
 11 UCLA, Los Angeles.
 12 Yip, J. C.-K. (2013). *Phonetic effects on the timing of gestural coordination in Modern Greek*
 13 *consonant clusters*. University of Michigan.
 14

15 **Appendix**

16 The number of consonant gestures by condition parsed by tangential vs. component velocities.

		C1		C2	
		voiceless	voiced	voiceless	voiced
ϕusoku~ϕuzoku	tangential velocity	68	70	63	66
	component velocity	1	1	6	5
ʃutaise:~ʃudaika	tangential velocity	31	42	50	24
	component velocity	36	28	20	43
katsutoki~katsudo:	tangential velocity	50	55	50	50
	component velocity	23	11	20	19
hakusai~yakuzai	tangential velocity	50	55	70	67
	component velocity	23	11	0	3
masutaa~masuda	tangential velocity	70	29	70	35
	component velocity	0	10	0	4

17