

Recoverability-conditioned sensitivity to phonetic detail

James Whang

MARCS Institute for Brain, Behaviour & Development

ARC Centre of Excellence for the Dynamics of Language

research@jameswhang.net

Abstract

Japanese speakers systematically devoice or delete high vowels [i, u] between two voiceless consonants. Japanese listeners also report perceiving the same high vowels between consonant clusters even in the absence of a vocalic segment. Although perceptual vowel epenthesis has been described primarily as a phonotactic repair strategy, where a phonetically minimal vowel is epenthesized by default, few studies have investigated how the predictability of a vowel in a given context affects the choice of epenthetic vowel. The present study uses a forced-choice labeling task to test how sensitive Japanese listeners are to coarticulatory cues of high vowels [i, u] and non-high vowel [a] in devoicing and non-devoicing contexts. Devoicing contexts were further divided into high-predictability contexts, where only one of the high vowels is phonotactically legal, and low-predictability contexts, where both high vowels are allowed, to specifically test for the effects of predictability. Results reveal a strong tendency towards [u] epenthesis as previous studies have found, but the results also reveal a sensitivity to coarticulatory cues that override the default [u] epenthesis, particularly when the vowel is in a context that is less predictable. Recoverability-conditioned gestural coordination has been proposed for production, and this study provides evidence that perception is conditioned by recoverability as well.

Key words: perceptual repair, recoverability, phonotactics, Japanese

1. Introduction

Recoverability refers to a listener’s ease in accessing speech units that a speaker intended to transmit (e.g., [kæt̚, kæt̚^h] → /kæt̚/ ‘cat’; Mattingly, 1981; McCarthy, 1999; Chitoran et al., 2002). The recoverability of a speech unit is affected both by the perceptibility of its cues in the signal and its predictability in a given context. When neither perceptibility nor predictability are sufficient, recovery fails. The current study investigates Japanese listeners and the role of recoverability in how illicit consonant clusters are repaired. While it is commonly thought that Japanese listeners use [u] epenthesis by default because it is the shortest (and thus phonetically minimal) vowel (Dupoux et al., 1999, 2011), the current study proposes that the choice of epenthetic vowel in Japanese listeners rely on a combination of phonotactic predictability and attention to phonetic cues, based on experience with recovering high vowels that are systematically devoiced or deleted in their language. To distinguish the respective roles of phonotactic prediction and perception of phonetic cues, participants are presented with conflicting phonotactic and phonetic information, allowing insight into which of the two they prioritize.

1.1. *Effects of recoverability during production*

Previous studies have shown that phonetic cues are often enhanced in less perceptible contexts to aid recoverability. Chitoran et al. (2002) used EMMA (electromagnetic midsagittal articulometer) to investigate how gestural overlap of Georgian stop clusters are affected by the cluster’s position in word and how their places of articulation are ordered. The authors found that word-initial clusters show decreased gestural overlap; they also found that stop clusters are less overlapped

when the consonants ordered back-to-front in terms of place (e.g. /tp, kp, kt/) than when ordered front-to-back (e.g. /pt, pk, tk/). Both of these results are interpreted to show that gestural overlap is decreased when recoverability is at risk – word-initially because the only cues for C₁ comes from its release, and for a similar reason in back-to-front ordering of place since early closure of C₂ in a sequence like /tp/ would mask the release of C₁ (Browman & Goldstein, 2000). Silverman (1997) observes that a similar effort to maximize perceptibility of phonetic cues can be found when entire sound systems are considered as well, where languages prefer segments that have maximally contrastive cues. When a sound system contains seemingly inefficient contrasts, it is to preserve other efficient contrasts. For example, Mazatec is a language that uses aspiration contrastively with unaspirated, post-aspirated, and pre-aspirated stops in its phoneme inventory. The language additionally has a complex vowel system with modal, nasal, breathy, and creaky vowels. With the observation that breathy vowels often surface approximately as [aa], starting out breathy then becoming modal towards the end, Silverman notes that sequences such as [taa] where unaspirated stops are followed by (partially) breathy vowels are unattested in the language, presumably because it is perceptually too similar to attested sequences such as [t^ha], where post-aspirated stops are followed by modal vowels. Following the argument of Bladon (1986), Silverman proposes that the attested sequence with post-aspiration is preferred because it yields maximal perceptibility (and hence recoverability) of the laryngeal abduction that corresponds to aspiration.

It has also been shown that phonetic cues are weakened for segments that are predictable from a given context. Exemplar-based approaches to phonology (Bybee, 2006; Ernestus, 2011; Pierrehumbert, 2001) have long noted that it is often the most frequent lexical items that are targeted for devoicing due to their predictability. Building on this line of research, Hall et al. (in press) argue that phonological systems tend to reduce segments in predictable and/or perceptually weak positions because enhancing the cues would require additional effort while contributing little to successful lexical access (by the listener). For example, word-final coda contrasts are often neutralized cross-linguistically for two reasons. First, during lexical access, segments become more predictable as the listener processes more and more of the target item. This means that word-final codas contribute less to identifying the target lexical item. Second, codas are less perceptible than onset consonants. Rather than enhancing the weak cues of an already predictable segment, phonological systems choose to enhance cues of segments in perceptually strong and/or unpredictable positions instead, such as in the case of word-initial obstruent aspiration in English (e.g. /pik/ → [p^hik] vs. /spik/ → [spik]).

1.2. *Recoverability in perception*

If it is the case that speakers are varying the amount of phonetic cues depending on the target segment’s perceptibility or predictability in a given context, the question that naturally follows is whether listeners similarly vary their attention to phonetic cues based on context. Much of the studies pertaining to recoverability focus on production, under the assumption that the observed modulation of phonetic cues happen for the benefit of the listener. While it is true that segments

carrying more disambiguating information for words reduce less (van Son & Pols, 2003b, 2003a), it has also been shown that speakers reduce the second mention of a word even if the listener is hearing the word for the first time (Bard et al., 2000), leading to a less clear picture of speakers' concern for the listener.

So what is known about what listeners do in their recovery process? Listeners are attuned to contrasts that are native to their language. For example, Korean listeners are more sensitive to V-to-C formant transition cues than English listeners (Hume et al., 1999), because coda obstruents are obligatorily unreleased in Korean while they are optionally released in English (Kang, 2003), making the transitional cue more useful to Korean listeners for recovery of the coda consonant than to English listeners, who have the option of waiting for the release of the coda obstruent. Conversely, listeners are also often insensitive to phonetic cues that are not contrastive in their native language. For example, French listeners have difficulty contrasting short versus long vowels (Dupoux et al., 1999), English listeners have difficulty perceiving tonal contrasts (So & Best, 2010), Japanese listeners have difficulty contrasting /l/ versus /r/ because neither are phonemes of the language (Flege et al., 1996), and so on.

Selective sensitivity to phonetic cues show that recoverability is affected by expectations stemming from the listener's language experience. A study by Pitt and McQueen (1998) showed that listeners are biased towards identifying phonetically ambiguous segments as segments with higher phonotactic probability (i.e., phonotactically more predictable). Phonotactic knowledge also seems to play an important role in other domains as well. When processing nonce words, sequences with higher phonotactic probabilities are processed faster (Vitevitch et al., 1997; Vitevitch & Luce, 1999). On the other hand, lexical items with high phonotactic probabilities are processed slower than lexical items with low phonotactic probabilities, presumably because high phonotactic probability in lexical items means that there are also that many more similar lexical items, ultimately slowing down lexical access (Marslen-Wilson, 1987; McClelland & Elman, 1986; Norris, 1994; Vitevitch & Luce, 1998). The question then is, which process takes precedence? The answer to this question seems to depend on the task. In general, listeners seem to prioritize the use of lexical knowledge, relying on their phonotactic knowledge only when lexical activation fails (Shademan, 2006; Vitevitch & Luce, 1999). Additionally, a study by (Mattys et al., 2005) investigated whether participants pay more attention to lexical and sublexical (segmental and prosodic) segmentation cues when they are in conflict. The results showed again that lexical cues are prioritized and that listeners rely on sublexical cues when lexical context or information cannot be accessed due to noise or absence. The current study adds to this line of work by investigating the interaction between two sublexical information, namely phonotactic predictability and very fine-grained phonetic cues, using high vowel epenthesis in Japanese as a test case.

1.3. Perceptual recovery by Japanese listeners

In the now well-known study commonly referred to as the “*ebzo* test” (Dupoux et al., 1999), French and Japanese speakers were presented with acoustic stimuli with the high back rounded vowel [u] of

varying durations ranging from 0 ms to 90 ms occurring between two consonants (e.g., [ebzo] → [ebu:zo]). The stimuli were designed so that when there is no vowel in the stimuli, the result is a phonotactically legal sequence in French but illegal in Japanese. Their results showed that while French speakers could accurately distinguish the vowel-less from the vowel-ful tokens, Japanese speakers were essentially “deaf” to such differences, erring heavily towards misperceiving—or for the purposes of the current study, mistakenly *recovering*—what the authors call an “illusory” vowel. On the other hand, French speakers were unable to accurately perceive vowel length, with which the Japanese participants had little trouble perceiving. The authors propose that the results are due to phonotactic differences in French and Japanese, where Japanese listeners perceive a non-existent vowel between two consonants because Japanese phonotactics disallows heterorganic consonant clusters. French listeners, on the other hand, are insensitive to vowel length because it is not contrastive in French. The authors further argue that there is a “top-down” phonotactic effect on perception, where phonotactically illegal sequences are automatically perceived as the nearest legal sequence rather than repaired at a higher, abstract phonological level.

Dehaene-Lambertz et al. (2000) also tested the illusory vowel epenthesis effect in an event-related potential (ERP) study. In this study, Dehaene-Lambertz et al. carried out experiments similar to that of Näätänen et al. (1997), where electrophysiological responses have been shown to be sensitive to phoneme categories. Dehaene-Lambertz et al. looked at how mismatch negativity (MMN) responses in Japanese and French speakers differ in the absence versus presence of a vowel in the same kind of sequences as those in Dupoux et al. (1999). The experiments followed an oddball paradigm where in one trial a sequence that is legal in both languages was presented as the standard (e.g., [igumo]) and one that is illegal only in Japanese as the deviant (e.g., [igmo]). The reverse was presented in a separate trial. Although the results reported collapsed the trials, the ERP results generally showed that Japanese speakers are insensitive to the differences between the vowel-ful and vowel-less items, while French speakers are, supporting the behavioral results from the original study by Dupoux et al. (1999). A related fMRI study by Jacquemot et al. (2003), also found similar but slightly weaker results. Jacquemot et al. report that in an AAX task (A-stimulus presented twice before X-stimulus), neural activity increased whenever the X stimulus was different from the A stimulus for both Japanese and French participants. This was true regardless of whether or not the acoustic difference was phonologically contrastive in the language, although neural activation was significantly greater when the acoustic contrasts were also phonologically contrastive.

A more recent study by (Dupoux et al., 2011) aimed to further bolster the automatic perceptual repair idea by also investigating European Portuguese, Brazilian Portuguese, and Japanese listeners. The reason for choosing the two dialects of Portuguese was that European Portuguese allows the same types of clusters as French, but Brazilian Portuguese has a strict CVCV phonotactic structure, leading to the expectation that their perception would be similar to that of Japanese listeners. The crucial difference between Brazilian Portuguese and Japanese is that in the former, the default epenthetic vowel is reported to be /i/ as opposed to the Japanese /u/. Since the quality of the epenthetic vowels are different in the two epenthesizing languages, the experiments were modified

slightly from the 1999 study to enable identification of the perceived illusory vowels in the results. Like French listeners, the results showed that European Portuguese listeners did not have trouble distinguishing vowel-less from vowel-ful tokens. Japanese listeners, again, showed a tendency towards mistakenly recovering /u/ between consonant clusters. The results, however, additionally showed that Japanese listeners were also sensitive to [i]-coarticulation in the first consonant (i.e., *eb^jzo*), recovering /i/ rather than /u/. By comparison, Brazilian Portuguese listeners tended to perceptually recover /i/ between illegal consonant clusters by default as expected, but did not show the same degree of sensitivity to [u]-coarticulation. Although the reasons for the disparity in sensitivity to coarticulatory cues were not discussed, the difference is likely due to Brazilian Portuguese listeners having little experience with a systematic high vowel devoicing process, leading them to underutilize coarticulatory cues relative to Japanese listeners.

1.4. *Problems and solutions*

The series of studies discussed above collectively suggest that there is a top-down imposition of the listeners native phonotactic grammar during perception. The experiments, however, would benefit from two particular refinements when considering Japanese listeners: using stimuli that are less foreign to Japanese listeners and controlling for the effects of high vowel devoicing in how Japanese listeners perceive certain consonant clusters. First, the waveform and spectrogram examples of the stimuli used in the studies by Dupoux and colleagues reveal that the burst of C₁ (e.g., [b] in [ebzo, ebuzo]) were rather long, potentially biasing the participants to perceive a vowel. For example, Dupoux et al. (2011) shows that in a sequence like [agno], the voiced stop had a burst of at least 50 ms and contained formant-like structures. Japanese voiced stops, however, typically have burst durations of less than 20 ms (Kong et al., 2012). In addition, Japanese high vowels are inherently short, with an average duration of approximately 40 ms, but they can be as short as 20 ms (Han, 1994; Beckman, 1982). Taking the short burst and vowel durations of Japanese together, an atypically long burst with formant structures can be interpreted as containing a vowel, possibly confounding the independent effects of acoustic cues and phonotactic violations (Wilson et al., 2014). Furthermore, the stop closure of the stop is also nearly 100 ms, which is closer to the geminate range than the singleton range in Japanese (Kawahara, 2006). Geminate consonants are not known to affect high vowel devoicing in C₁ position, but geminate consonants in C₂ position have been shown to increase the likelihood of preceding vowels being phonated in Japanese regardless of whether the consonants are voiced (Maekawa & Kikuchi, 2005; Fujimoto, 2015). This means that stimuli with geminate-like obstruents in both C₁ and C₂ positions (e.g., [igba]) could have further biased Japanese participants toward expecting a vowel in the target context. While this is also a tendency that is language-specific and phonotactically driven, it is unclear whether the primary driving force behind perceptual epenthesis is the heterorganic clusters, the phonetic cues of geminate-like segments, or a combination of both.

Second, the stimuli used in the *ebzo* tests included a mix of environments in which high vowel devoicing is expected to occur in Japanese (i.e., when between two voiceless obstruents) as well as

non-devoicing environments. The results reported in these studies, however, make no distinction between the two types of environments. High vowel devoicing is a process, where high vowels lose at the least their phonation and at most delete completely (Shaw & Kawahara, 2018), and Japanese listeners’ life-long experience in recovering the devoiced vowels from between consonant cluster-like sequences is extremely relevant. Because high vowel devoicing is highly productive, it is very likely that this phonological process had an effect that is independent of phonotactic constraint violations in creating an expectation for a vowel. The most straightforward remedy is to test and analyze the two environments separately (e.g., [ezpo] vs. [espo]). Furthermore, the devoicing stimuli group can be divided into a low-predictability sub-group, where both high vowels can occur (e.g., [efpo] \rightarrow /efipo, efupo/), and a high-predictability sub-group, where only one of the high vowels is phonotactically legal (e.g., [espo] \rightarrow /esupo, *esipo/). Varden (2010) states what seems to be a prevalent assumption in the literature on Japanese high vowels, which is that since high vowels trigger allophonic variation for /t, s, h/ in the language (i.e., /t/ \rightarrow [tʃi, tsu]; /s/ \rightarrow [ʃi, su]; /h/ \rightarrow [çi, φu]), high vowels in these contexts are easily recoverable even if the vowel is phonetically deleted. To give a concrete example, [tski] ‘moon’ would be analyzed as /tuki/, because [ts] can only occur as an allophone of /t/ preceding /u/. In other words, only one high vowel is phonotactically legal in such contexts, so listeners can simply predict the high vowel that follows from the consonant alone. If true, then this would also mean that in other devoicing contexts, where both high vowels are possible (i.e., /p, k, ʃ/ \rightarrow [pi, pu, ki, ku, ʃi, fu]), Japanese listeners would be more inclined to pay closer attention to the phonetic cues. This idea, however, has never been tested systematically. The current study, therefore, presents a perception experiment that specifically controls phonotactic predictability and investigates how much more or less sensitive to coarticulatory cues Japanese listeners are depending on the phonotactic predictability of the vowel.

2. Materials and methods

There are three main ways in which phonotactics and phonetic cues are likely to interact: (i) bottom-up, where listeners prioritize phonetic cues, relying on phonotactic knowledge only when phonetic cues are insufficiently perceptible, (ii) top-down, where listeners prioritize phonotactic knowledge, relying on phonetic cues only when phonotactic predictability is insufficient, and (iii) listeners simply stick to either phonotactics or phonetic cues. The stimuli are in the form $V_1C_1(V_T)C_2V_2$, where V_T is the target vowel and C_1 and C_2 are determined based on the stimulus group the token belongs to. The stimuli were divided into three groups: non-devoicing (No-Devoice) where vowel devoicing is not expected, low predictability (Lo-Predict) where both high vowels can occur and devoice, making coarticulatory cues necessary for recovery of a devoiced vowel, and high predictability (Hi-Predict) where phonotactic predictability is sufficient for recovery of a devoiced vowel, making coarticulatory cues less important. Below in Table 1 are the stimuli.

Table 1: Stimuli for Experiment 2.

<i>No-Devoice</i>	eb_ko	ez_po	eg_to	ob_ke	oz_pe	og_te
<i>Lo-Predict</i>	ep_ko	ef_po	ek_to	op_ke	of_pe	ok_te
<i>Hi-Predict</i>	eφ_ko	es_po	eç_to	oφ_ke	os_pe	oç_te

There were 252 stimulus items in total. The stimulus forms shown in Table 1 were first recorded by a trained, non-Japanese-speaking, English-Hungarian bilingual phonetician in a sound-attenuated booth with stress on the initial vowel and with /i, u, a/ as target vowels (V_T). /a/ was included as a target vowel because it is a low vowel that typically does not devoice in Japanese, and also to test whether Japanese listeners are sensitive to coarticulatory cues of all vowels or just high vowels. Attempts were made to record the stimuli with two native Japanese speakers, but both speakers had difficulties keeping high vowels voiced in devoicing contexts, and even when they were successful in producing voiced high vowels in devoicing contexts, either the burst durations were too short to manipulate or the target vowel was stressed.

For each recording, the target vowels were manipulated by inserting or removing whole periods to achieve a duration of $\sim 40 \pm 5$ ms. From each of the recordings, four additional tokens were created by removing from right to left, half of V_T (splice-1), the remaining half of V_T (splice-2), half of the C_1 burst/frication noise (splice-3), then the remaining half of the C_1 burst/frication noise leaving only the closure for stops and ~ 15 ms for fricatives (splice-4). An example of how the splicing was done is shown in Figure 1 below with the token [ekuto].

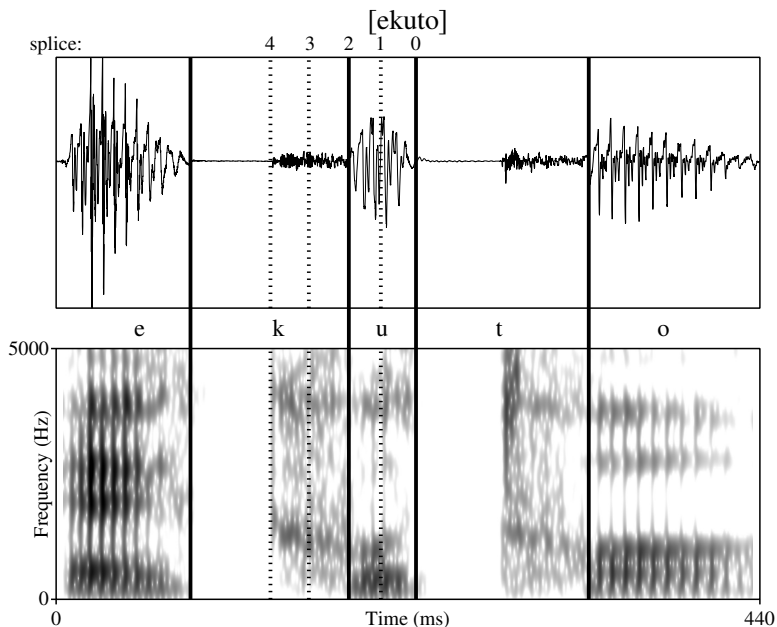


Figure 1: Example of token splicing: [ekuto].

The result of the splicing process is a gradual decrease of vowel coarticulatory information available

in the burst/frication noise of C_1 . Stop bursts in particular were manipulated to test whether it is phonotactic predictability or interpretation of phonetic information that drive illusory vowel epenthesis, since sensitivity to and interpretation of stop bursts as signaling the presence of a vowel is reported not just in Japanese (Furukawa, 2009; Whang, 2016) but in Korean (Kang, 2003) and English (Davidson & Shaw, 2012; Hsieh, 2013) as well.

Naturally produced, vowel-less tokens were also recorded for each stimulus form to test how it differs in perception from the spliced vowel-less stimuli (splice-2), which have traces of coarticulation from the target vowel on the surrounding consonants.

2.1. *Participants*

Twenty-nine monolingual Japanese listeners (16 women, 13 men) were recruited for the perception experiment in Tokyo, Japan. All participants were undergraduate students born and raised in the greater Tokyo area and were between the ages 18 and 24. Although all participants learned English as a second language as part of their compulsory education, none had resided outside of Japan for more than six months and have not been overseas within a year prior to the experiment. All participants were compensated for their time.

2.2. *Procedure*

The experiment follows the forced-choice vowel labeling task from Dupoux et al. (2011). The participants were told that they will be listening to foreign words over headphones and that they would have 5 seconds to choose a spelling choice that best matches the word they heard. The stimuli were presented through noise-isolating headphones, and answer choices that give the vowel-less and various vowel-ful spellings of the stimulus that just played were presented on screen simultaneously (e.g., [epuko] \rightarrow <epko>, <epako>, <epiko>, <epuko>). Participants selected their answer choices by using arrow keys on a keyboard (i.e., $\uparrow \downarrow \leftarrow \rightarrow$). A typical answer-choice screen is shown below in Figure 2.

```

      epako
    epiko      epuko
      epko
  
```

Figure 2: Answer choice screen for [epVko], where $V = /a, i, u, \emptyset/$.

While it is true that Japanese orthography is a syllabic system, most Japanese speakers are quite comfortable with the Latin alphabet, not only because of frequent exposure to loanwords but also because of the keyboards used for word processing. There are currently two main input methods—direct input (one key = one syllabic character) and conversion (QWERTY keyboard used to input CV combinations which are then converted to the corresponding syllabic character)—and the conversion method is commonly more preferred, and thus participants are expected to be comfortable with answer choices presented in the Latin alphabet. The experiment was designed to continue as soon as the participant makes an answer choice.

2.3. Analysis and predictions

All statistical analyses were performed by fitting linear mixed effects models using the *lme4* package (Bates et al., 2015) for R (R Core Team, 2016). The statistical analyses assess vowel detection and vowel identification. Detection refers to how often participants report perceiving any vowel at all both in the presence and absence of vocalic segments. Identification refers to whether the vowel the participants perceive is in agreement with the acoustic vocalic information contained in the stimuli.

In the case of detection, accuracy is expected to be higher for the *No-Devoice* group (e.g., [ez_po]) and lower in the *Lo-Predict* and *Hi-Predict* groups (e.g., [eʃ_po] and [es_po], respectively). Since the phonological process of high vowel devoicing is nearly obligatory (Vance, 1987), it could bias Japanese speakers toward mistakenly “recovering” a high vowel that is expected to be present between two voiceless consonants even when it is acoustically absent. Since the current experiment uses nonce-words, there is no underlying or lexical form to access. Devoiced and voiced sequences involving two voiceless obstruents in Japanese would map to the same phonotactically legal surface form (e.g., [esupo] \equiv [esupo] \equiv [espo] \rightarrow /esupo/). The devoiced and voiced sequences would all be regarded as legal, and the actual presence or absence of the vowel in the signal is readily ignored. It should be noted that if the stimuli being used were lexical items, reaction times would be predicted to differ depending on devoicing status of the stimuli. Ogasawara and Warner (2009) found that Japanese listeners are quicker to identify reducible lexical items when presented with devoiced overt forms. The interpretation of the results was that the most frequent overt form is considered to be the lexical form, and thus a direct mapping is possible from the most frequent overt form to the lexical form. For example, assuming that the lexical form of the loanword ‘star’ is ⟨sta:⟩, identification of the word would be the quickest when presented with the overt form [sta:] (\rightarrow ⟨sta:⟩). Conversely, rare but nevertheless equivalent forms would undergo repairs, leading to slower lexical access. For example, when presented with the rare, voiced overt form [suta:], it must be first mapped to the surface form /suta:/, then to the underlying form |sta:|, then to the lexical form ⟨sta:⟩.

While devoicing is possible in the *No-Devoice* environments, it is extremely rare (Maekawa & Kikuchi, 2005). Since only the vowel-ful token is legal in the language in non-devoicing contexts, the devoiced or vowel-less counterpart is not in an equivalence relationship (e.g., [sude] \neq *[sude] \neq *[sde] ‘barehand’). Thus Japanese speakers are expected to be more sensitive to the presence versus absence of a medial vowel. Furthermore, regardless of the stimulus group, higher accuracy is

expected in recognizing that there is no vowel as the burst/frication noise gets shorter, especially when there is no burst present.

In the case of identification, high accuracy is expected for the *Lo-Predict* group and lower accuracy for the *Hi-Predict* group. Japanese speakers have been shown to be sensitive to high vowel coarticulation in /ʃ/ (Beckman & Shoji, 1984), but this sensitivity is only useful when the vowel is unpredictable after a given C_1 (i.e., *Lo-Predict* group). Japanese listeners, therefore, should be sensitive to coarticulatory cues of at least [i, u]-coarticulation in the *Lo-Predict* group but biased towards a single high vowel that most frequently follows C_1 *Hi-Predict* group regardless of coarticulation. Since there are four answer choices <i, u, a, \emptyset >, identification rates are expected to be at least 50% in the *Lo-Predict* group and approximately 25% in the *Hi-Predict* group. Furthermore, since /a/ rarely devoices in Japanese, <a> responses should be relatively low even for [a]-coarticulated tokens, defaulting instead to the most phonotactically probable vowel. The *No-Devoice* group is expected to show some effects of coarticulation, as was the case in Dupoux et al. (2011), but like the *Hi-Predict* group, /a/ should show little effect.

These predictions contrast with the account given by Dupoux and colleagues. According to Dupoux and colleagues, there are two mechanisms at play during illusory vowel epenthesis. First, perceptual repair is a one-step process where phonotactically illegal sequences are perceived as their repaired counterparts rather than being perceived accurately first then repaired to their phonotactically legal counterparts. What this means is that listeners do not have access to the source language’s underlying form, making heterorganic C_1C_2 sequences and their repaired C_1VC_2 sequences equivalent for Japanese listeners. If this is correct, the prediction in terms of detection is that the rate of vowel detection between C_1C_2 and C_1VC_2 sequences should be statistically the same since the two sequences are equivalent.

Second, although Dupoux and colleagues argue that perceptual repair is triggered by phonotactic violations, the repair strategy employed is not a purely phonotactic one but also a phonetic one, where L2 listeners make phonetically minimal repairs to the phonotactically illegal sequence. In Japanese, /u/ is epenthesized because it is the shortest vowel in the language, whereas the epenthesized vowel is /i/ in Brazilian Portuguese for the same reason (Dupoux et al., 2011). What this means is that the choice of the epenthesized segment is not because it is the most phonotactically probable vowel between any given consonant cluster but because it results in the smallest possible phonetic change. Also, high vowels can delete in Japanese (Shaw & Kawahara, 2018), making many C_1C_2 and C_1i/uC_2 sequences equivalent. If the choice of the epenthetic segment is indeed based on the magnitude of phonetic change rather than phonotactic probability, no observable effect of phonotactic predictability is expected, since the phonotactic knowledge merely flags repair sites but is not involved in the repair itself. Vowel identification rates, therefore, are predicted to suffer across all contexts whenever the coarticulated vowel is not /u/. What this study aims to show, instead, is that minimizing phonetic changes are more relevant in low-predictability contexts and that phonotactic probability plays a more decisive role in high-predictability contexts. For example, given with [ek_±to], the surface form that involves the minimum phonetic change is /ek_±ito/, since

/i/-epenthesis retains velar-fronting. Mapping the same stimulus to /ekuto/ results in a larger degree of phonetic change despite /u/ being the “default” epenthetic vowel because /u/-epenthesis also removes velar-fronting. On the other hand, given the high-predictability token [eçto], the surface form it maps to is /eçito/, not because of phonetic change but because /çu/ is phonotactically less probable.

3. Results

Shown in Figure 3 below are the overall results of the experiment. Figure 3.A shows results for all C_1 and Figure 3.B for stop C_1 only, which consequently also results in the exclusion of all high-predictability tokens, since / ϕ , s, ç/ are all fricatives. The colors indicate the target vowels, and the solid and dashed lines indicate vowel detection and successful vowel identification rates, respectively. Vowel detection rates simply collapse all non-zero responses, whereas vowel identification rates only include cases where participant responses matched the coarticulated vowels in the stimuli (e.g., respond <epuko> for [epuko]). The smaller the distance between two lines of the same color, the higher the proportion of successful vowel identification.

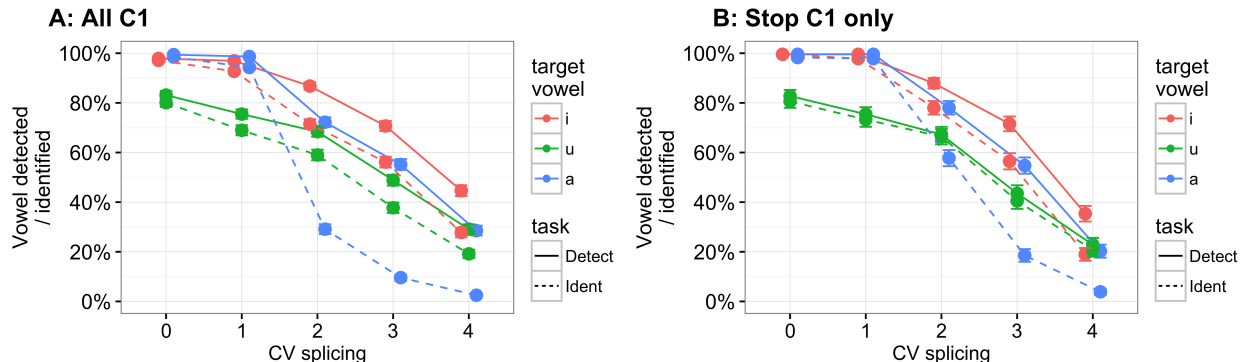


Figure 3: Vowel detection and identification rates with error bars by degree of splicing. CV splicing: 0 = full-CV, 1 = full-C half-V, 2 = full-C zero-V, 3 = half-C zero-V, 4 = zero-CV.

Figures 3.A and 3.B are qualitatively similar, where detection and identification rates fall as more of the C_1V_T information is spliced, and the most noticeable effect of including fricatives in 3.A is that identification rates are driven lower. In both figures, there are three things that stand out. First, detection rates for /u/ never reach 100% even when there is a full vowel of 40 ms present in the stimuli, suggesting that there is confusion between the presence and absence of /u/. Second, vowel detection rates never quite reach 0%, remaining above 20% even in the absence of any C_1 burst noise (Figure 3.B, splice-4), suggesting an overall confusion between vowel-fulness and vowel-lessness. Third, /a/ identification rates (blue dashed line) fall the most dramatically and are the lowest in tokens where the medial target vowel is spliced out, suggesting that only high vowels are potentially available for recovery.

Because the results of splice-1 and splice-3 tokens show no surprising trends, the rest of this paper will focus on the splice-0 (full-vowel), splice-2 (no vowel), and splice-4 (no vowel and no C_1 burst/frication) results. The splice-2 results will also be compared against naturally produced vowel-less tokens to test how the presence of coarticulatory cues affect the responses.

3.1. Tokens with full medial vowel

Shown below in Figure 4 are vowel identification rates for tokens with a full target vowel of 40 ms, broken down by context and by C_1 . As shown previously in Figure 3, the identification rates are surprisingly low for /u/.

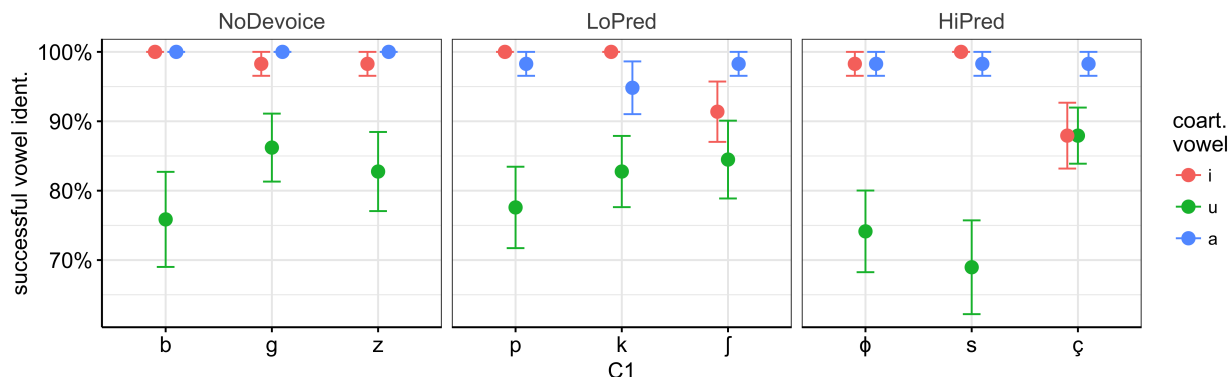


Figure 4: Successful vowel identification in VC_1VC_2V tokens with full medial vowel.

The most common wrong response by the participants for [u] identification was \emptyset for all C_1 as shown in Figure 5, meaning that the participants either heard the vowel accurately or confused [u] with \emptyset , but rarely confused the vowel with another vowel. The confusion specifically between $/C_1C_2/$ and $/C_1uC_2/$ sequences suggests two things. First, the overt sequence $[C_1uC_2]$ can be mapped to both $/C_1uC_2/$ and $/C_1C_2/$, although there is a bias towards the former. Second, fact that there is confusion between $/C_1uC_2/$ and $/C_1C_2/$ even when a vowel of 40 ms is fully present suggests that the distinction between the two sequences is weak, and that C_1C_2 and C_1uC_2 sequences are treated as more or less equivalent by Japanese listeners.

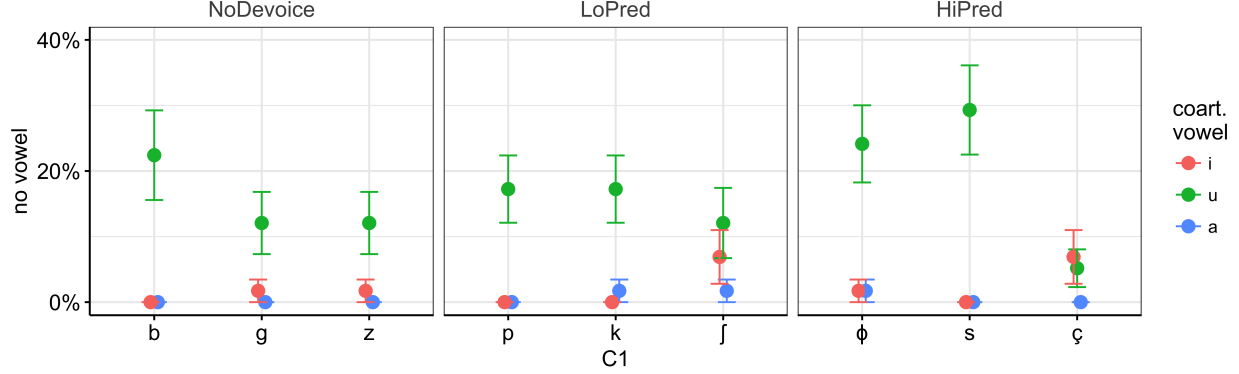


Figure 5: “No vowel” responses for VC_1VC_2V tokens with full medial vowel.

While this provides some support for the account presented by Dupoux and colleagues, the participants also exhibit some confusion between /i/ and \emptyset after /ʃ, ɕ/. The reason for this additional confusion most likely stems from the phonotactics of Japanese. Presented below in Table 2 are the observed/expected ratios for all pertinent C_1V biphones, calculated from the Corpus of Spontaneous Japanese. What the O/E ratios show is that /u/ is highly overrepresented in Japanese after most consonants. The exceptions are /ʃ, ɕ/ after which /i/ is the most common vowel, and /g/ after which /a/ is the most common vowel.

Table 2: Observed/expected (O/E) ratio of C_1V from CSJ. Highest O/E in bold.

	NoDevoice			LoPred			HiPred		
	b_	g_	z_	p_	k_	ʃ_	ɸ_	s_	ɕ_
_a	1.63	3.44	0.93	1.78	1.80	0.27	0.11	0.92	0.43
_i	0.79	0.31	0.00	0.65	1.12	6.28	0.10	0.04	6.28
_u	4.14	0.78	4.67	2.86	2.24	0.33	9.01	5.42	0.002
_e	1.24	0.75	2.30	0.49	0.97	0.003	0.12	0.90	0.006
_o	0.75	1.33	0.99	0.43	1.33	0.42	0.07	1.16	0.01

The results of the full-vowel tokens suggest that stimuli such as [epko, epuko, epuko] are possibly all being treated as equivalent to /epuko/. Because they all map to the same phonotactically legal structure, there is bidirectional repair, although with a bias towards vowel recovery. The fact that there is confusion for /u/ across the board, even for /g/ despite /a/ being the most common vowel to follow, provides some support to the phonetically minimal repair hypothesis presented by Dupoux and colleagues. However, the fact that there is also confusion for /i/ after /ʃ, ɕ/ additionally suggests that phonotactic probability affects perception as well.

3.2. Tokens with no medial vowel

This section compares the results of naturally vowel-less tokens and the splice-2 tokens where the medial, phonated vocalic material has been completely removed but C_1 burst/frication noise fully

remains. Acoustically, the difference between these tokens is that the naturally vowel-less tokens contain no obvious coarticulatory information, unlike the spliced tokens.

3.2.1. Naturally vowel-less tokens

The prediction in terms of vowel detection was that the rate of $\langle\emptyset\rangle$ responses should be highest for non-devoicing contexts since high vowel devoicing is rare in these contexts making Japanese listeners more sensitive to the presence versus absence of a medial vowel. Conversely, the rate of $\langle\emptyset\rangle$ responses was expected to be low in contexts where high vowel devoicing is expected, since high vowels can delete in these contexts, leading to a bias towards recovery. This bias should be especially high in high-predictability contexts because C_1 in these contexts are allophones that precede specific high vowels (e.g., $\zeta \rightarrow hi$).

Presented first below in Table 3 are the responses for naturally produced VCCV tokens. Bold numbers indicate the most frequent responses for a given C_1 . A chi-square test was performed using the *chisq.test()* function in R, to test whether the observed response rates were significantly different from chance. /a/ responses were excluded under the assumption that /a/ is not a candidate for recovery and also because /a/ responses were at or near 0% in most contexts. The results showed that the observed responses were significantly different from chance at $p < 0.01$ with the exception of /p/ ($p = 0.4909$).

Table 3: Responses for naturally produced VC_1C_2V tokens. Most frequent responses in bold.

	NoDevoice			LoPred			HiPred		
	ebko	egto	ezpo	epko	ekto	efpo	eϕko	espo	eçto
a	0.14	0.02	0.03	0.10	0.02	0.00	0.00	0.00	0.00
i	0.10	0.05	0.09	0.24	0.02	0.55	0.07	0.07	0.76
u	0.34	0.43	0.50	0.29	0.59	0.26	0.60	0.60	0.14
\emptyset	0.41	0.50	0.38	0.36	0.38	0.19	0.33	0.33	0.10

Overall, the results show that $\langle\emptyset\rangle$ responses are 50% or lower across all contexts, revealing an overall bias towards illusory epenthesis. However, the rate of $\langle\emptyset\rangle$ responses are highest for NoDevoice environments, suggesting that there indeed is an effect of high vowel devoicing. Additionally, $\langle\emptyset\rangle$ responses are lowest for HiPred environments, suggesting that predictability has an effect on the rate of repair as well.

The responses for the naturally vowel-less tokens also suggest that there is an effect of phonotactics that drives the choice of vowel that is recovered by Japanese listeners. The vowel recovered after [ʃ, ç] is, again, /i/ rather than /u/, further strengthening the account that the choice of the vowel used for phonotactic repair is not just merely a default, minimal vowel but rather chosen based on phonotactics. This is also in line with a recent finding by Durvasula and Kahng (2015), who also found in Korean listeners that the choice of recovered vowel is better predicted by the phonological alternations observed in the language rather than a phonetically minimal repair strategy.

3.2.2. Spliced vowel-less tokens (*Splice-2*)

Another prediction was that participants should be more sensitive to high vowel coarticulatory cues in contexts where high vowel devoicing is expected, and especially so in low-predictability contexts. A mixed logit model was fit using the *glmer* function of the *lme4* package of R, with successful vowel identification rates as the dependent variable. The statistical analysis compares the rate of correct identification of spliced vowels from coarticulatory cues, so naturally produced VCCV tokens, which should contain no vowel coarticulatory cues, are not included in the analysis. The fixed effects structure of the model consisted of target vowel, context, and their interaction. The model with a fully-crossed, maximal random effects structure failed to converge, hence the final random effects structure included by-participant and by-stimulus random intercepts as well as by-participant random slopes for target vowel and C₁. The interaction term was shown to be a non-significant contributor to the fit of the model ($p = 0.466024$), and thus was excluded from the final model. The results are shown below in Table 4 with spliced [i] tokens in LoPred contexts as the baseline. The diacritic for devoicing (i.e., ̥) is used throughout to indicate the vowels that have been spliced out.

Table 4: Mixed logit model results comparing successful vowel identification rates across difference predictability contexts. Spliced [i] tokens as baseline.

	Estimate	Std. Error	z	$\Pr(> z)$	
(Intercept)	2.3179	0.5048	4.591	4.40e-06	***
[u]	-0.9420	0.5879	-1.602	0.10906	
[a]	-2.9936	0.5460	-5.483	4.19e-08	***
NoDevoice	-0.9332	0.5121	-1.822	0.06844	.
HiPred	-1.6675	0.5146	-3.240	0.00119	**

The results show that the rates of high vowel identification were statistically comparable, but the rate of /a/ identification was significantly lower. Identification rates were also significantly lower in HiPred contexts than in LoPred contexts as predicted. Identification rates were also lower in NoDevoice contexts, although not significantly so. These results are also shown graphically in Figure 6 below.

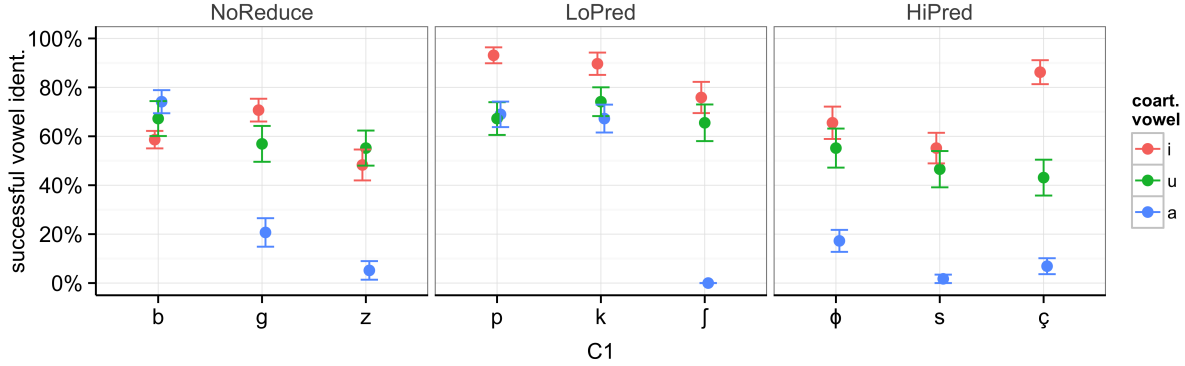


Figure 6: Successful identification rate of target vowel for spliced VCVCV tokens.

Since /a/ is a low vowel that is typically not targeted for devoicing in Japanese, another mixed logit model with the same full fixed and random effects structures was fit to the data with [a] tokens excluded in order to test how sensitive the participants were to high vowel cues specifically. Target vowel and its interaction with context were shown to be non-significant contributors to the fit of the model ($p = 0.07804$ and 0.36211 , respectively), and thus only context was retained as a predictor in the final model. Shown below in Table 5 are the results, with LoPred context as the baseline.

Table 5: Mixed logit model results comparing successful vowel identification rates across difference predictability contexts, excluding [a].

	Estimate	Std. Error	z	$\Pr(> z)$	
(Intercept)	1.7789	0.3597	4.945	7.6e-07	***
NoDevoice	-1.1143	0.4382	-2.543	0.01100	*
HiPred	-1.2712	0.4376	-2.905	0.00367	**

The results show that there is a clear effect of predictability on how successful Japanese listeners are in identifying coarticulated vowels when only high vowels are considered. Both NoDevoice and HiPred contexts have significantly lower identification rates, with HiPred being the lowest as predicted.

3.2.3. Comparison of naturally vowel-less and spliced vowel-less tokens

Naturally vowel-less tokens and spliced tokens by themselves tell only part of the story. Another prediction was that Japanese listeners should be able to recover high vowels from the coarticulatory information in spliced tokens, leading to differences between splice-2 and naturally vowel-less tokens. If it is the case that phonotactic violation alone is responsible for vowel epenthesis and that the choice of vowel is the phonetically minimal segment, namely /u/, then the presence of vowel coarticulatory information should do little to affect the choice of vowel.

Shown in Table 6 below are the results of a mixed logit model that compares detection rates for spliced tokens compared to a naturally vowel-less baseline. The results show that [i] coarticulation

drives up the vowel responses significantly, and nearly so for [a] coarticulation. The lack of an effect for [u] coarticulation again shows increased confusion in these contexts. Additionally, vowel responses are significantly higher for HiPred tokens compared to a NoDevice baseline.

Table 6: Mixed logit model results comparing vowel detection between VCCV and spliced VCVCV tokens.

	Estimate	Std. Error	z	$\Pr(> z)$	
(Intercept)	0.340915	0.361828	0.942	0.346089	
[i]	1.691540	0.480801	3.518	0.000435	***
[u]	0.634394	0.470474	1.348	0.177525	
[a]	0.789575	0.455128	1.735	0.082768	.
LoPred	0.673063	0.457778	1.470	0.141485	
HiPred	1.032983	0.454347	2.274	0.022993	*
[i]:LoPred	0.221310	0.700610	0.316	0.752092	
[u]:LoPred	-0.233197	0.646576	-0.361	0.718350	
[a]:LoPred	0.004041	0.648511	0.006	0.995028	
[i]:HiPred	-1.126838	0.670755	-1.680	0.092966	.
[u]:HiPred	-0.811062	0.650404	-1.247	0.212393	
[a]:HiPred	-1.552257	0.639344	-2.428	0.015187	*

For identification rates, spliced tokens are compared separately to naturally vowel-less tokens to make the effects of coarticulation for each vowel clearer. Presented below in Table 7 below are the responses for spliced [u] tokens. The rate of <u> responses is higher compared to naturally vowel-less tokens (Table 3 above).

Table 7: Responses for $VC_1(u)C_2V$ tokens with medial vowel spliced out. Most frequent response in bold.

	NoDevice			LoPred			HiPred		
	ebu <u>ko</u>	egu <u>to</u>	ezu <u>po</u>	epu <u>ko</u>	eku <u>to</u>	efu <u>po</u>	e <u>ph</u> u <u>ko</u>	esu <u>po</u>	e <u>ç</u> u <u>to</u>
a	0.02	0.00	0.02	0.00	0.00	0.00	0.02	0.00	0.02
i	0.02	0.00	0.12	0.00	0.00	0.09	0.03	0.05	0.47
u	0.67	0.57	0.55	0.67	0.74	0.66	0.55	0.47	0.43
∅	0.29	0.43	0.31	0.33	0.26	0.26	0.40	0.48	0.09

A mixed logit model was fit to the data with the rate of <u> responses as the dependent variable. <u> was chosen since it is regarded as the default epenthetic segment. The predictors were target vowel (i.e., /∅, i, u, a/), C_1 , and their interactions. C_1 was used as a predictor rather than context because the epenthetic vowel does not seem to be uniform across all contexts but rather depend on C_1 . By-participant and by-stimulus random intercepts were included. By-participant random slopes for target vowel and C_1 were also included. All predictors were significant contributors to the fit of the model. The results for <u> responses are shown below in Table 8, with ∅ tokens (i.e., naturally vowel-less) tokens as the baseline.

Table 8: Mixed logit model results comparing <u> responses between VCCV and spliced VC(u)CV tokens.

	Estimate	Std. Error	z	Pr(> z)	
(Intercept)	-0.84033	0.42766	-1.965	0.049418	*
[u]	1.89897	0.55110	3.446	0.000569	***
[g]	0.51123	0.51257	0.997	0.318574	
[z]	0.84678	0.53856	1.572	0.115880	
[p]	-0.39739	0.54914	-0.724	0.469270	
[k]	1.31600	0.53419	2.464	0.013757	*
[ʃ]	-0.45277	0.56776	-0.797	0.425184	
[ϕ]	1.58576	0.59324	2.673	0.007516	**
[ç]	-1.34398	0.68487	-1.962	0.049718	*
[s]	1.51309	0.57901	2.613	0.008969	**
[g]:[u]	-1.21609	0.73349	-1.658	0.097328	.
[z]:[u]	-1.62714	0.72768	-2.236	0.025347	*
[p]:[u]	0.50075	0.77601	0.645	0.518741	
[k]:[u]	-0.63586	0.78019	-0.815	0.415068	
[ʃ]:[u]	0.27845	0.76058	0.366	0.714290	
[ϕ]:[u]	-2.24928	0.78144	-2.878	0.003997	**
[ç]:[u]	-0.03906	0.78253	-0.050	0.960191	
[s]:[u]	-2.85498	0.77463	-3.686	0.000228	***

The model shows that there indeed is a significant raising effect of the coarticulated vowel on the rate of <u> responses. The mere presence of [k, ϕ, s] also drive up the rate of <u> responses significantly, while [ç] significantly lowers the rate of <u> responses, presumably because the expected vowel after [ç] is /i/. Interestingly, the expected vowel after [ʃ] is also /i/, but no lowering effect of <u> responses are observed. The interaction terms also show that the raising effect of [u] coarticulation is mitigated significantly after [z, ϕ, s] and nearly so after [g]. This is perhaps because the rates of <u> responses after these C₁ were already high in the naturally vowel-less tokens. These observations are also shown graphically in Figure 7 below.

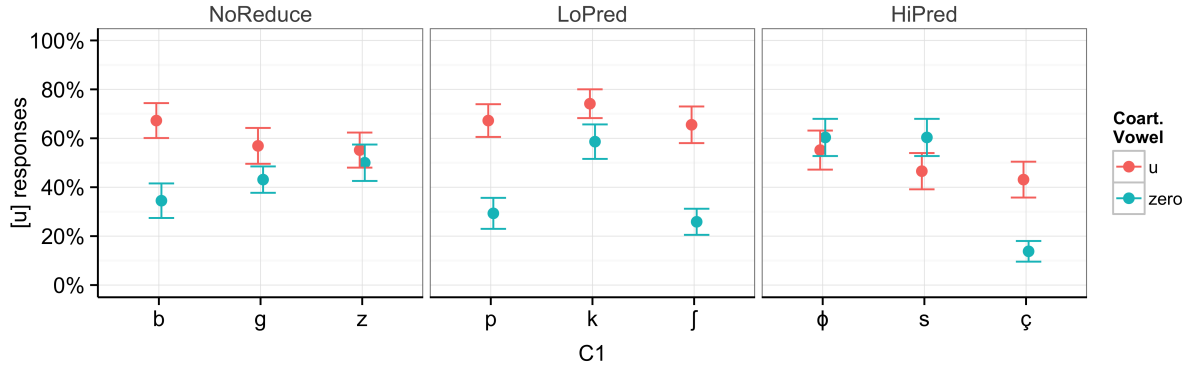


Figure 7: <u> responses for naturally vowel-less vs. spliced [u] tokens.

<i> responses are also driven up by [i] coarticulation. Shown below in Table 9 is a summary of the responses for spliced [i] tokens.

Table 9: Responses for VC₁(i)C₂V tokens with medial vowel spliced out.

	NoDevoice			LoPred			HiPred		
	ebiko	egito	ezipo	epiko	ekito	efipo	eφiko	esipo	eçito
a	0.09	0.03	0.00	0.00	0.00	0.00	0.00	0.00	0.00
i	0.59	0.71	0.48	0.93	0.90	0.76	0.66	0.55	0.86
u	0.10	0.10	0.41	0.03	0.03	0.10	0.19	0.24	0.03
∅	0.22	0.16	0.10	0.03	0.07	0.14	0.16	0.21	0.10

As was the case with [u] tokens, vowel coarticulation has a significant effect on which vowel participants report to hearing. A similar model as in Table 8 was fit, with the same predictors and random effects structure. The dependent variable was <i> responses with naturally vowel-less tokens as the baseline. The results are shown in Table 10 below.

Table 10: Mixed logit model results comparing <i> responses between VCCV and spliced VC(i)CV tokens.

	Estimate	Std. Error	z	Pr(> z)	
(Intercept)	-2.6791	1.0210	-2.624	0.00869	**
[i]	3.6475	1.4304	2.550	0.01077	*
[g]	-1.2952	1.5484	-0.836	0.40290	
[z]	-0.9058	1.4768	-0.613	0.53962	
[p]	0.3428	1.4628	0.234	0.81470	
[k]	-8.3719	4.7751	-1.753	0.07956	.
[ʃ]	2.9869	1.4244	2.097	0.03600	*
[φ]	-0.9942	1.5300	-0.650	0.51584	
[ç]	4.1493	1.4175	2.927	0.00342	**
[s]	-1.2364	1.5057	-0.821	0.41157	
[g]:[i]	1.9034	2.0904	0.911	0.36252	
[z]:[i]	-0.1457	2.0147	-0.072	0.94233	
[p]:[i]	2.6507	2.0799	1.274	0.20250	
[k]:[i]	16.1672	5.3645	3.014	0.00258	**
[ʃ]:[i]	-2.2153	1.9577	-1.132	0.25780	
[φ]:[i]	0.9448	2.0502	0.461	0.64491	
[ç]:[i]	-2.3320	1.9793	-1.178	0.23871	
[s]:[i]	0.6056	2.0328	0.298	0.76576	

In addition to the raising effect of [i] coarticulation on the rate of <i> responses, [ʃ, ç] also drives up the rate of <i> responses in naturally vowel-less tokens. This is also shown graphically in Figure 8 below. As discussed previously, the most common vowel after these consonants is /i/ in Japanese. The interaction term [k]:[i] also shows that the raising effect of [i] coarticulation is significantly higher after [k]. This is most likely due to the fact that /k/ is fronted before /i/, surfacing as [kʲ].

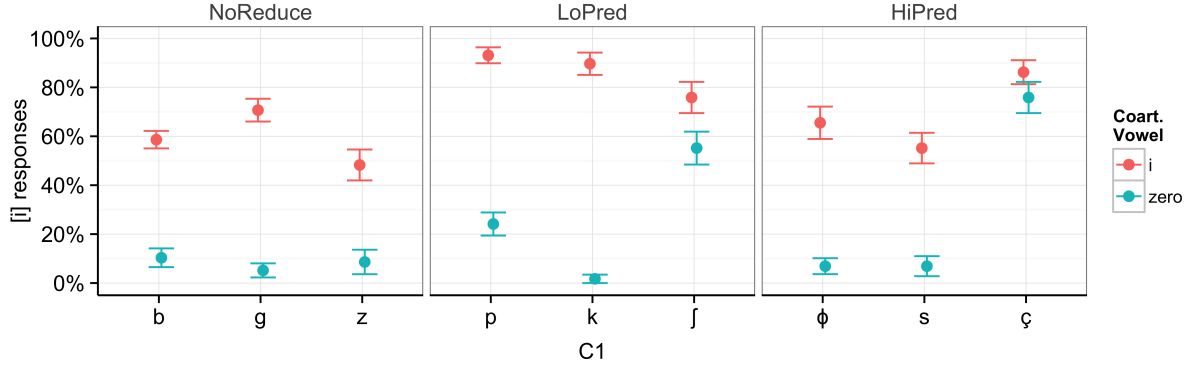


Figure 8: <i> responses for naturally vowel-less vs. spliced [i] tokens.

Thus far, the results suggest that the choice of epenthetic vowel for Japanese listeners is not simply a default /u/, but rather that the choice of vowel is sensitive to the acoustic cues in the signal. /u, i/ are both high vowels that are targeted for devoicing in Japanese, so this is perhaps not surprising. Japanese listeners have had a lifetime of practice attending to subtle coarticulatory cues to recover devoiced high vowels. Then what about a vowel like /a/, which rarely undergoes devoicing? The responses to spliced [a] tokens are shown in Table 11 below.

Table 11: Responses for VC₁(a)C₂V tokens with medial vowel spliced out.

	NoDevoice			LoPred			HiPred		
	eb̥ako	egato	ezapo	ep̥ako	ekato	efapo	eϕ̥ako	esapo	eçato
a	0.74	0.21	0.05	0.69	0.67	0.00	0.17	0.02	0.07
i	0.00	0.03	0.09	0.00	0.00	0.57	0.10	0.03	0.52
u	0.12	0.38	0.55	0.09	0.19	0.26	0.21	0.47	0.28
∅	0.14	0.38	0.31	0.22	0.14	0.17	0.52	0.48	0.14

Although limited to post-stop environments (i.e., [b, g, p, k]), the results show that participants can recover the spliced [a] vowel at relatively high rates. Bilabial place also seems to have a facilitatory effect. Given that <a> responses were generally low in naturally vowel-less tokens, the raising effect even in the limited environments is surprising. The fact that /a/ identification is limited to stops may be due to the articulatory differences between stops and fricatives. Because stops have a portion in which there is no airflow, coarticulation with the following vowel can be more complete by the time the stop burst/aspiration occurs. This is also true of bilabial place, where the lack of lingual gesture allows the following vowel to be coarticulated earlier. This is less true of fricatives where the transition into a fricative is more gradual, and coarticulation with the following vowel occurs towards the end of the segment rather than throughout. Since /a/ is a low vowel that a Japanese listener does not often have to recover, it may be that the beginning of the fricative already leads to the listener anticipating a high vowel and ignore to the low vowel cue towards the end.

A mixed logit model with the same predictors and random effects structure as in Tables 8 and 10 was fit. The interaction between target vowel and C_1 was not a significant contributor to the fit of the model and thus was excluded in the final model. The results are shown below in Table 12. The dependent variable was <a> responses with naturally vowel-less tokens as the baseline. Responses to [j] tokens were removed from the model because <a> responses are at 0% for both the naturally vowel-less and spliced [a] tokens, resulting in no meaningful difference. When included in the model, [j] tokens had an extremely low intercept of -27, but an absurdly high standard error of 22,246, both of which are most likely errors stemming from an absolute lack of difference between participants.

Table 12: Mixed logit model results for <a> responses.

	Estimate	Std. Error	z	$\Pr(> z)$	
(Intercept)	-3.0933	0.8330	-3.713	0.000205	***
[a]	4.6969	0.7104	6.612	3.80e-11	***
[g]	-5.5653	1.2935	-4.303	1.69e-05	***
[z]	-9.7734	2.2923	-4.264	2.01e-05	***
[p]	-0.2094	0.9610	-0.218	0.827515	
[k]	-1.2300	0.9920	-1.240	0.215034	
[ϕ]	-3.8215	1.1929	-3.204	0.001357	**
[ζ]	-8.7425	2.2688	-3.853	0.000116	***
[s]	-12.6671	2.6739	-4.737	2.17e-06	***

The results confirm that indeed [a] coarticulation does have a significant raising effect on the rate of <a> responses. With [b] as the baseline C_1 , the rate of <a> responses are significantly lower for [g, z, ϕ , ζ , s]. The results are also shown graphically in Figure 9 below. The figure shows that the rate of <a> responses is indeed much higher in the spliced tokens than in the naturally vowel-less tokens.

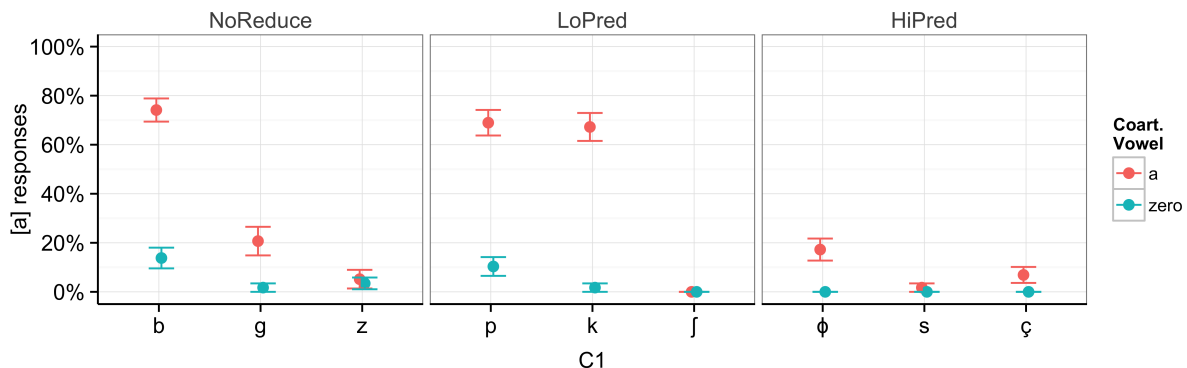


Figure 9: <a> responses for naturally vowel-less vs. spliced [a] tokens.

If the responses to naturally vowel-less tokens are taken as the default for phonotactically illegal

sequences, the responses to spliced tokens show that there is an additional effect of sensitivity to phonetic cues. Phonotactically illegal consonant clusters are indeed repaired, but the epenthetic vowel is chosen due to a combination of phonotactic predictability and sensitivity to phonetic cues. Japanese listeners seem more sensitive to coarticulatory information of high vowels across all environments but also non-high vowels like /a/ in contexts where coarticulatory information is easier to detect.

3.3. Tokens with no vowel and no burst/short frication noise

The results discussed in §3.2 for spliced vowel-less but burst-ful tokens (splice-2) show that Japanese listeners are biased towards perceiving a vowel between heterorganic consonant clusters, and that the choice of vowel is sensitive to the coarticulatory cues present in the C_1 burst/frication noise. Numerous studies have shown that the presence of a stop burst or frication noise in phonotactically illegal sequences are often interpreted as signaling the presence of a vowel (see Davidson & Shaw, 2012, Hsieh, 2013 for English; Furukawa, 2009, Whang, 2016 for Japanese; Kang, 2003 for Korean). This section therefore discusses the results of splice-4 tokens, where the target vowel has been spliced out completely and C_1 also has been spliced out leaving just the closure for stop C_1 and <15 ms of frication noise for fricative C_1 .

The responses to all splice-4 tokens are summarized in Table 13 below. A mixed logit model was fit to test whether the rates of <∅, i, u, a> responses were significantly affected by the identify of the vowel that was spliced out. Stop C_1 and fricative C_1 were analyzed separately. The results revealed that the responses were not significantly different regardless of the target vowel, with the exception of spliced [u] tokens where C_1 was /b/, which drove up <u> responses ($p = 0.002333$). Because the effect was limited to a single consonant, this section collapses the responses across all target vowels and focuses more on vowel detection.

Table 13: Responses for VC_1C_2V tokens with medial vowel and C_1 burst/frication noise spliced out.

	NoDevoice			LoPred			HiPred		
	eb`ko	eg`to	ez`po	ep`ko	ek`to	ej`po	eφ`ko	es`po	eç`to
a	0.05	0.02	0.02	0.01	0.01	0.01	0.01	0.00	0.01
i	0.08	0.13	0.16	0.07	0.02	0.36	0.09	0.10	0.45
u	0.32	0.17	0.34	0.07	0.09	0.11	0.11	0.14	0.10
∅	0.55	0.68	0.47	0.85	0.87	0.52	0.78	0.76	0.45

The results show first and foremost that the rate of <∅> responses never reaches 100%. This is perhaps expected for fricative C_1 , since there was ~15 ms of frication remaining in the tokens. Factors contributing to the results for stop C_1 , on the other hand, are less obvious. A mixed logit model was fit separately for the stops and fricatives since the the fricative tokens had a short frication noise remaining whereas the stop tokens had no burst at all. The full model for both data subsets had the following structures. The fixed effects included context, V_1 , and their interaction.

All stimuli used in the experiment had the form $V_1C_1(V)C_2V_2$, where the order of V_1 - V_2 was always either [e-o] or [o-e]. V_1 was included as a predictor to test whether the ordering of the initial and final vowels had a significant effect on vowel detection, which would suggest that there might be V-to-V coarticulatory cues that the participants are picking up on. The random effects included by-participant and by-stimulus random intercepts as well as by-participant random slopes for context, V_1 , and their interaction.

Shown first below in Table 14 is the result of the final model for the stop-only subset. Since the HiPred context had no stops, the subset only includes NoDevice and LoPred contexts with the latter as the baseline. The interaction term was shown to be a non-significant contributor to the fit of the model ($p = 0.5463$) and thus was removed.

Table 14: Mixed logit model result for vowel detection in spliced vowel-less and burst-less stop tokens.

	Estimate	Std. Error	z	$\Pr(> z)$	
(Intercept)	-1.8118	0.3162	-5.730	1.00e-08	***
$V_1 = [o]$	-0.4413	0.3369	-1.310	0.19	
NoDevice	1.5019	0.3493	4.299	1.71e-05	***

The results show that V_1 did not have a significant effect, but the rate of vowel detection was significantly higher for NoDevice tokens than LoPred tokens. A possible explanation for this effect is that the C_1 in NoDevice tokens had consistent phonation during closure, as shown in Figure 10 below.

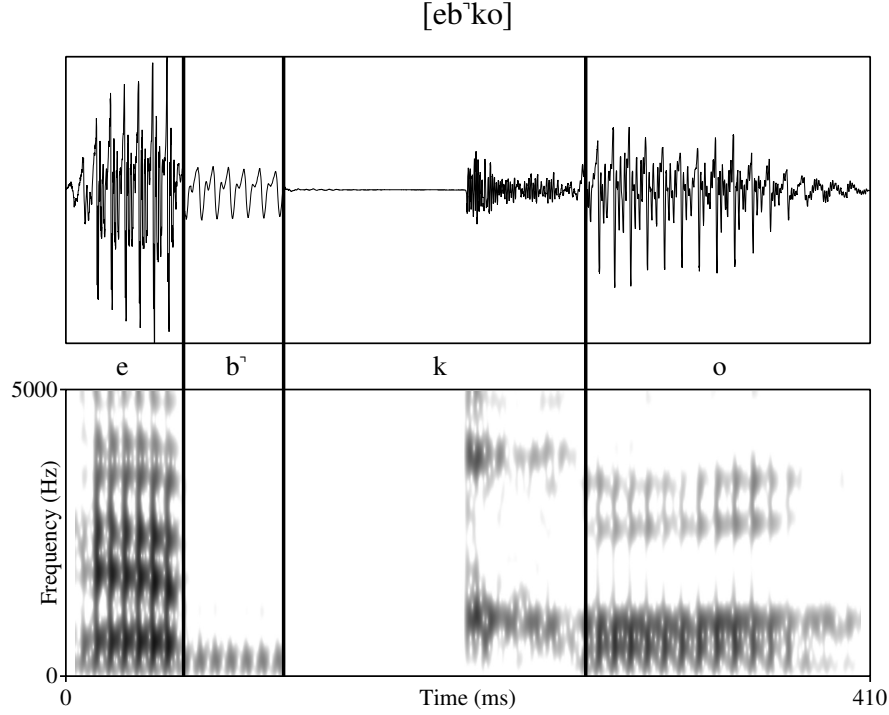


Figure 10: Spliced vowel-less, burst-less token created from [ebako].

The mixed logit model for the fricative-only subset also shows that the vowel detection rate for the NoDevice fricative [z] is significantly higher than for HiPred fricatives although not higher than the LoPred fricative [ʃ] ($p = 0.658$). For the fricatives, only context was a significant contributor to the fit of the model, and thus V_1 ($p = 0.81919$) and $V_1:\text{Context}$ ($p = 0.82666$) were excluded from the fixed effects structure of the final model. The results are shown below in Table 15.

Table 15: Mixed logit model result for vowel detection in splice-4 fricative tokens.

	Estimate	Std. Error	z	$\Pr(> z)$	
(Intercept)	-0.9063	0.3023	-2.998	0.00272	**
LoPred	0.7666	0.5488	1.397	0.16243	
NoDevice	1.0569	0.5065	2.087	0.03691	*

Although the fact that vowel detection rates never fall to 0% can be easily explained by the presence of prevoicing for NoDevice tokens and the 15 ms frication noise for the fricatives, the 10+% of vowel detection for the LoPred stops [p̥, k̥] is still somewhat puzzling. Without a vowel and without a burst between C_1 and C_2 , a token such as [ep̥ko] contains a doubly long stop closure, much like a geminate medial consonant as in [ekko]. Geminate consonants are phonotactically legal in Japanese and require no repair. Nevertheless, participants report perceiving a vowel some of the time. It is possible that some participants are picking up on the mismatch between the transitional cues out of V_1 and into V_2 . This seems unlikely, however, in that transitional cues into a vowel

often outweighs transitional cues out of a vowel for Japanese listeners (Fujimura et al., 1978) and that Japanese listeners rely more on centroid spectral cues than on formant transitions (Hirai et al., 2005). Perhaps a more likely explanation is one of task effect. Although the stimuli sounded as though they contain a geminate obstruent, there was no geminate option given as a possible answer. This might have kept the participants from fully eliminating the vowel-ful answer choices, and having been exposed to numerous vowel-ful tokens (both acoustically and perceptually) during the task, the participants might have assumed that a vowel should be present at least some of the time.

3.4. *Summary of main findings*

There were five main findings in the perceptual experiment. First, Japanese listeners seem to sometimes confuse the high vowel that is phonotactically the most likely after a given C_1 with \emptyset even when the high vowel is 40 ms long and fully phonated. This sort of confusion was not observed with the low vowel /a/, which typically does not devoice in Japanese. Second, results from naturally vowel-less tokens revealed that the vowel most often perceptually epenthesis between illicit clusters is /u/, largely due to the fact that it is phonotactically the most probable vowel after most obstruents in Japanese. This is further supported by the finding that after /ʃ, ʧ/, which is most often followed by /i/ rather than /u/, the choice of epenthetic vowel is in fact /i/. Third, participants successfully identified spliced high vowels in splice-2 tokens (full C_1 with target vowel completely spliced out) at rates significantly higher than the baseline rates observed in naturally vowel-less tokens. Identification rates of spliced /a/ were significantly lower and limited to after stops. Fourth, related to the third finding, identification rates of high vowels were lowest in HiPred contexts, suggesting that listeners are less sensitive to low-level coarticulatory cues in contexts where the phonotactics typically is sufficient for identifying the target vowel. Lastly, $\langle \emptyset \rangle$ responses never quite reach 100% even for splice-4 tokens where both C_1 and target vowel were fully spliced out.

4. Discussion

The aim of the current study was to test whether Japanese listeners are more sensitive to coarticulatory cues in low predictability contexts during perception than in high predictability contexts. Broadly speaking, overt consonant clusters were shown to be mapped to a phonotactically legal CVC sequence, neutralizing the contrast between CC and CVC sequences as Dupoux and colleagues have shown. However, the specific vowel recovered is modulated by CV co-occurrence probabilities in Japanese, as well as by detailed phonetic information.

First, the perception of full-vowel tokens showed that there is confusion between /u/ and \emptyset , even when there is a 40 ms-long, phonated [u]. It is possible that this confusion arises because /u/ is indeed the default epenthetic vowel in Japanese, making it equivalent to \emptyset . However, a survey of biphone co-occurrence probabilities in the Corpus of Spontaneous Japanese revealed that /u/ also happens to be the most common vowel after most consonants, making it difficult to attribute the

seemingly default status of /u/ as stemming simply from its shortness (Dupoux et al., 1999, 2011). Furthermore, similar confusion with \emptyset is observed for /i/ after /ʃ, ʒ/, suggesting that the choice of epenthetic vowel must be conditioned by the phonotactic probabilities of the language.

Second, the perception of vowel-less tokens further suggests that Japanese listeners confuse vowel-ful and vowel-less tokens with a tendency towards vowel-fulness. The results for splice-4 (vowel-less and burst-less) tokens in particular showed that Japanese listeners interpret even the most minute acoustic cues such as prevoicing of stops as signaling the presence of a vowel (§3.3). However, participants do not seem to simply perceive a default vowel. A comparison between naturally vowel-less and spliced vowel-less tokens showed that spliced tokens drive up the rate of target vowel responses significantly. This suggests that while heterorganic C_1C_2 sequences are perceived as being equivalent to C_1VC_2 as Dupoux and colleagues argue, the particular vowel is again not simply the “default” but the result of sensitivity to the acoustic information in the signal as dictated by the listener’s native language. The participants, therefore, are recovering the vowel that is the most probable based on the phonetic cues contained in the burst/frication noise of C_1 .

Third, the rate of high vowel identification was above chance at 40% across all contexts in spliced vowel-less tokens. Specifically, recovery rates were the highest in LoPred contexts as predicted, and the recovery rates were significantly lower for HiPred contexts, also as predicted. Recovery rates in NoDevoice contexts fell somewhere between the two devoicing contexts. The high rates of recovery suggest that Japanese listeners are hypersensitive to vowel coarticulatory cues, and the lower rate of recovery in HiPred contexts additionally suggests that sensitivity to coarticulatory cues are conditioned by phonotactic predictability.

Lastly, sensitivity to coarticulatory cues in Japanese listeners is limited primarily to high vowels. The participants were worst at identifying /a/. Non-high vowels are typically not devoiced in Japanese, and thus Japanese listeners have relatively little experience recovering them.

5. Conclusion

Based on the results discussed above, perhaps the terms perceptual epenthesis and “illusory” vowel epenthesis should be not used interchangeably. The “default” vowel is not [u] in Japanese simply because it is the shortest, but because it is the most common high vowel that Japanese listeners have been trained to recover all their lives. Phonotactic repair in Japanese listeners, therefore, is more akin to perceptual repair, where they use phonotactic and phonetic processes to choose the most probable vowel. In contrast, Brazilian Portuguese lacks a similar systematic devoicing process, and thus phonotactic repairs by Brazilian Portuguese listeners as reported by Dupoux et al. (2011) might be more “illusory” in nature, triggered primarily by phonotactic violations.

Also, the results suggest that language-specific sensitivity to phonetic cues are also context-specific. It has been observed for some time that listeners attend to the types of cues that are the most informative in their language. Korean listeners are more sensitive to V-to-C formant transitions than English speakers (Hume et al., 1999) because coda obstruents are obligatorily unreleased in Korean but optionally so in English (Kang, 2003). English listeners tend to interpret

the presence of nasality during a vowel as a coarticulatory cue for an upcoming nasal consonant, whereas Bengali speakers interpret the same cue as signaling a nasal vowel in a non-nasal context (Lahiri & Marslen-Wilson, 1991). Native Spanish listeners fail to show sensitivity to F1 difference between [i, ɪ] (Kondaurovaa & Francis, 2010) across all contexts. What the Japanese listeners in the current study also show is that listeners can lack sensitivity to phonetic cues that *are* contrastive in their native language, if low-level cues lead to the recovery of the same target as higher level processes. In other words, phonotactic knowledge can have a similar “top-down” effect as lexical knowledge in enhancing or dampening the perception to phonetic cues.

Acknowledgments

This material is based upon work supported by the National Science Foundation under Grant No. BCS-1524133.

References

- Bard, E., Anderson, A., Sotillo, C., Sotillo, Aylett, M., Doherty-Sneddon, G., & Newlands, A. (2000). Controlling the intelligibility of referring expressions in dialogue. *Journal of Memory and Language*, 42, 1-22.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48.
- Beckman, M. (1982). Segmental duration and the ‘mora’ in Japanese. *Phonetica*, 39, 113–135.
- Beckman, M., & Shoji, A. (1984). Spectral and perceptual evidence for CV coarticulation in devoiced /si/ and /syu/ in Japanese. *Phonetica*, 41, 61–71.
- Bladon, A. (1986). Phonetics for hearers. In G. McGregor (Ed.), *Language for hearers* (p. 1-24). Oxford, Pergamon Press.
- Browman, C. P., & Goldstein, L. (2000). Competing constraints on intergestural coordination and self-organization of phonological structures. *Les Cahiers de l’ICP. Bulletin de la Communication Parlée*, 5, 25–34.
- Bybee, J. (2006). From usage to grammar: The mind’s response to repetition. *Language*, 82(4), 711–733.
- Chitoran, I., Goldstein, L., & Byrd, D. (2002). Gestural overlap and recoverability: Articulatory evidence from Georgian. In N. Warner & C. Gussenhoven (Eds.), *Papers in Laboratory Phonology VII*. Berlin: Mouton de Gruyter.
- Davidson, L., & Shaw, J. (2012). Sources of illusion in consonant cluster perception. *Journal of Phonetics*, 40(2), 234-248.
- Dehaene-Lambertz, G., Dupoux, E., & Gout, A. (2000). Electrophysiological correlates of phonological processing: a cross-linguistic study. *Journal of Cognitive Neuroscience*, 12, 635-647.

- Dupoux, E., Kakehi, K., Hirose, Y., Pallier, C., & Mehler, J. (1999). Epenthetic vowels in Japanese: a perceptual illusion? *Journal of Experimental Psychology: Human Perception & Performance*, 25, 1568-1578.
- Dupoux, E., Parlato, E., Frota, S., Hirose, Y., & Peperkamp, S. (2011). Where do illusory vowels come from? *Journal of Memory and Language*, 64, 199-210.
- Durvasula, K., & Kahng, J. (2015). Illusory vowels in perceptual epenthesis: The role of phonological alternations. *Phonology*, 32(3), 385-416.
- Ernestus, M. (2011). Gradience and categoricity in phonological theory. In M. van Oostendorp, C. J. Ewen, E. Hume, & K. Rice (Eds.), *The Blackwell Companion to Phonology* (p. 2115-36). Wiley-Blackwell.
- Flege, J. E., Takagi, N., & Mann, V. (1996). Lexical familiarity and English-language experience affect Japanese adults' perception of / r/ and / l/. *Journal of the Acoustical Society of America*, 99, 1161-1173.
- Fujimoto, M. (2015). Vowel devoicing. In H. Kubozono (Ed.), *Handbook of Japanese Phonetics and Phonology* (chap. 4). Mouton de Gruyter.
- Fujimura, O., Macchi, M., & Streeter, L. (1978). Perception of stop consonants with conflicting transitional cues: A cross-linguistic study. *Language and Speech*, 21(4), 337-346.
- Furukawa, K. (2009). *Perceptual similarity in loanword adaptation between Japanese and Korean* (Unpublished master's thesis). University of Toronto.
- Hall, K. C., Hume, E., Jaeger, F., & Wedel, A. (in press). *The message shapes phonology*. (Forthcoming)
- Han, M. S. (1994). Acoustic manifestations of mora timing in Japanese. *Acoustical Society of America*, 96, 73-82.
- Hirai, S., Yasu, K., Arai, T., & Iitaka, K. (2005). Acoustic cues in fricative perception for Japanese native speakers. *Technical Report of Institute of Electronics, Information and Communication Engineers*, 104(696), 25-30.
- Hsieh, C.-H. (2013). *The perception of epenthetic vowels in voiced and voiceless contexts in Japanese* (Unpublished master's thesis). University of Kansas.
- Hume, E., Johnson, K., Seo, M., Tserdanelis, G., & Winters, S. (1999). A cross-linguistic study of stop place perception. In *Proceedings of the 14th International Congress of Phonetic Sciences* (p. 2069-2072).
- Jacquemot, C., Pallier, C., LeBihan, D., Dehaene, S., & Dupoux, E. (2003). Phonological grammar shapes the auditory cortex: A functional magnetic resonance imaging study. *The Journal of Neuroscience*, 23(29), 9541-9546.
- Kang, Y. (2003). Perceptual similarity in loanword adaptation: English postvocalic word-final stops in Korean. *Phonology*, 20(2).
- Kawahara, S. (2006). A faithfulness ranking projected from a perceptibility scale: The case of [+voice] in Japanese. *Language*, 83(2), 536-574.
- Kondaurovaa, M. V., & Francis, A. L. (2010). The role of selective attention in the acquisition

- of English tense and lax vowels by native Spanish listeners: comparison of three training methods. *Journal of Phonetics*, 38(4), 569-87.
- Kong, E. J., Beckman, M., & Edwards, J. (2012, November). Voice onset time is necessary but not always sufficient to describe acquisition of voiced stops: The cases of Greek and Japanese. *Journal of Phonetics*, 40(6), 725-744.
- Lahiri, A., & Marslen-Wilson, W. D. (1991). The mental representation of lexical form: A phonological approach to the recognition lexicon. *Cognition*, 38, 245-294.
- Maekawa, K., & Kikuchi, H. (2005). Corpus-based analysis of vowel devoicing in spontaneous Japanese: an interim report. In J. van de Weijer, K. Nanjo, & T. Nishihara (Eds.), *Voicing in Japanese*. Mouton de Gruyter.
- Marslen-Wilson, W. D. (1987). Functional parallelism in spoken word-recognition. *Cognition*, 25((1-2)), 71-102.
- Mattingly, I. G. (1981). Phonetic representation and speech synthesis by rule. In J. Myers, J. Laver, & A. J. (Eds.), *The cognitive representation of speech* (p. 415-420). North-Holland Publishing Company.
- Mattys, S. L., White, L., & Melhorn, J. F. (2005). Integration of multiple speech segmentation cues: a hierarchical framework. *Journal of Experimental Psychology: General*, 134(4), 477-500.
- McCarthy, J. J. (1999). Sympathy and phonological opacity. *Phonology*, 16, 331-399.
- McClelland, J. L., & Elman, J. L. (1986). The TRACE model of speech perception. *Cognitive Psychology*, 18, 1-86.
- Näätänen, R., Lehtokoski, A., Lennes, M., Cheour, M., Huottilainen, M., Iivonen, A., ... Alho, K. (1997). Language-specific phoneme representations revealed by electric and magnetic brain responses. *Nature*, 385, 432-434.
- Norris, D. (1994). Shortlist: A connectionist model of continuous speech recognition. *Cognition*, 52, 189-234.
- Ogasawara, N., & Warner, N. (2009). Processing missing vowels: Allophonic processing in Japanese. *Language and Cognitive Processes*, 24(3), 376-411.
- Pierrehumbert, J. (2001). Exemplar dynamics: word frequency, lenition, and contrast. In J. Bybee & P. Hopper (Eds.), *Frequency effects and the emergence of linguistic structure* (p. 137-157). Amsterdam: John Benjamins.
- Pitt, M., & McQueen, J. (1998). Is compensation for coarticulation mediated by the lexicon? *Journal of Memory and Language*, 39, 347-370.
- R Core Team. (2016). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria.
- Shademan, S. (2006). Is phonotactic knowledge grammatical knowledge? In D. Baumer, D. Montero, & M. Scanlon (Eds.), *Proceedings of the 25th west coast conference on formal linguistics* (p. 371-379).
- Shaw, J., & Kawahara, S. (2018). The lingual articulation of devoiced /u/ in Tokyo Japanese. *Journal of Phonetics*, 66, 100-119.

- Silverman, D. (1997). *Phasing and Recoverability* (Unpublished doctoral dissertation). UCLA, Los Angeles.
- So, C. K., & Best, C. T. (2010). Cross-language perception of non-native tonal contrasts: Effects of native phonological and phonetic influences. *Language and Speech*, 53(2), 273-293.
- Vance, T. (1987). *An Introduction to Japanese Phonology*. New York: SUNY Press.
- van Son, R., & Pols, L. (2003a). An acoustic model of communicative efficiency in consonants and vowels taking into account context distinctiveness. In *Proceedings of icphs 15* (p. 2141-2144). Barcelona.
- van Son, R., & Pols, L. (2003b). Information structure and efficiency in speech production. In *Proceedings of eurospeech 2003* (p. 769-772). Geneva.
- Varden, J. K. (2010, March). Acoustic correlates of devoiced Japanese vowels: velar context. *The Journal of English and American Literature and Linguistics*, 125, 35-49.
- Vitevitch, M. S., & Luce, P. A. (1998). When words compete: Levels of processing in spoken word perception. *Psychological Science*, 9, 325-329.
- Vitevitch, M. S., & Luce, P. A. (1999). Probabilistic phonotactics and neighborhood activation in spoken word recognition. *Journal of Memory and Language*, 40, 374-408.
- Vitevitch, M. S., Luce, P. A., Charles-Luce, J., & Kemmerer, D. (1997). Phonotactics and syllable stress: Implications for the processing of spoken nonsense words. *Language and Speech*, 40, 47-62.
- Whang, J. (2016). Perception of illegal contrasts: Japanese adaptations of Korean coda obstruents. In *Proceedings of Berkeley Linguistics Society* (Vol. 36).
- Wilson, C., Davidson, L., & Martin, S. (2014). Effects of acoustic-phonetic detail on cross-language speech production. *Journal of Memory and Language*, 77, 1-24.