

1     **Speech rate and the perception of consonant and**  
2     **vowel length in Japanese: neural entrainment and**  
3             **episodic memory**  
4

5 Timothy Gadanidis,<sup>1</sup> Yoonjung Kang,<sup>2,1,a</sup>

6 *<sup>1</sup> Department of Linguistics, University of Toronto, Toronto, ON, Canada*

7 *<sup>2</sup> Department of Language Studies, University of Toronto Scarborough, Toronto, ON, Canada*

8 *[timothy.gadanidis@mail.utoronto.ca](mailto:timothy.gadanidis@mail.utoronto.ca), [yoonjung.kang@utoronto.ca](mailto:yoonjung.kang@utoronto.ca)*

9

---

<sup>a</sup> Author to whom correspondence should be addressed.

10 Neural entrainment (persistence of neural oscillation that captures speech rhythm)  
11 and episodic memory-based perception (exemplar theory and belief-updating model)  
12 have both been proposed to account for speech rate-dependent perception. This  
13 paper probes these models by testing rate-dependent perception of Japanese vowel  
14 and stop length contrasts. Contextual speech rate was varied by manipulating all  
15 segments, vowels only, or consonants only in the carrier sentence. Robust rate  
16 effects were found across the board, supporting the neural entrainment account, and  
17 partial support was found for the episodic account: the rate effect was modulated by  
18 manipulation method for stops, but not for vowels.

## 19 1. Introduction

20           Speech is highly variable. Some sources of this variation are random, but  
21 many are structured and systematic, such as the effects of speakers' physiological or  
22 social differences and linguistic factors such as prosody, adjacent segments, and  
23 speech style. At issue in the present study is variation by speech rate: as the rate of  
24 speech increases or decreases, segments are shortened, or lengthened. For  
25 durationally-signaled contrasts, shortening and lengthening of segments requires rate-  
26 dependent speech perception. Listeners have been shown to take into account the  
27 rate of speech when identifying the length of a given stimulus — identical stimuli are  
28 more likely to be heard as 'long' when embedded in fast speech and 'short' in slow  
29 speech (Miller, Green & Reeves, 1986), and the rate effect is found even when the  
30 rate cue is distal, i.e., in the same utterance but one or more syllables away (Reinisch  
31 et al., 2011, Heffner et al., 2017).

32           Neural entrainment has been proposed as the cause of rate-dependent  
33 perception. Kösem et al. (2018) demonstrate that neural oscillation tracks the speech  
34 rhythm, and the oscillation persists a few cycles past stimulation (i.e., the speech rate  
35 cue), shaping the perception of temporal contrasts in subsequent speech. However,  
36 neural entrainment is unlikely to be the only mechanism at play in rate normalization:  
37 it cannot explain context-dependent perception involving non-temporal cues (e.g.,  
38 vowel quality, fricative place, and pitch) and cannot stand as a general account of  
39 context-dependent perception. Also, it has been shown that listeners can track  
40 speaker-specific speech rates and use them in subsequent perception distinctly

41 (Maslowski et al., 2019, Reinisch, 2016), which is not straightforward to explain  
42 under the neural entrainment account.

43         On the other hand, models have been proposed to account for listeners'  
44 ability to perceive the target structure accurately despite variability along diverse  
45 acoustic dimensions and linguistic and non-linguistic contexts. The exemplar-based  
46 model of speech perception (Pierrehumbert, 2001) holds that speakers store rich  
47 acoustic detail about past production of a structure so that speaker- and situation-  
48 specific information can be used in later perception. The belief-updating model of  
49 speech adaptation (Kleinschmidt and Jaeger, 2015) holds that when listeners  
50 recognize consistencies in the speech signal for a given situation, they will track  
51 situation- or speaker-specific distributions of acoustic cues. Listeners apply their  
52 knowledge of these context-specific cue distributions in perception when the context  
53 is recognized again. In these models, rate normalization arises when listeners infer  
54 whether a segment they heard is 'long' or 'short', not based on the absolute duration,  
55 but based on their knowledge of speech rate-specific distributions of the long and  
56 short segments.

57         In this study, we tested these accounts of rate-dependent perception by  
58 investigating whether the category-specific durational distributions in the ambient  
59 speech affect subsequent length perception in Japanese. Japanese is an ideal  
60 environment for testing multiple durational contrasts because it has a phonemic  
61 length distinction for both (voiceless) stops and vowels (e.g., kito = 'returning', kitto  
62 = 'surely', kiito = 'raw silk thread'). Both of these contrasts are potentially affected  
63 by speech rate variation: the durations of long vowels at fast speech rates overlap  
64 with the durations of short vowels at slow speech rates (Hirata, 2004), and the same

65 is true for stops (Hirata and Whiton, 2005). This allows us to probe how different  
66 segment types are affected by rate normalization when we vary the way fast and slow  
67 speech rates in the ambient speech are achieved, namely by either altering the  
68 duration of all segments equally or by manipulating only one category of segments. If  
69 neural entrainment is the sole or primary mechanism of rate normalization, we may  
70 see little or no difference based on which segments are manipulated. On the other  
71 hand, if rate normalization is due to context-specific inference of target structure, we  
72 expect different rate normalization patterns depending on the category-specific  
73 durational cues.

74         We report the results of five perception experiments on length perception of  
75 Japanese vowels and stops under varying speech rate conditions. We first confirm  
76 that rate normalization does occur for both categories (Experiment 1). We then find  
77 that how the rate is manipulated affects the perception of stops but not vowels: we  
78 find a weaker rate effect for stops when rate is varied by manipulating vowels only  
79 rather than all segments equally, while no difference was found for vowels  
80 (Experiment 2). This asymmetry holds even when the pitch cue for length is  
81 removed, which is a strong cue for vowel length but not for stop length (Experiment  
82 3) or when the manipulation conditions are presented in separate blocks  
83 (Experiment 4). The consonant-vowel asymmetry also holds in a consonant-only  
84 manipulation condition: the rate effect for stops is stronger when the rate is changed  
85 by manipulating consonants only rather than all segments equally, but no difference  
86 is found for vowels (Experiment 5).

## 87 **2. Methods**

88           We built five online experiments using *jsPsych* (de Leeuw, 2015) and recruited  
89 participants from prefectures throughout Japan via CrowdWorks, a Japanese  
90 crowdsourcing site. The gender breakdown (F:M) of the participants for the five  
91 experiments was 13:17, 16:24, 20:20, 20:19, and 8:31+11, respectively. The range and  
92 the median of the participants' age were comparable across the experiments: Exp. 1:  
93 27–62 (37.5), Exp. 2: 27–64 (42), Exp. 3: 20–62 (40), Exp. 4: 21–65 (42), and Exp. 5:  
94 24–61 (40). The source code for the experiments as well as the data, manipulation  
95 and analysis scripts, and stimuli are provided in the Data Availability Statement. We  
96 created duration continua based on four vowel-length minimal pairs (/kádo/ 'corner'  
97 ~ /ká:do/ 'card'; /bíru/ 'building' ~ /bí:ru/ 'beer'; /médo/ 'goal' ~ /mé:do/  
98 'brightness'; /tófo/ 'books' ~ /tó:fo/ 'beginning') and four consonant-length  
99 minimal pairs (/méta/ 'meta' ~ /mét:a/ 'thoughtless'; /táte/ 'height' ~ /tát:e/  
100 'stand up'; /kákó/ 'the past' ~ /kák:o/ 'parentheses'; /íka/ 'below' ~ /ík:a/  
101 'family'), some of which were drawn from Hirata (2004).

102           A male native speaker of Kanto Japanese produced both members of each  
103 target pair embedded in a carrier sentence, which was designed not to include any  
104 long segments: *Takeuchi-san wa totemo odayaka ni [word] to hatsuonshita* "Mr. Takeuchi  
105 very calmly pronounced [word]." One production of the carrier sentence was chosen  
106 and all target words were spliced into it at the same location to create the final  
107 versions of each baseline stimulus. Then, the duration continua were created from  
108 tokens of both a long and short member of each pair using Praat's PSOLA algorithm  
109 (Boersma and Weenink, 2021). The duration steps for each pair of continua were  
110 chosen to make estimated responses comparable across all word pairs based on the  
111 results of a 10-listener pilot test, also conducted via CrowdWorks: the midpoint of

112 each continuum was the point at which the odds of a participant in the pilot test  
113 hearing it as “short” or “long” were predicted to be even based on a regression  
114 model, with the ends of each continuum being the points at which the response was  
115 predicted to approach near-100% short and near-100% long. For each continuum,  
116 we manipulated the duration of the carrier sentence to be either 20% shorter (“fast”) or  
117 or 20% longer (“slow”) than the original production. Participants listened to the  
118 stimuli embedded in a carrier sentence and selected the normal orthographic form of  
119 the word they heard by key press. Except for Experiment 4, items were presented in  
120 a pseudorandom order.

121         Statistical analysis was conducted using R (R Core Team, 2021). Mixed-  
122 effects logistic regression models were fitted using the *glmer()* function from the *lme4*  
123 package (Bates et al., 2015b). Each model was built to predict the frequency of a  
124 ‘long’ (1) response as opposed to a ‘short’ (0) response. Figures were created using  
125 the *ggplot2* package (Wickham, 2016), using the *pilot* theme (Hawkins 2024). The  
126 predictors used in the experiments are summarized in Table 1. We include  
127 interactions between predictors motivated by our questions; significant interactions  
128 of interest were followed up using post-hoc tests using the *emmeans* package (Lenth  
129 2024). For model selection, we first aimed to fit the model with the maximal  
130 random-effects structure (Barr et al., 2013; Bates et al., 2015a), given the data, in all  
131 cases. We then used the stepwise regression functionality of the *buildmer* package  
132 (Voeten, 2020) to trim this maximal model down to the optimal model based on  
133 likelihood ratio tests.

134         Table 1. Summary of predictors across experiments.

<b>predictor</b>	<b>summary</b>	<b>contrasts</b>
step_n	The duration step of the segment	Numeric (centered)
original_length	The pre-manipulation segment length	short = -0.5, long = 0.5
rate	The carrier speech rate	fast = -0.5, slow = 0.5
segment_type	The target segment type	stop = -0.5, vowel = 0.5
manipulation_type	The segments manipulated to achieve carrier rate changes	Exp. 2–4: all = -0.5, vowels = 0.5 Exp. 5: all = -0.5, consonants = 0.5
subject_id	Unique participant label	random effect
word_pair	The target minimal pair	random effect

135

### 136 3. Results

137 In Experiment 1, we aimed to confirm that perceptual rate normalization  
138 occurs for both stops and vowels in Japanese. We created a 9-step durational  
139 continuum for each base word (8 pairs \* 2 words) and manipulated the duration of  
140 the carrier sentence to be 20% slower or 20% faster than the original production (2  
141 rates), for 288 trials in total. If there is rate normalization, we would expect ‘fast’



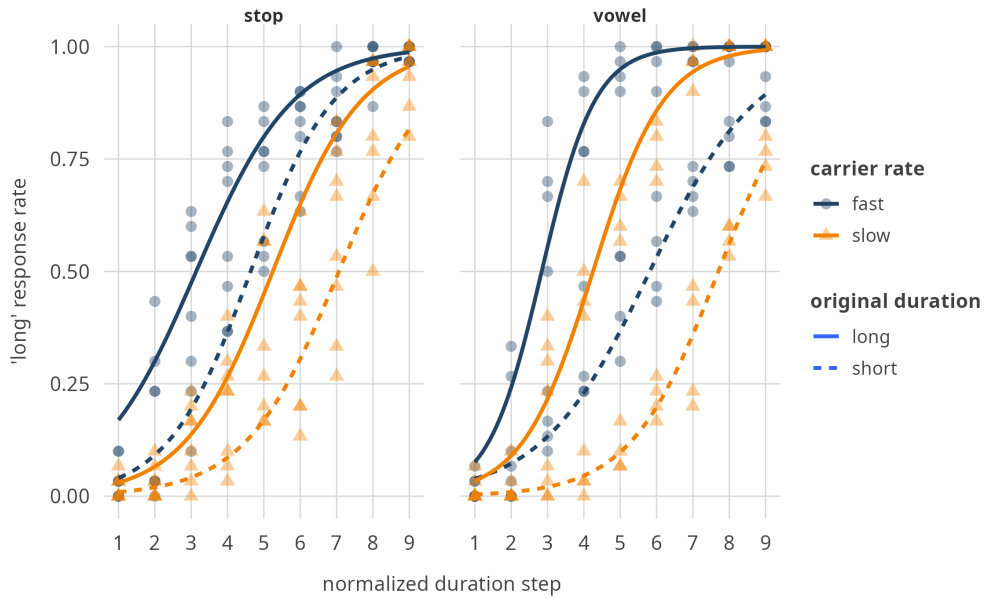
142 carriers to predict more ‘long’ responses, and we would expect ‘slow’ carriers to  
143 predict more ‘short’ responses.

144           The results of the linear mixed-effects regression model are displayed in  
145 Table 2, and the results are visualized in Figure 1. We find that rate plays a role, but  
146 that there is an interaction between rate and segment category, such that the effect of  
147 rate is weaker for vowels than for stops, even though vowels and consonants  
148 (including stops) are manipulated the same way in the carrier. Whether the stimuli  
149 were created from the original long or short word production also plays a large role,  
150 and this effect interacts with segment type, such that it is much stronger for vowels  
151 than for stops. This is likely due to pitch accent information: originally long vowel  
152 stimuli have a steeper pitch fall than originally short vowel stimuli (Kozasa, 2005).  
153 Post-hoc testing confirms that the effect of rate is significant across both vowels and  
154 stops and long and short baselines.

155           Table 2. Mixed-effects model output for Experiment 1. *Random effects kept by*  
156 *buildmer*: subject\_id and word\_pair, and by-subject and by-continuum random slopes  
157 for step\_n, rate, and original\_duration. *Nakagawa’s R2*: conditional = 0.780, marginal  
158 = 0.683.

	coef.	s.e.	p-value
(Intercept)	-0.06	0.22	0.8
step_n	1.05	0.04	<0.01
rate (slow-fast)	-2.03	0.08	<0.01
segment_type (vowel-stop)	-0.09	0.22	0.7
original_length (long-short)	2.54	0.21	<0.01

rate:seg_type	0.41	0.15	0.01
rate:original_length	0.43	0.14	<0.01
seg_type:original_length	1.97	0.4	<0.01
rate:seg_type:original_length	0.07	0.27	0.81



159 Fig. 1. The mean ‘long’ response rate for a given word pair at each duration  
160 rate<sup>1</sup>, segment type, and original length for Experiment 1. The lines are the LOESS  
161 curves.

162 In Experiment 2, we investigated whether listeners attend to category-  
163 specific duration variation across rate conditions or use the global speech rate for  
164 rate normalization. To test this, we added a new condition in which the rate of the  
165 carrier sentence was increased or decreased by 20% by manipulating vowels only.  
166 This 20% overall change was achieved by changing the vowel duration by 35.5%,  
167 while keeping the duration of consonants the same. Each participant completed both

<sup>1</sup>Plots based on raw duration are also available in our OSF repository (see Data Availability Statement).

168 rate manipulation conditions. This doubled the number of trials, so to make the  
 169 experiment more manageable, we reduced the number of steps in the durational  
 170 continuum from nine to five, for 360 trials in total. If listeners normalize by  
 171 attending to category-specific duration, we would expect the rate effect to differ  
 172 depending on the manipulation condition and the segment type: perception of  
 173 vowels would show a stronger effect when only vowels are altered than when all  
 174 segments are manipulated (same as Experiment 1), while stops would show a weaker  
 175 effect for the former than the latter.

176 The results of the linear mixed-effects regression model are displayed in  
 177 Table 3(a), and the results are visualized in Figure 2(a). We find that manipulation  
 178 type interacts with rate and segment type. A post-hoc test comparing the fast to slow  
 179 rate conditions shows that, as with Experiment 1, rate affects perception for both  
 180 vowels and stops, and this effect holds across both manipulation types, suggesting  
 181 that the rate effect can occur in the absence of category-specific duration change for  
 182 stops. As in Experiment 1, original length also has a strong main effect, which is  
 183 stronger for vowels than for consonants (not shown in Figure 2 (a) for readability).

184 Table 3. Mixed-effects model output for Experiments 2, 3 and 5. *Random*  
 185 *effects kept by buildmer*: intercepts for subject\_id and word\_pair, and by-subject and by-  
 186 continuum slopes for step\_n, rate (E5 only), and original\_duration (E2, E3 only).

	(a) Experiment 2			(b) Experiment 3			(c) Experiment 5		
Nakagawa's R2	Cond.		Marg.	Cond.		Marg.	Cond.		Marg.
	0.791		0.705	0.741		0.633	0.785		0.679
	coef.	s.e.	p-value	coef.	s.e.	p-value	coef.	s.e.	p-value

(Intercept)	-0.15	0.19	0.41	3.83	0.38	<0.01	3.8	0.32	<0.01
step_n	2	0.08	<0.01	1.8	0.09	<0.01	1.9	0.1	<0.01
manipulation_type (subset-all)	-0.02	0.06	0.69	-0.1	0.05	0.06	0.04	0.06	0.49
rate (slow-fast)	-2.33	0.18	<0.01	-2.11	0.06	<0.01	-2.62	0.23	<0.01
segment_type (vowel-stop)	-0.63	0.16	<0.01	-0.74	0.21	<0.01	-0.45	0.2	0.03
original_length (long-short)	2.44	0.07	<0.01	1.09	0.18	<0.01	2.18	0.07	<0.01
manip:rate	0.13	0.11	0.24	0.23	0.11	0.04	-0.58	0.11	<0.01
manip:seg_type	0.05	0.11	0.66	0.03	0.11	0.81	0.01	0.11	0.94
rate:seg_type	0.05	0.29	0.87	0	0.12	0.99	0.76	0.48	0.11
seg_type:orig_length	1.41	0.13	<0.01	-0.06	0.37	0.87	1.41	0.13	<0.01
manip:rate:seg_type	-0.55	0.23	0.02	-0.57	0.22	0.01	0.73	0.22	<0.01

187

188 Fig. 2. The mean ‘long’ response rate for a given word pair at each duration

189 step<sup>2</sup> by rate, segment type, and manipulation type for (a) Experiment 2 and (b)

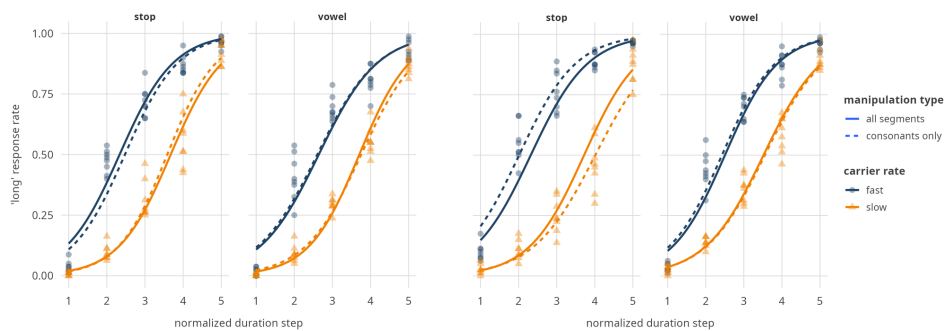
190 Experiment 5. The lines are the LOESS curves.

---

<sup>2</sup>Plots based on raw duration are also available in our OSF repository (see Data Availability Statement).

191 (a)

(b)



192

193 A second post-hoc test comparing the effect across segment types and  
194 manipulation types shows that rate has a weaker effect for stops in the vowel-only  
195 manipulation condition compared to the overall manipulation condition. In  
196 Figure 2(a), for stops, the two solid lines (overall manipulation) are further apart than  
197 the two dotted lines (vowel-only manipulation). However, the rate effect on vowels  
198 does not vary significantly by manipulation type: the two solid lines are no less  
199 separated from each other than the two dotted lines. If rate normalization is  
200 category-dependent, why do stops show sensitivity to manipulation type but not  
201 vowels?

202 One hypothesis is that because pitch accent is a strong cue to vowel length, it  
203 may play a larger role for vowels, making them less sensitive to speech-rate variation  
204 in general, so that the finer distinction between the two manipulation conditions  
205 does not affect perception. We test this hypothesis in Experiment 3, which replicates  
206 Experiment 2, but with all pitch information removed by flattening the pitch of all  
207 stimuli. Since the pitch cue is no longer available, we expect that we should find a  
208 difference between the vowel-only and overall rate manipulation conditions for  
209 vowels. The results of the linear mixed-effects regression model are displayed in

210 Table 3(b). We no longer find an interaction of original length with segment type.  
211 That is, vowels are no more affected by original length than stops, indicating that the  
212 loss of pitch cue for length in vowels makes vowels and stops more comparable in  
213 this regard. Otherwise, however, the results are essentially the same as in Experiment  
214 2. There is still a difference between manipulation types for stops, and no difference  
215 between manipulation types for vowels. In other words, pitch information seems to  
216 explain the stronger effect of original length for vowels but cannot explain the lack  
217 of a manipulation-type effect for vowels.

218 Another possibility is that, because the two manipulation conditions were  
219 intermingled and presented in pseudorandom order, participants did not cue in to  
220 the difference between the conditions. We test this hypothesis in Experiment 4 (not  
221 shown), which replicates Experiment 3 with a blocked design. Yet again, we find the  
222 results are largely the same as Experiments 2 and 3; whether the stimuli were  
223 presented in a blocked or unblocked order did not make a difference.

224 One final possibility is that the two manipulation types are not actually  
225 equally different with respect to vowels and consonants. For consonants, the  
226 difference is qualitative: there is no change in consonant duration in the vowel-only  
227 manipulation, versus there being some change in consonant duration in the global  
228 manipulation. For vowels, though, the difference is quantitative: there are changes in  
229 vowel duration in both manipulation conditions, with the difference just being  
230 between smaller and larger changes. It is perhaps not surprising, then, that  
231 consonants are more noticeably affected by the change in manipulation methods,  
232 while vowels are not. Experiment 5 followed up on this by inverting the way  
233 consonants and vowels are manipulated: rather than one condition where only

234 vowels are manipulated and another where both segment types are manipulated, we  
235 had one condition where consonants are manipulated and another where both  
236 segment types are manipulated. If listeners normalize using category-specific  
237 duration information, we would expect perception of vowels, not stops, to show  
238 different rate effects by manipulation type. The results of the linear mixed-effects  
239 regression model are displayed in Table 3(c), and the results are visualized in  
240 Figure 2(b). Post-hoc testing finds a significant rate effect for both vowels and stops  
241 in both manipulation conditions, which means that vowels are affected by rate even  
242 when only consonants are altered between the slow and fast rates, similar to the rate  
243 effect for stops in the vowel-only manipulation condition in Experiments 2–4. The  
244 result of main interest is an interaction between rate, manipulation type and segment  
245 type, which is significant. Post hoc testing shows that this is due to a difference in  
246 responses to stop segment trials: the effect of rate for stops is greater when only  
247 consonants are manipulated than when all segments are manipulated, while the rate  
248 effect for vowels remains the same. In other words, the different effects of  
249 manipulation for vowels and stops found in Experiments 2–4 persisted regardless of  
250 whether the manipulation was varied quantitatively (more or less duration change  
251 between speech rates) or qualitatively (change or no change).

#### 252 **4. Discussion**

253 This study has demonstrated that perception of both vowels and stops is  
254 affected by rate normalization in Japanese and found a robust rate effect across all  
255 experiments and all conditions, even when rates were varied by altering vowels or  
256 consonants only, leaving the other segment types unchanged. For vowels, we only

257 found the effect of the overall speech rate and no effect of category-specific duration  
258 variation. These findings are consistent with the neural entrainment account of rate  
259 effect (e.g., Kösem et al., 2018), whereby the neural oscillation tracks the overall  
260 rhythmic fluctuation of the speech signal and creates predictions for the upcoming  
261 speech. However, we also found a small but significant effect of manipulation type,  
262 but only for stops. The effect of rate normalization is weaker for stops when only  
263 vowel durations are manipulated (Experiments 2–4); and it is stronger when stop  
264 durations are manipulated to a higher extent (Experiment 5). This suggests that one  
265 of the cues listeners attend to when perceiving the length of a stop is the duration of  
266 other stops in that category in the speech stream, and this cue-specific duration  
267 variation can modulate the overall speech rate effect. This cue-specific effect is  
268 consistent with episodic accounts of rate normalization, which hold that listeners  
269 attend not (just) to overall rate, but to the duration of specific segments.

270         What remains puzzling is the asymmetry between vowels and stops. One  
271 possibility is the different ways that vowel and consonant gestures are organized in  
272 speech production. A body of articulatory phonetic studies has established that  
273 vocalic gestures are articulated contiguously through intervening consonantal  
274 gestures, while consonantal gestures are superimposed on top of a sequence of  
275 continuous vocalic gestures (Öhman, 1966; Carney and Moll, 1971; Fowler, 1983;  
276 Browman and Goldstein, 1990; Gafos, 1996). Under this conception of speech  
277 production, even when speech rate is modulated by stretching the duration of  
278 consonants only, the speaker-turned-listener reconstructs the vocalic gestures present  
279 underneath the superimposed consonantal gestures, and thus perceives the  
280 lengthening and shortening of vocalic gestures along with the consonantal ones.



281           Another possibility is that listeners start off with different expectations about  
282 the duration of the target vowels and consonants — and how they would vary by  
283 speech rate — based on their lifetimes of exposure, and the speech material in the  
284 carrier sentence can only modulate this expectation, not completely change it. If so,  
285 then what prior belief would listeners have to hold to exhibit the vowel-stop  
286 asymmetry observed in our study? Some have claimed that vowels generally vary  
287 more by speech rate than consonants in production in Japanese (Kuwabara, 1996)  
288 and other languages (Gay, 1978; Lo and Sóskuthy, 2023). If this is the case, the  
289 listeners may have more confidence about adjusting their vowel perception when  
290 they recognize fast or speech rate compared to stops and for stops, they may rely less  
291 on speech rate per se of the ambient speech, but the absolute duration of stops is  
292 relatively more important compared to vowels. The stronger rate effect for stops  
293 than vowels found in Experiment 1 may seem to contradict this prediction, but note  
294 that our rate manipulation, which varies the duration of vowels and consonants  
295 equally, may artificially exaggerate the consonant duration difference by rate, if they  
296 vary less than vowels in naturalistic fast and slow speech.

297           Indeed, rate normalization is not necessary or consistently attested in all  
298 durationally-signaled contrasts. Nakai and Scobbie (2016) found that the category  
299 boundary for the English voicing distinction does not meaningfully vary by speech  
300 rate and that there is no overlap between the VOT of voiced and voiceless, obviating  
301 the need for rate normalization. Theodore et al. (2009) found that speakers varied  
302 widely in the magnitude of VOT change by speech rate, with some exhibiting a large  
303 difference by speech rate while others showing almost no difference; this cross-  
304 speaker variation makes rate normalization less reliable. Relatedly, Ting and Kang

305 (2023) found that after listening to a short dialogue between two speakers, listeners  
306 did not appear to use each speaker's habitual speech rate as a cue for distinguishing  
307 /pi/ and /bi/, unlike German /a/ and /a:/ (Reinisch 2016). Future studies need to  
308 investigate the rate effects in the production of specific segment types in specific  
309 languages, and how (or if) their degree of variability in production informs their rate  
310 normalization in perception.

### 311 **Acknowledgments**

312 [to be added in the event of acceptance]

### 313 **Author Declarations**

#### 314 *Conflict of Interest*

315 The authors have no conflicts to disclose.

#### 316 *Ethics Approval*

317 The human-subjects component of the research was approved by the  
318 University of Toronto Office of Research Ethics. Informed consent was obtained  
319 from all participants.

### 320 **Data Availability**

321 The data that support the findings of this study, along with the experimental  
322 materials, are openly available in an Open Science Framework (osf.io) repository at  
323 <https://doi.org/10.17605/OSF.IO/SB8UF>, reference number SB8UF.

324

325 **References and Links**

- 326 Amano, S., and Hirata, Y. (2010). "Perception and production boundaries between  
327 single and geminate stops in Japanese," *The Journal of the Acoustical Society*  
328 *of America*, **128**, 2049–2058.
- 329 Barr, D. J., Levy, R., Scheepers, C., and Tily, H. J. (2013). "Random effects structure  
330 for confirmatory hypothesis testing: Keep it maximal," *Journal of memory*  
331 *and language*, **68**, 255–278.
- 332 Bates, D., Kliegl, R., Vasishth, S., and Baayen, H. (2015a). "Parsimonious mixed  
333 models," arXiv preprint arXiv:1506.04967.
- 334 Bates, D., Mächler, M., Bolker, B., and Walker, S. (2015b). "Fitting linear mixed-  
335 effects models using lme4," *Journal of Statistical Software*, **67**, 1–48.  
336 doi:[10.18637/jss.v067.i01](https://doi.org/10.18637/jss.v067.i01)
- 337 Browman, C. P., and Goldstein, L. (1990). "Tiers in articulatory phonology, with  
338 some implications for casual speech," *Papers in laboratory phonology I:*  
339 *Between the grammar and physics of speech*, **1**, 341–397.
- 340 Carney, P. J., and Moll, K. L. (1971). "A cinefluorographic investigation of fricative  
341 consonant-vowel coarticulation," *Phonetica*, **23**, 193–202.
- 342 de Leeuw, J. R. (2015). "jsPsych: A JavaScript library for creating behavioral  
343 experiments in a Web browser," *Behavior research methods*, **47**, 1–12.
- 344 Fowler, C. A. (1983). "Converging sources of evidence on spoken and perceived  
345 rhythms of speech: Cyclic production of vowels in monosyllabic stress feet."  
346 *Journal of Experimental Psychology: General*, **112**, 386.
- 347 Gafos, A. I. (1996). "V/C planar segregation and long-distance spreading revisited,"  
348 *Selected papers on theoretical and applied linguistics*, **9**, 25–37.

349 Gay, T. (1978). "Effect of speaking rate on vowel formant movements," The Journal  
350 of the Acoustical Society of America, **63**, 223–230.

351 Hawkins, O. (2024). *Pilot: A minimal ggplot2 theme with an accessible discrete color palette*. R  
352 package version 4.0.1. Retrieved from <https://github.com/olihawkins/pilot>

353 Heffner, C. C., Newman, R. S., and Idsardi, W. J. (2017). "Support for context  
354 effects on segmentation and segments depends on the context," Attention,  
355 Perception, & Psychophysics, **79**, 964–988.

356 Hirata, Y. (2004). "Effects of speaking rate on the vowel length distinction in  
357 Japanese," Journal of Phonetics, **32**, 565–589.

358 Hirata, Y., and Lambacher, S. G. (2004). "Role of word-external contexts in native  
359 speakers' identification of vowel length in Japanese," *Phonetica*, **61**, 177–200.

360 Hirata, Y., and Whiton, J. (2005). "Effects of speaking rate on the single/geminate  
361 stop distinction in Japanese," The Journal of the Acoustical Society of  
362 America, **118**, 1647–1660.

363 Idemaru, K., and Guion-Anderson, S. (2010). "Relational timing in the production  
364 and perception of Japanese singleton and geminate stops," *Phonetica*, **67**,  
365 25–46.

366 Kawahara, S., Kato, M., and Idemaru, K. (2022). "Speaking rate normalization across  
367 different talkers in the perception of Japanese stop and vowel length  
368 contrasts," *JASA Express Letters*, **2**.

369 Kleinschmidt, D. F., and Jaeger, T. F. (2015). "Robust speech perception: Recognize  
370 the familiar, generalize to the similar, and adapt to the novel," *Psychological*  
371 *Review*, **122**, 148.

372 Kösem, A., Bosker, H. R., Takashima, A., Meyer, A., Jensen, O., and Hagoort, P.  
373 (2018). “Neural entrainment determines the words we hear,” *Current*  
374 *Biology*, **28**, 2867–2875.

375 Kozasa, T. (2005). “An acoustic and perceptual investigation of long vowels in  
376 Japanese and Pohnpeian,” Ph.D. dissertation, University of Hawai’i at  
377 Manoa, Honolulu, Hawai’i.

378 Kuwabara, H. (1996). “Acoustic properties of phonemes in continuous speech for  
379 different speaking rate,” *Proceeding of fourth international conference on*  
380 *spoken language processing. ICSLP’96, IEEE*, 2435–2438.

381 Lahiri, A., Gewirth, L., and Blumstein, S. E. (1984). “A reconsideration of acoustic  
382 invariance for place of articulation in diffuse stop consonants: Evidence from  
383 a cross-language study,” *The Journal of the Acoustical Society of America*,  
384 **76**, 391–404.

385 Lenth, R. V. (2024). *Emmeans: Estimated marginal means, aka least-squares means*. R  
386 packaRetrieved from <https://CRAN.R-project.org/package=emmeans>

387 Lo, R. Y.-H., and Sóskuthy, M. (2023). “Articulation rate in consonants and vowels:  
388 Results and methodological challenges from a cross-linguistic corpus study,”  
389 In R. Skarnitzl and J. Volín (Eds.), *Proceedings of the 20th international*  
390 *congress of phonetic sciences, Guarant International*, 3206–3210.

391 Miller, J. L., Green, K. P., and Reeves, A. (1986). “Speaking rate and segments: A  
392 look at the relation between speech production and speech perception for  
393 the voicing contrast,” *Phonetica*, **43**, 106–115.

394 Nakai, S., and Scobbie, J. M. (2016). “The VOT category boundary in word-initial  
395 stops: Counter-evidence against rate normalization in English spontaneous

396 speech,” *Laboratory Phonology: Journal of the Association for Laboratory*  
397 *Phonology*, **7**, 1–31.

398 Öhman, S. E. (1966). “Coarticulation in VCV utterances: Spectrographic  
399 measurements,” *The Journal of the Acoustical Society of America*, **39**, 151–  
400 168.

401 Pierrehumbert, J. (2001). “Lenition and contrast,” *Frequency and the emergence of*  
402 *linguistic structure*, **45**, 137.

403 Port, R. F. (1981). “Linguistic timing factors in combination,” *The Journal of the*  
404 *Acoustical Society of America*, **69**, 262–274.

405 Port, R. F., Al-Ani, S., and Maeda, S. (1980). “Temporal compensation and universal  
406 phonetics,” *Phonetica*, **37**, 235–252.

407 R Core Team (2021). *R: A language and environment for statistical computing*, R Foundation  
408 for Statistical Computing, Vienna, Austria. Retrieved from [https://www.R-](https://www.R-project.org/)  
409 [project.org/](https://www.R-project.org/)

410 Reinisch, E., and Sjerps, M. J. (2013). “The uptake of spectral and temporal cues in  
411 vowel perception is rapidly influenced by context,” *Journal of Phonetics*, **41**,  
412 101–116.

413 Theodore, R. M., Miller, J. L., and DeSteno, D. (2009). “Individual talker differences  
414 in voice-onset-time: Contextual influences,” *The Journal of the Acoustical*  
415 *Society of America*, **125**, 3974–3982.

416 Ting, C., and Kang, Y. (2023). “The effect of habitual speech rate on speaker-specific  
417 processing in english stop voicing perception,” *Language and Speech*.

418 Voeten, C. C. (2020). “Buildmer: Stepwise elimination and term reordering for  
419 mixed-effects regression,” R package version 2.11. Retrieved from  
420 <https://CRAN.R-project.org/package=buildmer>  
421 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*, Springer-Verlag New  
422 York. Retrieved from <https://ggplot2.tidyverse.org>