

Perceptual normalization for speaking rate occurs below the level of the syllable

Margaret Cychosz¹ and Rochelle S. Newman¹

*Department of Hearing and Speech Sciences, University of Maryland, College Park,
College Park, MD*

1 Because speaking rates are highly variable, listeners must use cues like phoneme
2 or sentence duration to normalize speech across different contexts. Scaling speech
3 perception in this way allows listeners to distinguish between temporal contrasts,
4 like voiced and voiceless stops, even at different speech speeds. It has long been
5 assumed that this speaking rate normalization can occur over small units such as
6 phonemes. However, phonemes lack clear boundaries in running speech, so it is not
7 clear that listeners can rely on them for normalization. To evaluate this, we isolate
8 two potential processing levels for speaking rate normalization—syllabic and sub-
9 syllabic—by manipulating phoneme duration in order to cue speaking rate, while also
10 holding syllable duration constant. In doing so, we show that changing the duration of
11 phonemes both with unique spectro-temporal signatures (/ka/) and more overlapping
12 spectro-temporal signatures (/wi/) results in a speaking rate normalization effect.
13 These results suggest that when acoustic boundaries within syllables are less clear,
14 listeners can normalize for rate differences on the basis of sub-syllabic units.

¹⁵ I. INTRODUCTION

¹⁶ Speaking rate varies widely between and within speakers. Yet many phonological con-
¹⁷ trasts of the world’s languages rely on temporal cues, such as vowel length or voice onset
¹⁸ time (VOT), whose raw values vary by speech rate. Consequently, listeners must percep-
¹⁹ tually *normalize* for speaking rate, remapping acoustic cues across different contexts and
²⁰ speakers, in order to comprehend speech and acquire language.

²¹ Listeners employ perceptual normalization (or compensation) for speaking rate over a
²² variety of levels in the speech signal.¹ For example, to categorize a temporally-cued contrast
²³ like /k-g/, listeners could use *proximal* information in the speech signal, like the duration
²⁴ of vowels or consonants that are adjacent to the target phoneme (Diehl and Walsh, 1989;
²⁵ Miller and Liberman, 1979; Newman and Sawusch, 1996; Summerfield, 1981). Listeners
²⁶ could also use *distal* information in the speech signal like the rate of the overall sentential
²⁷ context, another talker’s habitual or situational speaking rate, or even the duration of non-
²⁸ speech stimuli like tones (Maslowski et al., 2019; Newman and Sawusch, 2009; Reinisch,
²⁹ 2016; Wade and Holt, 2005). In both cases, for a contrast like /k-g/, shorter-duration cues
³⁰ (e.g., a shorter adjacent consonant or a faster sentence) suggest a faster speaking rate and
³¹ therefore bias listeners to the positive VOT phoneme /k/. In contrast, longer-duration cues
³² bias listeners to the negative or neutral VOT phoneme /g/.

³³ Research on proximal information for speaking rate normalization has focused on cues
³⁴ such as the duration of phones preceding or following the target segment. As a result, we now
³⁵ know that although there are more degrees of freedom in vowel than consonant duration

³⁶ (Crystal and House, 1988), both vowels and consonants can provide rate normalization
³⁷ cues (Diehl and Walsh 1989; Summerfield 1981; see Toscano and McMurray 2012 for an
³⁸ alternative interpretation). There is evidence both for long-term, distal cues being employed
³⁹ during speaking rate normalization (Baese-Berk et al., 2014; Maslowski et al., 2019; ?; ?),
⁴⁰ as well as adjacency biases as listeners normalize over limited temporal windows of single
⁴¹ adjacent phonemes or syllables under typical listening conditions (Newman and Sawusch,
⁴² 1996; Sawusch and Newman, 2000).

⁴³ While careful experimental manipulations have led us to understand which cues listen-
⁴⁴ ers *can* use during perceptual normalization for speaking rate, less is known about the
⁴⁵ specific units that listeners employ. This gap in our understanding of rate normalization
⁴⁶ processes is relevant for a number of reasons, theoretical and applied. Research into prox-
⁴⁷ imal cues for rate normalization has traditionally assumed that phonemes are the basic
⁴⁸ unit over which speaking rate can be normalized. But this assumption may be prema-
⁴⁹ ture. For one thing, it is difficult for listeners to isolate phonemes in the comprehension
⁵⁰ of spontaneous, running speech. Articulatory undershoot and hypoarticulation compromise
⁵¹ phonological contrasts (Johnson et al., 1993; Lindblom, 1990). Coarticulation blurs acoustic
⁵² boundaries between adjacent phones as speakers consistently anticipate upcoming speech
⁵³ sounds (Whalen, 1990). The ability to resist coarticulatory pressures from adjacent phones
⁵⁴ decreases with increased lingual contact on the palate. The result is that some manners of
⁵⁵ articulation with more lingual contact, such as glides or laterals, are especially susceptible to
⁵⁶ coarticulation with adjacent phones (Recasens, 1985) and that some phonemes, particularly
⁵⁷ voiced, non-strident phonemes, are not reliably discriminable, meaning that their bound-

58 aries with adjacent sounds could be less clear. Syllables, however, are sometimes classified
59 as relatively more temporally-based (?) and temporal encoding is highly-discriminable even
60 in noise ([Giannella Samelli and Schochat, 2008](#)) and by young infants (6-12 months, with
61 sensitivity increasing through middle childhood) ([Trehub et al., 1995](#)). It is therefore plau-
62 sible that listeners would instead normalize speaking rates over units that variably carry
63 stress, like syllables, or other segments with more well-defined acoustic boundaries which
64 may not straightforwardly correspond to linguistic representations.

65 Rate normalization has often been considered a low-level, domain-general auditory pro-
66 cess ([Bosker, 2017](#); [Miller and Dexter, 1988](#)): it is involuntarily activated after milliseconds of
67 exposure to a speech- ([Reinisch, 2016](#)) or non-speech-like stimulus and has been documented
68 in non-human (avian) species ([Welch et al., 2009](#)). However, it is also increasingly apparent
69 that several higher-level constructs such as language experience ([Baese-Berk et al., 2016](#)),
70 listener familiarity with the speaker ([Kleinschmidt, 2016](#); [Reinisch, 2016](#)), and some aspects
71 of language-specific structure such as intonation ([Steffman, 2019](#)) also mediate rate normal-
72 ization. It is thus possible that rate normalization interacts with additional higher-level
73 units, such as the syllable, although this has not been empirically tested.

74 Thus, understanding how rate normalization unfolds has clear relevance to theories of
75 speech perception and learning. However, understanding the units used in rate normaliza-
76 tion is also relevant for more applied, artificial intelligence applications. From a machine
77 learning perspective, invariance in the speech signal is a central obstacle to achieving higher-
78 performing speech-to-text and automatic speech recognition applications. Understanding
79 appropriate mechanisms for normalization, including rate normalization, in human listeners

80 may facilitate machine performance, as it may be simpler to program normalization on the
81 basis of signals that *tend* to be relatively more acoustically-driven (such as syllables) than
82 those that tend to be more linguistically-driven (phonemes). If human listeners reliably
83 normalize for speaking rate at the phonemic level, even in the absence of explicit acoustic
84 boundaries, it would suggest that phonemic structure could be incorporated into natural lan-
85 guage processing algorithms to benefit machines' learning of speech (though the mechanisms
86 may vary by the type of speech, e.g., read versus spontaneous).

87 **A. Cue integration as an alternative**

88 Some work on proximal cues for phoneme classification has suggested that listeners may
89 not normalize for speaking rate via temporal cues such as syllable or VOT duration but
90 rather *integrate* acoustic cues that overlap with speaking rate to classify phonemes during
91 real-time speech processing (Massaro and Cohen, 1983; Toscano and McMurray, 2010, 2012,
92 2015). For example, duration-dependent phonetic decisions, such as stop voicing, could be
93 made sequentially by computing VOT and *then* vowel duration (vowel duration is likewise
94 a cue to stop voicing as the burstiness of voiceless stops can cause the the following vowel
95 to de-voice slightly at onset, leading to shorter perceived duration (Allen and Miller, 1999),
96 but see Turk et al. (2006)). Evidence from the visual world paradigm, as well as phoneme
97 decision tasks, suggests that listeners do indeed integrate multiple phonetic cues in this way,
98 sequentially, as they become available in the speech signal (McMurray et al., 2008; Miller
99 and Dexter, 1988; Toscano and McMurray, 2015). This result provides evidence against

¹⁰⁰ a speaking rate normalization account because such accounts would predict simultaneous
¹⁰¹ integration of VOT and vowel duration.

¹⁰² It was not the goal of this study to contrast cue integration and rate normalization ac-
¹⁰³ counts to explain proximal effects upon phonetic boundary shifts—and the results of [Toscano](#)
¹⁰⁴ and [McMurray \(2012\)](#), among others, do convincingly demonstrate that vowel length inte-
¹⁰⁵ gration, not normalization for speaking rate, explains proximal effects upon stop voicing
¹⁰⁶ classifications. Nevertheless, should we find an effect of consonant duration in the current
¹⁰⁷ studies, we believe that this could be interpreted as rate normalization and not the more
¹⁰⁸ straightforward acoustic cue integration. This is because our target contrast for both studies,
¹⁰⁹ /ʃ-tʃ/, will be cued by the duration of the following *consonant*, not vowel (/k/ in /ʃkas/-
¹¹⁰ /tʃkas/ for Exp. 1 and /w/ in /ʃwɪb/-/tʃwɪb/ for Exp. 2; these stimuli will be explained in
¹¹¹ more detail in the following section). But more importantly, there is no evidence that stop
¹¹² or glide duration reliably indicates fricative-affricate classification. And unlike the effect of
¹¹³ stop aspiration upon perceived vowel length (aspiration causes vowel de-voicing), there is
¹¹⁴ no phonetic reason to assume that fricatives and affricates would have different effects on
¹¹⁵ /k/ or /w/ duration or voicing. Consequently, should the current study find an effect of
¹¹⁶ consonant duration/speaking rate upon the phonetic boundary shift between /ʃ/ and /tʃ/,
¹¹⁷ it could indicate rate normalization, not cue integration.

¹¹⁸ **B. Current study**

¹¹⁹ The present experiments were designed to investigate the effects of acoustic separability,
¹²⁰ or the ability to distinguish between two adjacent phonemes, on speaking rate normaliza-

tion. Here and throughout the manuscript we will refer to “rate normalization” broadly, though we wish to emphasize that our results concern backward, proximal rate effects. The overarching goal is to understand the processing level (syllabic or sub-syllabic) involved in the perceptual normalization of speaking rate. In a pair of phoneme category rating experiments, we asked whether phones differing in acoustic separability (acoustically-distinct /ka/ versus overlapping /wɪ/) would result in separate rate normalization effects or in a single combined rate normalization effect. We chose to evaluate the effects of speaking rate upon the perception of the /ʃ-ʈʃ/ contrast in American English as this contrast has demonstrated a rate normalization effect in prior research (Newman and Sawusch, 1996; Repp et al., 1978) and its primary acoustic cue is temporal. For example, Repp et al. (1978) manipulated the duration of noise (frication) and silence between words in the phrase “gray ship” and found that shorter noise intervals predisposed listeners to hear word-initial /ʈʃ/ or “gray chip.” Similarly, in Newman and Sawusch (1996), the authors were able to trigger a /ʃ-ʈʃ/ phonetic boundary shift in a nonce word series ranging from /ʃkas/ - /ʈʃkas/ (“shkas” to “chkas”) by adjusting the duration of /k/ in the stimuli. Ambiguous stimuli, with a longer /k/ duration, suggested a slower speaking rate and biased listeners to perceive /ʈʃ/ while a shorter-duration /k/ suggested a faster speaking rate and biased listeners to perceive /ʃ/.

A limitation of previous work on this topic, including Newman and Sawusch (1996), is that changes to the duration of a single phoneme like /k/ also rendered changes to the duration of the surrounding syllable (e.g. /ka/) and word (e.g. /ʃkas/): a longer-duration /k/ resulted in a longer /ka/ syllable and /ʃkas/ word. As a result, any rate normalization

¹⁴² effect could just as easily be attributed to the duration of the manipulated phoneme as the
¹⁴³ duration of the entire syllable or word.

¹⁴⁴ To isolate sub-syllabic information as the potential processing unit in speaking rate nor-
¹⁴⁵ malization, Experiment 1 uses the same /ʃkas/-/tʃkas/ series as previous work but varies
¹⁴⁶ the syllable nucleus /a/ duration in the opposite direction of /k/. This manipulation leads
¹⁴⁷ to a /ʃkas/-/tʃkas/ series with consistent syllable and word, but varying phoneme, dura-
¹⁴⁸ tions. Although this adjustment to the original stimuli design is small, it has important
¹⁴⁹ consequences: all previous work that attempted to identify the units that listeners use in
¹⁵⁰ rate normalization had a confound in the stimuli design and as such, to date, we have not
¹⁵¹ been able to identify the units that listeners use. In the current design, any rate normal-
¹⁵² ization effect cannot be due to syllable or word durations, as the series did not differ in
¹⁵³ these respects.² Instead, the normalization effect could only be caused by variation in the
¹⁵⁴ manipulated phoneme /k/. Finding a rate normalization effect would suggest that the /k/
¹⁵⁵ was treated as a separate unit from the following vowel and that rate normalization took
¹⁵⁶ place over sub-syllabic, potentially phoneme-sized, units.

¹⁵⁷ Varying the nucleus duration in the opposite direction of the consonant is unlikely to
¹⁵⁸ cancel out any potential effect of the consonant's duration because duration effects are (1)
¹⁵⁹ weighted by distance (and /k/ is linearly closer to the target contrast in the /ʃkas/-/tʃkas/
¹⁶⁰ series) and (2) proportional (and /k/ is much shorter than /a/ so similar durational changes
¹⁶¹ (e.g. 20 ms) have disproportionate impacts upon /k/ and /a/). Indeed, we do find a
¹⁶² rate normalization effect in Experiment 1, suggesting that changing the vowel duration in
¹⁶³ the opposite direction did *not* cancel out any consonant duration effect. Finding a rate

¹⁶⁴ normalization effect in the Experiment 1 stimuli leads us to conduct Experiment 2 where we
¹⁶⁵ again test for rate normalization effects but using syllables that contain less discriminable
¹⁶⁶ phones. We use a similar nonce word series ranging from /ʃwɪb/ - /tʃwɪb/ (“shwihb” to
¹⁶⁷ “chwihb”) where we manipulated the duration of /w/ in /wɪ/. Although we did find an effect
¹⁶⁸ of /k/ duration upon perception of the initial /ʃ-tʃ/ contrast in Experiment 1, suggesting
¹⁶⁹ sub-syllabic level processing during rate normalization, we hypothesized that we may *not*
¹⁷⁰ find this same effect of /w/ duration on the same /ʃ-tʃ/ contrast in Experiment 2, suggesting
¹⁷¹ higher-level (syllable or word) processing for sequences with less-discriminable phones.

¹⁷² II. EXPERIMENT 1

¹⁷³ A. Methods

¹⁷⁴ *a. Participants.* Twenty-one members of the University of Iowa community partici-
¹⁷⁵ pated in this experiment for course credit. All listeners were native speakers of American
¹⁷⁶ English, and had no reported history of a speech or hearing impairment. An additional 8
¹⁷⁷ listeners completed the experiment, but their data were removed because in later question-
¹⁷⁸ ing they were found not to be a native English speaker ($n=1$) or they failed to respond on
¹⁷⁹ at least 80% of trials ($n=7$) due to an automatic 3000 ms trial timeout. We cannot be sure
¹⁸⁰ why these participants failed to respond, but it could be due to task fatigue or boredom.
¹⁸¹ The remaining participants responded to, on average 91.3% of trials ($SD=5.92$).

¹⁸² *b. Stimuli.* An adult native English-speaking man was recorded producing the syllable
¹⁸³ /fkas/ in carrier phrases (“He said the word X.”). His speech was digitized via a 12-bit,

184 analog-to-digital converter at a 10-kHz sampling rate, low-pass filtered at a 4.8 kHz, and
185 amplified. The initial consonant /ʃ/ was then separated from the remainder of the syllable,
186 with the boundary being the onset of closure for the following /k/. A continuum of ten
187 items, /ʃ/-/tʃ/, was then created by removing successive 10-ms sections from the /ʃ/ onset.
188 A linear amplitude ramp, with duration varying along with frication duration, was used over
189 the initial portion of each token to give the items a more natural attack. The duration of
190 the ramp varied from 6 to 60 ms, with a 9 ms step. The resulting series ranged from 60-145
191 ms in duration, with the longer frication sounding more similar to a /ʃ/ and the shorter
192 frication sounding more similar to a /tʃ/. Further details on the original stimulus creation
193 can be found in [Newman and Sawusch \(1996\)](#).

194 The remainder of the word—the syllable /kas/—was edited to create two new syllables,
195 one with a shorter /k/ (and longer /a/) and one with a longer /k/ (and shorter /a/). We
196 interpreted the /k/ to include the closure, burst, aspiration, and first four pitch pulses (which
197 appeared to correspond to the transition of the first formant). The duration of this base
198 /k/ was between 1/3 and 1/2 that of the vowel (see Figure 1). Thus, an equivalent amount
199 of change in duration for /k/ and /a/ will be much larger proportionately for /k/.

200 The duration of /k/ was altered by removing or reduplicating pitch pulses and sections of
201 burst and aspiration. Only short, nonadjacent sections of burst and aspiration were deleted
202 or reduplicated so as to maintain the general amplitude profile and prevent the perception
203 of frozen noise. No change was made to the closure duration; although closures do tend
204 to vary slightly with speaking rate, this variability is typically quite small ([Crystal and](#)
205 [House, 1988; Gay, 1978](#)), and thus unlikely to have a substantial perceptual effect. For

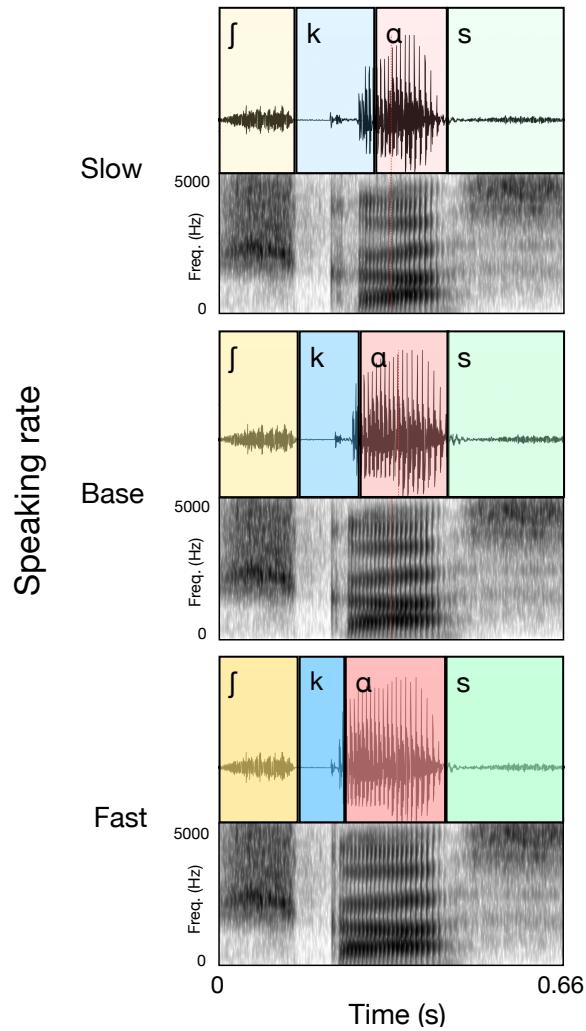


FIG. 1.

Speaking rate manipulations on the basis of /k/ duration and stimuli duration for first step of

series: Experiment 1.

the short /k/ stimulus, two pitch pulses were removed, as well as 17.2 ms of the burst
 and aspiration; for the long /k/ stimulus, four pitch pulses and 22 ms of the burst and
 aspiration were reduplicated. The number of pitch pulses was modeled on the number that
 the model speaker used when asked to speak quickly and slowly. The vowel duration was

TABLE I. Stimuli duration (ms): Experiments 1 and 2. The fricative represents the first point on the continuum (most /ʃ/-like).

	Speaking rate /ʃ/ /k/ /tʃ/ /s/
Slow	145 119 202 198
Base	98 223
Fast	86 235
	Speaking rate /ʃ/ /w/ /ɪ/ /b/
Slow	102 62 151 174
Base	52 161
Fast	35 178

210 similarly adjusted by removing or reduplicating nonadjacent pitch pulses, so as to make
 211 the absolute amount of change in the vowel as close as possible to the absolute amount of
 212 change in the stop consonant. The original /k/ stimulus served as the intermediate duration
 213 stimulus resulting in a 3-way /k/-duration series (short-k/fast speaking rate, intermediate-
 214 k/base speaking rate, and long-k/slow speaking rate), although we make no claims as to the
 215 baseline item actually being half-way between the other two stimuli perceptually. The short
 216 /k/, base /k/, and long /k/ versions of the syllable were then appended to each member of
 217 the 10-item /ʃkas/-/tʃkas/ series. See table I for additional details.

218 *c. Procedure.* Participants completed 1 practice/training block of 60 trials and 4 test
 219 blocks of 90 trials each. The four 90-trial test blocks were comprised of three repetitions of

²²⁰ each of the 30 stimuli (3 /k/ durations X 10-step /ʃ-ʈʃ/ continuum) for a total of 360 trials
²²¹ per participant, or 12 repetitions of each stimulus. Trials in the training block were identical
²²² to those in the test block, but comprised of two repetitions of each stimulus. Responses from
²²³ the training block were not analyzed.

²²⁴ The stimuli were presented to listeners via a lab-created software program that random-
²²⁵ ized stimulus presentation within each block on a Macintosh 7100/AV computer. Stimuli
²²⁶ were presented at a comfortable listening level over Audiotecnica ATH-M40 headphones.
²²⁷ Listeners were prompted with each stimulus and asked to rate the quality of the initial
²²⁸ phoneme on a six-point scale, ranging from “an excellent *sh*” to “an excellent *ch*”, by press-
²²⁹ ing the appropriate button on a computer-controlled response box. Specifically, listeners
²³⁰ were told that they should use 1 for a good, clear *sh*, 2 for an okay *sh*, 3 if they were guess-
²³¹ ing it was *sh*, 4 if they were guessing *ch*, 5 for an okay *ch*, and 6 for a good, clear *ch*. A label
²³² was also posted above the response box to match number to the category endpoints (‘1’ for
²³³ excellent *sh* and ‘6’ for excellent *ch*). Presentation pace depended on the subject’s response
²³⁴ rate. Each trial began 1000 ms after the listener had responded to the previous trial, or
²³⁵ after an interval of 3000 ms following stimulus onset, whichever came first. The experiment
²³⁶ lasted approximately 45 minutes.

²³⁷ **B. Results**

²³⁸ Data were analyzed in the RStudio computing environment (version: 1.4.1103; [RStu-](#)
²³⁹ [dioTeam 2020](#)). Visualizations were created with [ggplot2](#) ([Wickham, 2016](#)). Modeling was
²⁴⁰ conducted and presented using the [lme4](#) ([Bates et al., 2015](#)), [lmerTest](#) ([Kuznetsova et al.,](#)

²⁴¹ 2017), and `broom.mixed` (Bolker and Robinson, 2020) packages. Data analysis decisions
²⁴² (modeling) were not formally pre-registered, but were planned prior to data viewing and are
²⁴³ thus confirmatory and not exploratory. Code to replicate these analyses is available in the
²⁴⁴ project's GitHub repository (<https://github.com/megseekosh/rate-normalization>).

²⁴⁵ To illustrate a possible effect of phoneme duration on rate normalization, we first visualize
²⁴⁶ (1) the proportion of /ʃ/ responses and (2) overall /ʃ/-ness ratings. For the proportion of
²⁴⁷ /ʃ/ responses, the summed proportion of “1”, “2”, and “3” responses (indicating better /ʃ/)
²⁴⁸ were calculated for each participant, for each stimulus item, by dividing the number of
²⁴⁹ “1-3” ratings by all ratings for a given participant/stimulus (Figure 2A). /ʃ/-ness ratings
²⁵⁰ were simply computed for each individual stimulus item presented (item-level) (Figure 2B).
²⁵¹ Overall /ʃ/-ness ratings are only plotted for illustration; statistical modeling was performed
²⁵² on binomial /ʃ/ responses, so effects are reported in log odds.

²⁵³ Figures 2A and 2B suggest that the presence of a rate normalization effect from phoneme
²⁵⁴ duration manipulations. The confidence intervals surrounding the speaking rate conditions
²⁵⁵ (Slow, Base, Fast) do not overlap in the middle, ambiguous section of the continuum. More
²⁵⁶ specifically, we see the effect in the expected direction: slower speaking rates bias /tʃ/
²⁵⁷ responses and higher /tʃ/ ratings, while faster rates bias /ʃ/ responses and higher /ʃ/ ratings.

²⁵⁸ To examine a potential rate normalization effect, we fit a mixed effects model with a logis-
²⁵⁹ tic linking function to predict the log-odds of a /ʃ/ response. This logistic function accounts
²⁶⁰ for the binomially distributed categorical outcome variable (Quené and van den Bergh,
²⁶¹ 2008). Ratings of 1-3 indicate an /ʃ/-bias response and 4-6 indicate /tʃ/-bias response, in
²⁶² line with the instructions that participants received when completing the task. The depen-

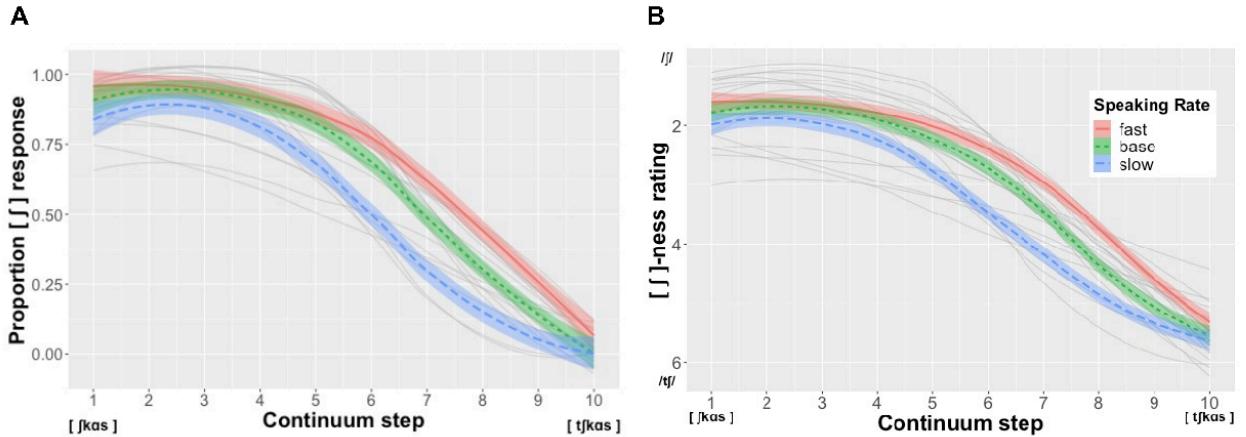


FIG. 2. A: Spaghetti plot of percentage /ʃ/ response by series step and speaking rate: /k/ duration manipulation. B: Spaghetti plot of /ʃ/-ness ratings (1=good /ʃ/, 6=good /tʃ/) by series step and speaking rate: /k/ duration manipulation. Both: Thick, color lines represent group averages by speaking rate and lighter lines represent individual participant responses. Ribbons represent 95% confidence intervals.

dent variable was subsequently re-coded to $/ʃ/ = 1$ and $/tʃ/ = 0$, so positive model coefficients in the summary indicate more /ʃ/ responses. The maximal random effects structure was fit and then backwards pruned until the model converged (Barr et al., 2013). Backwards pruning began by eliminating correlations and then slopes, preferencing the removal of random effects with the smallest variance. Fixed effects were added step-wise, though we additionally evaluated fixed effects via backward fitting and concluded with the same models; this is explained in further detail below. Model parameter significance was determined via a combination of likelihood ratio tests between models, AIC estimations, and p-values (under $\alpha < .05$ criterion) from model summaries.

TABLE II. Model predicting /ʃ/ responses: Experiment 1

Parameter	Estimate	S.E.	z-statistic	p-value	95% CI
Intercept	0.85	0.14	6.00	p<.001	0.58 – 1.13
Rate:Fast	0.53	0.09	6.24	p<.001	0.36 – 0.7
Rate:Slow	-0.81	0.08	-9.57	p<.001	-0.98 – -0.64
Continuum Step	-0.74	0.02	-43.14	p<.001	-0.77 – -0.71

272 The final random effect-only model included Participant-level intercepts; random slopes
 273 of Speaking Rate and Continuum Step by Participant, as well as Speaking Rate by Con-
 274 tinuum Step, did not converge and were removed from the maximal random effect-only
 275 model. The fixed effect of Speaking Rate (modeled categorically with simple coding “Slow,”
 276 “Base,”_(reference) and “Fast”) improved upon the random effects-only model as did Continuum
 277 Step (modeled as a continuous variable and centered at 0 by subtracting the mean) (Table
 278 II). The interaction of Continuum Step and Speaking Rate did not improve upon a model
 279 where these parameters were modeled independently.

280 Unsurprisingly, the proportion of /ʃ/ responses decreased with increased steps along the
 281 /ʃ/-/tʃ/ continuum ($\beta=-0.74$, $z=-43.14$, $p<.001$). For Speaking Rate, there was a greater
 282 proportion of /ʃ/ responses in the Fast condition than the Base condition ($\beta=0.53$, $z=6.24$,
 283 $p<.001$) and a lower percentage of /ʃ/ responses in the Slow condition than Base ($\beta=-0.81$,
 284 $z=-9.57$, $p<.001$), suggesting a rate normalization effect.

285 Overall, these results demonstrate that manipulating /k/ duration, while holding the
286 syllable duration constant, significantly affected the proportion of /ʃ/ responses, suggesting
287 that listeners can normalize for speaking rate over sub-syllabic units such as phonemes.

288 **C. Interim discussion**

289 Experiment 1 demonstrated that two phonemes with obvious acoustic boundaries, /k/
290 and /a/, were treated as separate units during rate normalization. This result implies that
291 the processing unit during rate normalization is something smaller than a syllable. However,
292 /k/ and /a/ are fairly acoustically distinct and separable during running speech. It could
293 be that listeners only rely on sub-syllabic structure to normalize for speaking rate when
294 syllables have a well-defined internal structure. Do listeners likewise normalize for speaking
295 rate over sub-syllabic units that are more difficult to distinguish?

296 Experiment 2 examines a syllable containing phonemes that are much more difficult
297 to segment acoustically: a glide and a vowel. To examine this, we chose a nonce word
298 series that ranged from /ʃwɪb/-/tʃwɪb/. Previous work on similar stimuli—a /swæb/-/twæb/
299 continuum—demonstrated that varying the /w/ duration while leaving the vowel constant,
300 and varying the /æ/ duration while leaving the glide constant, both lead to a change in
301 category boundary location for the initial /s-t/ contrast (Newman and Sawusch, 1996). Yet,
302 as outlined in the justification for experiment 1, this effect could have been driven by the
303 duration of a unit larger than the phoneme, because changing the /w/ duration while leaving
304 the /æ/ constant results in the combined syllabic unit also being longer. Consequently, as
305 in Experiment 1, we again varied the /w/ duration while also altering the /ɪ/ duration in

306 the opposite direction, leading to a series with consistent syllable and word durations. If
307 /w/ and /ɪ/ are treated as separate units during rate normalization like /k/ and /a/ were,
308 despite the acoustic inseparability between /w/ and /ɪ/, then manipulating the duration of
309 /w/ should lead to a rate normalization effect in this series.

310 **III. EXPERIMENT 2**

311 **A. Methods**

312 *a. Participants.* Twenty-two members of the University of Iowa community partici-
313 pated in this experiment for course credit. All were native English speakers with no re-
314 ported history of a speech or hearing impairment and had not participated in experiment 1.
315 Three participants did not respond on at least 80% of the trials, so their data were removed
316 from analysis leaving 19 participants. The remaining participants responded to, on average
317 92.78% of trials (SD=5.18).

318 *b. Stimuli.* Stimulus creation was nearly identical to that in Experiment 1. The same
319 speaker produced the syllable /ʃwɪb/ in the same manner previously described. The initial
320 fricative was separated from the remainder of the syllable, with the boundary being the
321 zero-crossing preceding the first pitch pulse of the /w/. A series of ten items ranging from
322 /ʃ/ to /tʃ/ was created in a similar manner as Experiment 1, by removing successive sections
323 of approximately 10 ms from the onset of the /ʃ/.

324 The syllable /wɪ/ was edited in the same manner as the /ka/ syllable in Experiment 1.
325 Based on spectral analysis, the first 7 vocal pulses were considered part of /w/ rather than

326 the /ɪ/, because these pulses appeared to constitute the /w/ formant transitions (especially
327 those of the first formant). We lengthened and shortened the /w/ and /ɪ/ durations by
328 reduplicating or deleting nonadjacent pitch pulses in the same manner as before, again
329 modeled on the number that the speaker used when asked to speak quickly and slowly. For
330 the shorter /w/, three pitch pulses were removed, whereas four pulses were reduplicated to
331 create the long /w/ (and pitch pules from the vowel were likewise removed or reduplicated
332 in the same manner to keep the syllable duration constant). The original items served as
333 the intermediate duration. The /w/ duration was shorter than that of the /ɪ/, so the same
334 amount of absolute change resulted in a larger change proportionately for the /w/ than for
335 the vowel. The short /w/, baseline /w/, and long /w/ versions of the syllable were then
336 appended to each member of the 10-item /ʃ/-/tʃ/ series. This resulted in three /w/-duration
337 series with a constant syllable and word duration, but varying /w/ (and vowel) durations.

338 c. *Procedure*. The procedure was identical to that of Experiment 1.

339 **B. Results**

340 As in experiment 1, the percentage of /ʃ/ response was calculated for each participant
341 (Figure 4A) and /ʃ/-ness ratings were computed for each individual stimulus item (Figure
342 4B). The visualizations suggest an effect of speaking rate (/w/ duration) upon /ʃ/ responses
343 and /ʃ/ ratings in the same direction as Experiment 1: slower speaking rates bias more /tʃ/
344 responses.

345 For the modeling, to evaluate a potential rate normalization effect, we again fit a mixed
346 effects model with a logistic linking function to predict /ʃ/ responses. All variables were

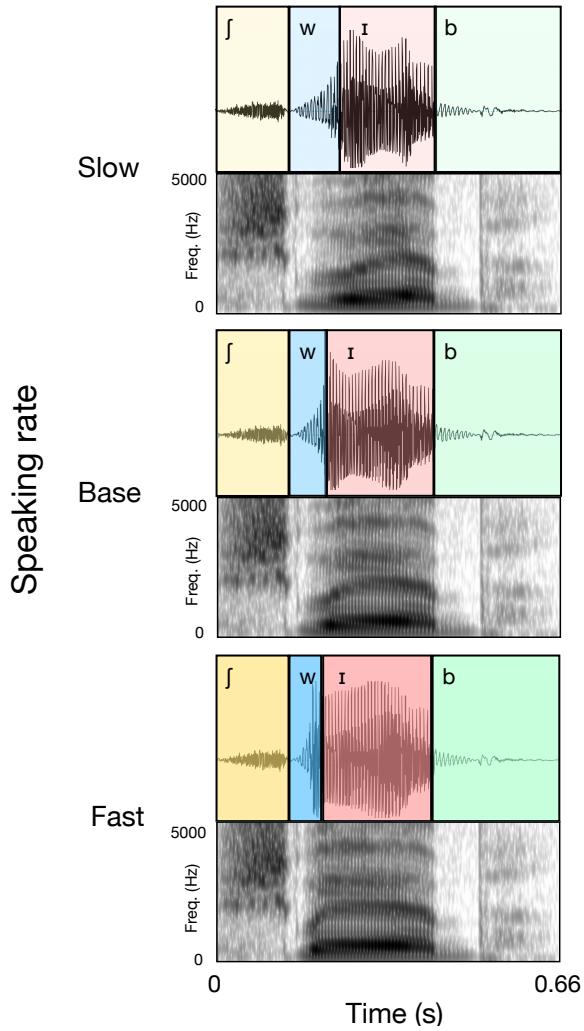


FIG. 3.

Speaking rate manipulations on the basis of /w/ duration and stimuli duration for first step of

series: Experiment 2.

³⁴⁷ coded as an Experiment 1. The model fitting procedure was likewise the same. The final

³⁴⁸ random effect-only model only included intercepts for Participant; more complex random

³⁴⁹ effects structures did not converge. Both the fixed effects of Speaking Rate and Continuum

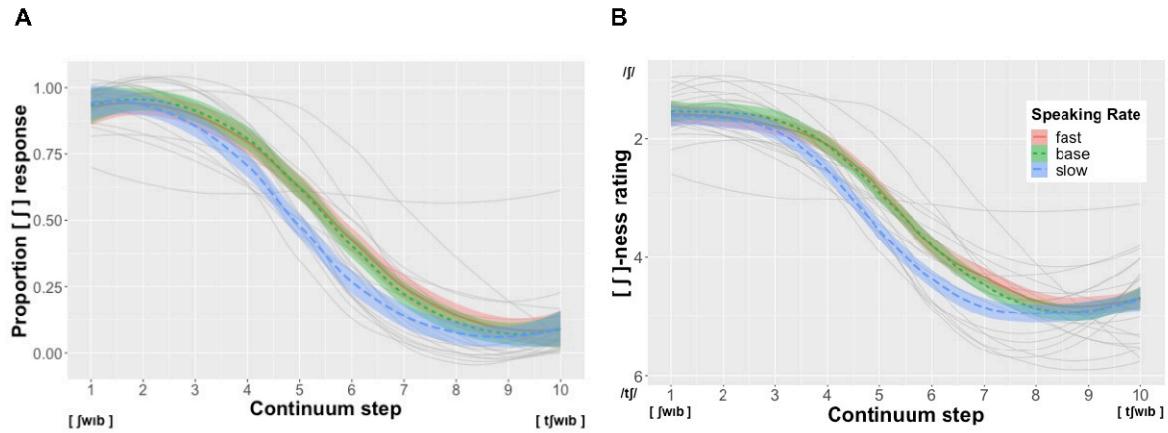


FIG. 4. A: Spaghetti plot of percentage /ʃ/ response by series step and speaking rate: /w/ duration manipulation. B: Spaghetti plot of /ʃ/-ness ratings (1=good /ʃ/, 6=good /tʃ/) by series step and speaking rate: /w/ duration manipulation. Both: Thick, color lines represent group averages by speaking rate and lighter lines represent individual participant responses. Ribbons represent 95% confidence intervals.

350 Step, but not their interaction, improved upon a model fit (as in experiment 1). See Table
 351 III for summary.

352 Once again, unsurprisingly, the proportion of /ʃ/ responses decreased with increased steps
 353 along the /ʃ/-/tʃ/ continuum ($\beta=-0.77$, $z=-42.94$, $p<.001$). For Speaking Rate, there were
 354 fewer /ʃ/ responses in the Slow condition than Base, ($\beta=-0.43$, $z=-4.86$, $p<.001$), suggesting
 355 a rate normalization effect in that direction. However, there were no reliable differences
 356 in /ʃ/ responses between the Fast and Base speaking rate conditions. Consequently, the
 357 results from Experiment 2 show an effect of speaking rate (/w/ duration) upon the perceived
 358 phonetic boundary between /ʃ/ and /tʃ/ for *slow* speaking rates, but not fast, indicating

TABLE III. Model predicting /ʃ/ responses: Experiment 2

Parameter	Estimate	S.E.	z-statistic	p-value	95% CI
Intercept	0.11	0.13	0.85	0.40	-0.15 - 0.38
Rate:Fast	0.05	0.09	0.52	0.61	-0.13 - 0.22
Rate:Slow	-0.43	0.09	-4.86	p<.001	-0.6 - -0.26
Continuum Step	-0.77	0.02	-42.94	p<.001	-0.8 - -0.73

³⁵⁹ that normalization can occur over phonemes in sequences without clear acoustic boundaries
³⁶⁰ (/wi/), but perhaps only when the speaking rate is slow enough to delineate the the units.

³⁶¹ IV. GENERAL DISCUSSION

³⁶² To comprehend speech and language, listeners must compensate for variation across dif-
³⁶³ ferent speakers, in different contexts. Normalization for speaking rate is one important
³⁶⁴ example of this process: it allows listeners to maintain temporal contrasts, such as VOT or
³⁶⁵ vowel length, across different speech speeds and between different speakers. In a pair of ex-
³⁶⁶ periments, we evaluated whether listeners could use information from sub-syllabic units like
³⁶⁷ phonemes—which coarticulation and hypoarticulation often render undefined in the acoustic
³⁶⁸ signal—instead of syllables to normalize for speaking rate (backwards, and proximally). Lis-
³⁶⁹ teners did normalize over phonemes, including acoustically-overlapping phonemes, to factor

³⁷⁰ out speaking rate, demonstrating that sub-syllabic information is used during rate normal-
³⁷¹ ization processes.

³⁷² Work on proximal information in the speech signal for rate normalization has long argued
³⁷³ that normalization occurs over individual phones (Diehl and Walsh, 1989; Newman and
³⁷⁴ Sawusch, 1996). Empirical support was lacking, however, because previous work altered the
³⁷⁵ duration of the carrier syllable and word in addition to the phone. Here we compensated for
³⁷⁶ changes in consonant duration by also changing the nucleus duration. This step allowed us
³⁷⁷ to maintain a consistent syllable duration, avoid the previous experimental confound, and
³⁷⁸ isolate effects of sub-syllabic information on rate normalization. Since we replicated previous
³⁷⁹ work in finding an effect of phoneme duration on this phonetic boundary shift, we can now
³⁸⁰ more definitively say that listeners can use units below the level of the syllable, such as
³⁸¹ phonemes, to compute speaking rate during online speech processing. Furthermore, by also
³⁸² evaluating the effects of acoustic distinctiveness on rate normalization, we were additionally
³⁸³ able to show that this phoneme processing for rate normalization even occurs in sequences
³⁸⁴ like /wi/ that share several acoustic features (periodicity, dynamicity, continuous dynamic
³⁸⁵ formant structure) and are thus less separable. As such, this work expands upon Newman
³⁸⁶ and Sawusch (1996) because the vowel was co-varied in this study and it was manipulated to
³⁸⁷ a much smaller extent ([52 to 197ms] in Newman and Sawusch (1996) versus [202 to 235ms]
³⁸⁸ in the current study), making it all the more interesting that a rate effect was still found.

³⁸⁹ Rate normalization can be activated after just milliseconds of exposure (Reinisch, 2016),
³⁹⁰ and is documented in human and non-human species alike (Welch et al., 2009), suggesting
³⁹¹ that this type of normalization is a low-level auditory process that could be partially domain-

392 general. Finding that listeners can compute speaking rate over sub-syllabic units such as
393 phonemes speaks directly to this idea. Phonemes do not relay a clear acoustic signal. They
394 are indistinct, coarticulated, and reduced—traits that are exacerbated when the features
395 (voicing, stridency) of adjacent phones overlap within syllables. If rate normalization were
396 exclusively or primarily domain-general, it is unclear how listeners could normalize over indi-
397 vidual phonemes. It is possible that listeners may *prefer* or *tend* to normalize over syllables,
398 or relatively more acoustically-reliable components of speech such as word boundaries, but
399 will compute over phonemes in the absence of higher-level information. Our experiments
400 were not designed to contrast listeners' preferred processing unit for rate normalization.

401 It is also possible, as [Bosker \(2017\)](#) suggests, that perceptual normalization for speaking
402 rate could be domain general for some lower-level constructs, such as phonetic boundary
403 shifts, but increasingly language-specific at higher levels such as determining the presence
404 of function words ([Dilley and Pitt, 2010](#)). Nevertheless, the fact that listeners could nor-
405 malize over sub-syllabic information in these experiments suggests that rate normalization
406 processes may be driven by some language-particular experience, instead of the raw acoustic
407 signal alone.

408 Although these results suggest the primacy of sub-syllabic information, such as phonemes,
409 for rate normalization, an alternative interpretation could be that in the presence of am-
410 biguous stimuli, listeners simply weigh information that is immediately available in the
411 signal more than distal, high-level information. (Or that, given that the syllable duration
412 is constant, listeners weigh longer consonants relative to shorter vowels or vice versa.) This
413 interpretation does not require the phoneme to be a perceptual unit. Furthermore, although

414 the current study manipulated phoneme duration, the experimental design still does not dis-
415 associate phonemes from other sub-syllabic entities for normalization, such as the duration
416 of diphthongs or phonetic cues like formant transitions. It could be that listeners are attuned
417 to information at ambiguous points in the signal—wherever that point of ambiguity lies. In
418 the current studies, the ambiguous point was a phoneme boundary but future work may
419 be able to manipulate ambiguity within a single phoneme and elicit a similar normalization
420 effect.

421 Finally, there is an interesting asymmetry in the results whereby listeners in experiment
422 1 used sub-syllabic information to normalize for speaking rate in all speaking rate conditions
423 but listeners in experiment 2 only used this information to normalize in the slow speaking
424 rate conditions. Consequently, it could be that, for less-discriminable syllables, listeners can
425 use sub-syllabic information in slower speaking rates, but not faster because faster speech
426 renders the phonemes too indiscriminable in sequences such as /wɪ/. Another, related topic,
427 could be to compare how normalization unfolds in different phonotactic sequences. Stimuli in
428 these studies were phonotactically illicit in American English—listeners were unaccustomed
429 to hearing onset clusters such as /tʃk/ and /ʃk/. However, it is possible that the default
430 processing strategy during rate normalization could change based on the listener's history
431 with a particular sequence. Perhaps rate normalization occurs more globally, at the lexical
432 or supra-syllabic level, for sequences with high phonotactic probability but more locally
433 at the sub-syllabic level for sequences with lower phonotactic probability. If processing
434 strategies for rate normalization do vary by listener experience, this would be one reason to
435 study the emergence of rate normalization skills in infants—especially given that infants as

436 young as two months are sensitive to duration manipulations (formant transitions) (Eimas
437 and Miller, 1980)—and older children where we may see changes in the default processing
438 strategy changing as a function of child age, vocabulary size, or phonological neighborhood
439 composition.

440 The results of these experiments open up several avenues for future research. First, these
441 experiments only tested American English listeners listening to mostly singleton consonants
442 and monophthongal vowels embedded in nonce words. But other works have found clear
443 effects of language structure and experience on rate normalization (Baese-Berk et al., 2016;
444 Steffman, 2019). Do listeners also normalize over units, like morae, geminates, or diph-
445 thongs that are heavier/larger than phonemes but smaller than syllables? As suggested
446 previously, phonotactic structure is another unexplored aspect of language structure that
447 may be relevant for understanding how listeners calculate speaking rate. Some languages,
448 such as Japanese, tend to have more acoustically “confusable” internal syllable structures,
449 only permitting nasal consonants, and not stops, in coda position for example. This element
450 of Japanese phonotactics renders the transition between nuclei and codas less discriminable,
451 given the shared acoustic properties of nasals and vowels, than in a language like English
452 where a much wider array of codas are permitted (e.g. /s/, /t/). Consequently, if, as in
453 Japanese, the acoustic signature within syllables tends to be more indistinct, listeners could,
454 over time, learn to rely less on individual phonemes for normalization.

455 It will also be important for future work to evaluate processing units for normalization in
456 faster and more naturalistic stimuli as perceptual normalization for speaking rate is likely id-
457 iosyncratic and dependent upon the context and speaker (Goldinger and Azuma, 2003). And

458 more naturalistic stimuli, that contain multiple, co-varying phonetic cues (i.e formant tran-
459 sition duration and frequency), have previously been shown to mitigate rate normalization
460 effects (Shinn et al., 1985). Here we originally hypothesized that listeners would normalize
461 over syllables or other supra-phonemic chunks because both spectral and temporal cues to
462 phonemes become highly confusable and indistinct, especially in fast, running speech, while
463 more global rhythmic cues to syllables may be robust in those settings. And while our ex-
464 periments instead showed reliable effects of sub-syllabic duration on the phonetic boundary
465 shift, the experimental stimuli clearly differed from what listeners would hear and process
466 in real-world contexts. For example, even the manipulated consonant in the “fast” speaking
467 rate condition in Experiment 1 was relatively slow (91 ms) compared to the word-medial stop
468 consonants that listeners might hear in everyday conversation. The duration of the syllable
469 stimuli was relatively long in comparison to typical speaking rates with syllable durations
470 closer to 250-400ms. For extremely fast speech, listeners might rely less on individual phones
471 and more on syllables or words. Faster, naturalistic speech also drives acoustic reduction
472 and heightened coarticulation (Fourakis, 1991; Gay, 1981). However, these acoustic cues did
473 not necessarily accompany the stimuli employed in these experiments as we wanted to iso-
474 late the effects of speaking rate. But extreme reduction in other, more naturalistic listening
475 conditions could lead listeners to normalize over different units.

476 V. CONCLUSION

477 Unlike previous work studying proximal effects on rate normalization, this study manipu-
478 lated speaking rate via phoneme duration while holding the duration of carrier syllables and

479 words constant. We still demonstrated rate effects upon the phonetic boundary shift be-
480 tween /ʃ/ and /tʃ/, both for syllables containing acoustically-distinct /ka/ and -overlapping
481 phonemes /wi/. These results present evidence that listeners process speaking rate over
482 sub-syllabic units, even in the absence of clear acoustic boundaries within syllables, sug-
483 gesting roles of linguistic structure and language experience for perceptual normalization of
484 speaking rate.

485 **ACKNOWLEDGMENTS**

486 The authors wish to thank Jessica Burnham, Jim Sawusch, and Jan Edwards for their
487 assistance with this work.

488 **DATA AVAILABILITY**

489 Analysis scripts to replicate modeling results are included in the affiliated GitHub repos-
490 itory (<https://github.com/megseekosh/rate-normalization>.

491 **CONFLICT OF INTEREST**

492 The authors declare that they have no conflicts of interest.

493 ¹Throughout the paper, we refer to normalization for speaking rate without implying that listeners normalize
494 for all contextual information during perception. We also do not use the term “normalization” to imply
495 that listeners eliminate vs. maintain rate-based information.

496 ²Although both consonant and vowel durations were manipulated, it is unlikely that an observed rate nor-
497 malization effect in the expected direction would be due to manipulations to the vowel. And even if it were
498 vowel-driven, the effect would proceed in the direction opposite to the one expected and would thus be
499 identifiable. Such an effect would run counter to well-known observations in the literature—that shorter
500 durations bias faster speaking rates—and would be illogical from a perceptual point of view.

501

502 Allen, J. S. and Miller, J. L. (1999). Effects of syllable-initial voicing and speaking rate on
503 the temporal characteristics of monosyllabic words. *The Journal of the Acoustical Society*

504 *of America*, 106(4):2031–2039.

505 Baese-Berk, M., Morrill, T., and Dilley, L. (2016). Do non-native speakers use context
506 speaking rate in spoken word recognition? volume 8, pages 979–983.

507 Baese-Berk, M. M., Heffner, C. C., Dilley, L. C., Pitt, M. A., Morrill, T. H., and McAuley,
508 J. D. (2014). Long-term temporal tracking of speech rate affects spoken-word recognition.

509 *Psychological Science*, 25(8):1546–1553. Publisher: SAGE Publications Inc.

510 Barr, D. J., Levy, R., Scheepers, C., and Tily, H. J. (2013). Random effects structure for
511 confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*,
512 68(3):255–278.

513 Bates, D., Maechler, M., Bolker, B., and Walker, S. (2015). Fitting linear mixed-effects
514 models using lme4. *Journal of Statistical Software*, 67(1):1–48.

515 Bolker, B. and Robinson, D. (2020). broom.mixed: Tidying methods for mixed models.

- 516 Bosker, H. R. (2017). Accounting for rate-dependent category boundary shifts in speech
517 perception. *Attention, Perception, & Psychophysics*, 79(1):333–343.
- 518 Crystal, T. H. and House, A. S. (1988). Segmental durations in connected-speech signals:
519 Current results. *The Journal of the Acoustical Society of America*, 83(4):1553–1573.
- 520 Diehl, R. L. and Walsh, M. A. (1989). An auditory basis for the stimulus-length effect
521 in the perception of stops and glides. *The Journal of the Acoustical Society of America*,
522 85(5):2154–2164.
- 523 Dilley, L. C. and Pitt, M. A. (2010). Altering Context Speech Rate Can Cause Words to
524 Appear or Disappear. *Psychological Science*, 21(11):1664–1670. Publisher: SAGE Publi-
525 cations Inc.
- 526 Eimas, P. D. and Miller, J. L. (1980). Contextual Effects in Infant Speech Perception.
527 *Science*, 209(4461):1140–1141.
- 528 Fourakis, M. (1991). Tempo, stress, and vowel reduction in American English. *The Journal*
529 *of the Acoustical Society of America*, 90(4):1816–1827.
- 530 Gay, T. (1978). Effect of speaking rate on vowel formant movements. *The Journal of the*
531 *Acoustical Society of America*, 63(1):223–230.
- 532 Gay, T. (1981). Mechanisms in the control of speech rate. *Phonetica*, 38:148–158.
- 533 Giannella Samelli, A. and Schochat, E. (2008). The gaps-in-noise test: Gap detection thresh-
534 olds in normal-hearing young adults. *International Journal of Audiology*, 47(5):238–245.
- 535 Goldinger, S. D. and Azuma, T. (2003). Puzzle-solving science: the quixotic quest for units
536 in speech perception. *Journal of Phonetics*, 31(3-4):305–320.

- 537 Johnson, K., Flemming, E., and Wright, R. (1993). The Hyperspace Effect: Phonetic Targets
538 Are Hyperarticulated. *Language*, 69(3):505–528. Publisher: Linguistic Society of America.
- 539 Kleinschmidt, D. F. (2016). *Perception in a variable but structured world: The case of speech*
540 *perception*. Unpublished doctoral dissertation, University of Rochester, Rochester, NY.
- 541 Kuznetsova, A., Brockhoff, P., and Christensen, R. (2017). lmerTest Package: Tests in linear
542 mixed-effects models. *Journal of Statistical Software*, 82(13):1–26.
- 543 Lindblom, B. (1990). Explaining Phonetic Variation: A Sketch of the H&H Theory. In
544 Hardcastle, W. J. and Marchal, A., editors, *Speech Production and Speech Modelling*, pages
545 403–439. Springer Netherlands, Dordrecht.
- 546 Maslowski, M., Meyer, A. S., and Bosker, H. R. (2019). How the tracking of habitual rate
547 influences speech perception. *Journal of Experimental Psychology: Learning, Memory, and*
548 *Cognition*, 45(1):128–138.
- 549 Massaro, D. W. and Cohen, M. M. (1983). Consonant/vowel ratio: An improbable cue in
550 speech. *Perception & Psychophysics*, 33(5):501–505.
- 551 McMurray, B., Clayards, M. A., Tanenhaus, M. K., and Aslin, R. N. (2008). Tracking the
552 time course of phonetic cue integration during spoken word recognition. *Psychonomic*
553 *Bulletin & Review*, 15(6):1064–1071.
- 554 Miller, J. L. and Dexter, E. R. (1988). Effects of speaking rate and lexical status on phonetic
555 perception. *Journal of Experimental Psychology: Human Perception and Performance*,
556 14(3):369–378.
- 557 Miller, J. L. and Liberman, A. M. (1979). Some effects of later-occurring information on the
558 perception of stop consonant and semivowel. *Perception & Psychophysics*, 25(6):457–465.

- 559 Newman, R. S. and Sawusch, J. R. (1996). Perceptual normalization for speaking rate:
- 560 Effects of temporal distance. *Perception & Psychophysics*, 58(4):540–560.
- 561 Newman, R. S. and Sawusch, J. R. (2009). Perceptual normalization for speaking rate III:
- 562 Effects of the rate of one voice on perception of another. *Journal of Phonetics*, 37(1):46–65.
- 563 Quené, H. and van den Bergh, H. (2008). Examples of mixed-effects modeling with crossed
- 564 random effects and with binomial data. *Journal of Memory and Language*, 59(4):413–425.
- 565 Recasens, D. (1985). Coarticulatory patterns and degrees of coarticulatory resistance in
- 566 Catalan CV sequences. *Language and Speech*, 28(2):97–114.
- 567 Reinisch, E. (2016). Speaker-specific processing and local context information: The case of
- 568 speaking rate. *Applied Psycholinguistics*, 37(6):1397–1415.
- 569 Repp, B. H., Liberman, A. M., Eccardt, T., and Pesetsky, D. (1978). Perceptual Integration
- 570 of Acoustic Cues for Stop, Fricative, and Affricate Manner. *Journal of Experimental*
- 571 *Psychology: Human Perception and Performance*, 4(4):621–637.
- 572 RStudioTeam (2020). *RStudio: Integrated Development for R*. RStudio, Inc., Boston, MA.
- 573 Sawusch, J. R. and Newman, R. S. (2000). Perceptual normalization for speaking rate II:
- 574 Effects of signal discontinuities. *Perception & Psychophysics*, 62(2):285–300.
- 575 Shinn, P. C., Blumstein, S. E., and Jongman, A. (1985). Limitations of context conditioned
- 576 effects in the perception of [b] and [w]. *Perception & Psychophysics*, 38(5):397–407.
- 577 Steffman, J. (2019). Intonational structure mediates speech rate normalization in the per-
- 578 ception of segmental categories. *Journal of Phonetics*, 74:114–129.
- 579 Summerfield, Q. (1981). Articulatory rate and perceptual constancy in phonetic perception.
- 580 *Journal of Experimental Psychology: Human Perception and Performance*, 7(5):1074–1095.

- 581 Place: US Publisher: American Psychological Association.
- 582 Toscano, J. C. and McMurray, B. (2010). Cue Integration With Categories:
583 Weighting Acoustic Cues in Speech Using Unsupervised Learning
584 and Distributional Statistics. *Cognitive Science*, 34(3):434–464. _eprint:
585 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1551-6709.2009.01077.x>.
- 586 Toscano, J. C. and McMurray, B. (2012). Cue-integration and context effects in speech:
587 Evidence against speaking-rate normalization. *Attention, Perception, & Psychophysics*,
588 74(6):1284–1301.
- 589 Toscano, J. C. and McMurray, B. (2015). The time-course of speaking rate compensation:
590 effects of sentential rate and vowel length on voicing judgments. *Language, Cognition and*
591 *Neuroscience*, 30(5):529–543.
- 592 Trehub, S. E., Schneider, B. A., and Henderson, J. L. (1995). Gap detection in infants,
593 children, and adults. *The Journal of the Acoustical Society of America*, 98(5):2532–2541.
- 594 Turk, A., Nakai, S., and Sugahara, M. (2006). Acoustic segment durations in prosodic re-
595 search: A practical guide. In Sudhoff, S., Lenertová, R., Meyer, S., Pappert, P., Augurzky,
596 I., Mleinek, N., and Richter, J. S., editors, *Methods in empirical prosody research*, pages
597 1–28. Walter de Gruyter, Berlin, New York.
- 598 Wade, T. and Holt, L. L. (2005). Perceptual effects of preceding nonspeech rate on temporal
599 properties of speech categories. *Perception & Psychophysics*, 67(6):939–950.
- 600 Welch, T. E., Sawusch, J. R., and Dent, M. L. (2009). Effects of syllable-final segment
601 duration on the identification of synthetic speech continua by birds and humans. *The*
602 *Journal of the Acoustical Society of America*, 126(5):2779–2787.

- 603 Whalen, D. H. (1990). Coarticulation is largely planned. *Journal of Phonetics*, 18:3–35.
- 604 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New
- 605 York, New York.