

**A Reanalysis of the Voicing Effect in English: With implications for featural specification**

Rebecca L. Morley and Bridget J. Smith

The Ohio State University 1712 Neil Ave.

Columbus OH 43210

614 292-4052

morley.33@osu.edu

## Abstract

The voicing effect is among the most studied and most robust of phonetic phenomena. Yet there remains a lack of consensus on why vowels preceding voiced obstruents should be longer than vowels preceding voiceless obstruents. In this paper we provide an extensive review of roughly seventy years of literature, an analysis of the voicing effect in a corpus of natural speech, and production data from a metronome-timed word repetition study. From this evidence we conclude: that vowel duration differences follow from consonant duration differences; that the voicing effect is largely limited to words of especially long duration; and that preceding vowel duration does *not* reliably cue obstruent voicing under the following circumstances: when obstruent voicing or duration cues conflict; for lax or unstressed vowels; and for most conversational speech. We show that this behavior can be modeled using a competing-constraints framework, where all segments resist expanding or compressing past a preferred duration. Inherent segment elasticity determines the degree of resistance, and duration is ultimately determined by the interaction of these segmental constraints with constraints on target rhyme duration (as a measure of speaking rate), and a preferred C/V duration ratio. Because these constraints are implemented as continuous probability distributions, under- and over- shoot of target duration is possible, and explicit temporal compensation is not required. This account of the voicing effect has a number of implications for phonological theory, especially the central role that the concept of prominence plays in the analysis of underlying features.

*Keywords:* voicing effect, vowel lengthening, final lengthening, temporal compensation, enhancement, Articulatory Phonology

## **A Reanalysis of the Voicing Effect in English: With implications for featural specification**

The terms “vowel lengthening” and “voicing effect” are used to refer to the highly-replicated empirical finding that vowels preceding voiced obstruents tend to be longer than those preceding voiceless obstruents (e.g., Sweet, 1880; House and Fairbanks, 1953; Denes, 1955; Peterson and Lehiste, 1960; House, 1961; Sharf, 1962; Chen, 1970; Raphael, 1972; Klatt, 1973; Lisker, 1974; Raphael, 1975; Raphael et al., 1975; Umeda, 1975; Klatt, 1976; Port, 1976; Fox and Terbeek, 1977; Javkin, 1977; Lisker, 1978; Derr and Massaro, 1980; Fitch, 1981; Walsh and Parker, 1981; Crystal and House, 1982; Krause, 1982; Port and Dalby, 1982; Ohala, 1983; Hillenbrand et al., 1984; Luce and Charles-Luce, 1985; Lisker, 1986; Van Summers, 1987; Kluender et al., 1988; Fischer and Ohde, 1990; De Jong, 1991; Laeuffer, 1992; Crowther and Mann, 1992; Braunschweiler, 1997; Smith, 2002; De Jong, 2004; Kulikov, 2012; Ko, 2018; Tanner et al., 2019; Sanker, 2019; Coretta, 2019; Beguš, 2017). The bulk of the literature focuses on varieties of English, but voicing effects have been documented in a number of different languages. In some, such as Arabic (De Jong and Zawaydeh, 2002), Catalan (Cuartero Torres, 2002), and Russian (Kulikov, 2012), the effect appears to be quite small, while in other languages, such as French (Abdelli-Beruh, 2004; Mack, 1982; Laeuffer, 1992), Swedish (Elert, 1965), and Korean (Chen, 1970), larger differences have been found. It is generally agreed that English (in many of its varieties) exhibits one of the strongest voicing effects, where pre-voiced vowels can be up to 50% longer than their pre-voiceless counterparts (e.g., Chen, 1970; Harris and Umeda, 1974; Mack, 1982). English speaking listeners also exhibit a robust categorical perception effect for final voicing based on preceding vowel duration alone (e.g., Raphael, 1972; Crowther and Mann, 1992; Klatt, 1976; Hillenbrand et al., 1984; Denes, 1955). Other cues to voicing have been shown to be unnecessary in such experiments, or, in fact, to be secondary to vowel duration (e.g., Raphael, 1972; Crowther and Mann, 1992), and preceding vowel duration has been described as the most reliable cue to voicing on final obstruents (Raphael, 1972; Raphael et al., 1975; Luce and Charles-Luce, 1985). Because stops are often unreleased in word-final position, it has also been suggested that a sound change has occurred (or is underway) in which

the contrastive relationship between words like “bad” (bæd) and “bat” (bæt) has shifted away from the final obstruent itself, to be expressed in the duration of the preceding vowel (bæɾ̃ vs. bæɾ̃), at least in phrase-final position (e.g., Klatt, 1976).

In this paper, however, we will argue that the primary cue to obstruent voicing in coda position is the duration of the obstruent itself. Short consonants are perceived as members of the phonologically voiced category, and long consonants, as members of the voiceless. Because the long/short distinction is relative, preceding vowel duration, as an indicator of speaking rate, affects the perception of voicing. While the aerodynamics of voicing maintenance are likely to have been the original source for the duration difference between the obstruents (Ohala, 1983), we find that the duration of the preceding vowel is better predicted by the following obstruent duration, rather than by its voicing (whether phonetic or phonological). We argue that the two phonological categories are distinguished by their inherent elasticity: the less elastic a segment, the more it resists both lengthening and shortening.

The paper is organized as follows. In the next section we provide a comprehensive literature review on the voicing effect and related phenomena. Section 2 contains a corpus study on American English. In Section 3, predictions of our hypothesis are tested using a variable-rate production task. In Section 4 we model the results using continuously violable constraints on duration at the segment and syllable level. Our hypothesis is elaborated in Section 5, where we account for the perceptual side of the voicing effect. We summarize and conclude in Section 6, where we discuss the implications of the present work for theories of phonological contrast.

## **1 Background**

Despite the large amount of research on the phenomenon, the underlying source of the voicing effect remains an open question. There is little consensus on what acoustic or articulatory properties give rise to the observed duration differences. Nor is the effect even consistently described as lengthening, but sometimes as shortening before voiceless consonants, or “pre-fortis clipping” (Gimson, 1970; Wells, 1982). The majority of the perception literature does not even

discuss possible explanations for the effect, or assumes an articulatory source without further discussion (e.g., Raphael, 1975; House, 1961; Ohala, 1983; Klatt, 1976). Belasco (1958), however, speculates that there is a trade-off in the force of production between the vowel and coda consonant of a syllable. When the consonant requires more energy, or effort, the vowel is altered to require less, and vice versa. Thus, voiceless stops, involving forceful release and aspiration, condition shorter vowels, which require less energy. Similarly, Delattre (1962) argues that anticipation of an effortful articulation should shorten the preceding vowel. However, Moreton (2004) and Schwartz (2010) argue essentially the opposite: that it is the spread of “fortisness”, or “hyper-articulation” that shortens the preceding vowel.<sup>1</sup> It has also been claimed that careful, and therefore slower, movements of the vocal cords are required to avoid spontaneous voicing under reduced pressure (Halle and Stevens, 1967); or that the transition from vowel to voiceless obstruent is more rapid than the transition to voiced (Chen, 1970); or that the glottal opening gesture tends to occur a bit earlier for voiceless final consonants to ensure that there is no residual voicing on the consonant (Klatt, 1976). However, clear evidence of differences in energy, effort, or precision between voiced and voiceless obstruents has not been forthcoming. Additionally, as Lisker (1974) points out, the cause and the effect for many articulation-based explanations cannot be assumed. Voiceless stops may involve earlier glottal opening, and a more rapid transition from the vowel, but those facts only explain the shorter duration of the vowel if such shortening is an unavoidable consequence of those properties of articulation. It could as easily be said that such articulatory properties are explained as the consequence of producing the desired voiceless stop.

On the auditory side, Kluender et al. (1988) suggest that longer vowels occurring with shorter voiced obstruents, and shorter vowels with longer voiceless obstruents, is an enhancement effect, reinforcing the length differences of the obstruents, and thus the voicing contrast (Jessen (2001) also proposes an auditory enhancement account). Javkin (1977) posits that vowels are consistently perceived as longer before voiced than voiceless consonants because listeners mis-attribute the glottal pulsing at the beginning of the consonant to the end of the vowel.

---

<sup>1</sup> In Moreton (2004) this also serves to explain why such vowels are more phonetically dispersed, or peripheralized.

However, there seems to be little evidence to support the latter hypothesis, and enhancement explanations are unable to account for why it is preceding vowel length, and not obstruent length, degree of voicing, presence of audible release, or aspiration that are used to enhance the contrastiveness of the obstruents themselves. Sanker (2020) proposes an explanation based on the interaction between acoustics and articulation: a subset of the features in the vowel that are affected by the voicing of the following obstruent (spectral tilt, and intensity contour) also affect perception of vowel duration, presumably for unrelated articulatory reasons. Thus, in the presence of those cues, independent of the presence of a following obstruent, listeners perceive the vowel as being longer/shorter than some baseline duration. This proposal hinges on the source of the duration percept being independent of the voicing, which has not been established, since no articulatory explanation for why spectral tilt and intensity contour affect perceived duration has been proposed.

## **1.1 Production**

The above explanations encounter more problems when the voicing effect is examined more closely. Those of the articulation-based hypotheses that rely on actual vocal fold vibration cannot account for the fact that voicing effects occur even when voiced obstruents are phonetically devoiced (e.g., Walsh and Parker, 1981; Chen, 1970; Fox and Terbeek, 1977). A universal basis for the effect is also called into question by the apparent absence of a lengthening effect in certain languages (Flege, 1979; Hillenbrand et al., 1984; Keating, 1979, 1985).<sup>2</sup> Even in English, with one of the most robust voicing effects measured, durational differences are not found in all contexts. Production studies typically consist of either word lists or brief sentences read by participants in a laboratory setting. In sentence contexts, the target words are often in absolute final position. Such words are also typically monosyllabic, which entails that the target vowel receives primary stress. When some or all of these factors are varied, the voicing effect can be significantly reduced, or disappear altogether: in phrase-medial position (versus phrase-final)

---

<sup>2</sup> Although it should be kept in mind that the cross-language comparisons are not always of like items in these studies.

(Umeda, 1975; Smith, 2002; Crystal and House, 1988; Luce and Charles-Luce, 1985), polysyllabic words (versus monosyllabic) (Umeda, 1975; Port, 1981; Klatt, 1973), lax vowels (versus tense) (Crystal and House, 1988; Luce and Charles-Luce, 1985; Peterson and Lehiste, 1960), unstressed vowels (versus stressed vowels) (Van Summers, 1987; De Jong, 2004), and fast speaking rates (versus slow speaking rates) (Port, 1976; Smith, 2002; Ko, 2018). What all these contexts have in common is that syllables are shorter than in their complementary contexts. We will argue that the voicing effect is highly dependent on absolute duration, because it is only at the longest durations that an appreciable difference in duration between voiced and voiceless obstruents occurs, a difference that is mirrored in preceding vowel duration differences.

## **1.2 Compensation**

In sentence-final position, Luce and Charles-Luce (1985) report voiceless closure durations on average 25% longer than voiced.<sup>3</sup> For monosyllables spoken in isolation, Chen (1970) reports closure durations up to 50% longer. Such large obstruent duration differences occur in the same contexts in which significant vowel duration differences are found, and in the opposite direction (e.g., Klatt, 1976; Umeda, 1975; Miller and Volaitis, 1989; Chen, 1970; Luce and Charles-Luce, 1985). Obstruent duration differences also reduce or disappear in some of the same contexts that vowel duration differences reduce or disappear, calculated either in absolute terms, or as percentages (Luce and Charles-Luce, 1985; Crystal and House, 1982; Miller et al., 1986). These two facts suggest a compensation-based explanation for the voicing effect.

The inverse correlation between differences in obstruent duration and differences in preceding vowel duration was noted early on (Kozhevnikov and Chistovich, 1965; Catford, 1977). However, temporal compensation as an explanation for the voicing effect has been explicitly considered and rejected on a number of separate occasions. Chen (1970), for example, found that syllable duration was not uniform across CVC and CVCC words, indicating that vowels in the latter type of syllable were not shortening to compensate for the added duration of the second

---

<sup>3</sup> This was calculated by taking the average over the two tense vowels and the three places of articulation.

coda consonant. This led him to rule out compensatory mechanisms altogether. Braunschweiler (1997) reached a similar conclusion based on the large duration differences between VC sequences containing a short versus a long vowel. Keating (1985) also rejected a compensation account, based on Polish data in which closure duration, but not preceding vowel duration, varied across voiced-, and voiceless-final syllables. In English, CVC syllables tend to be longer when closed with a voiced obstruent, than with a voiceless, which can also be taken as evidence against a compensation account (e.g., Jacewicz et al., 2009; Luce and Charles-Luce, 1985). Such arguments are based on the assumption that temporal compensation arises from a pressure to keep syllable length uniform (isochrony), meaning that compensation should be total, or near total.<sup>4</sup>

Syllable-level isochrony was originally hypothesized to apply in so-called “syllable-timed” languages like English (e.g., Pike, 1945), and to be the source of a number of apparently compensatory effects. However, while isochronic tendencies exist,<sup>5</sup> it has become clear that uniform timing for syllables is not consistently enforced in English, where syllable duration varies significantly by vowel type,<sup>6</sup> or in any other language that has been investigated (see Krivokapić (2020) for a review).

---

<sup>4</sup> There are, in fact, a handful of studies that report vowel duration differences that are very close to closure duration differences across minimal pairs (in English: Lisker, 1957b; Sharf, 1962; Davis and Van Summers, 1989; in Polish: Coretta (2019); in Georgian: Beguš (2017)). However, across studies, the measured stops were in word-medial position, or in polysyllabic words. Most were produced phrase-medially. Some stops appeared in post-stress position; and for some stimuli, there may have been a syllable boundary between the consonant and the vowel. For all these reasons, the effect sizes were quite small, with duration differences for both vowels and stop closures ranging between 8 and 35 ms.

<sup>5</sup> One possible reason for imperfect compensation could be that isochrony operates at a higher prosodic level than the one being measured. Port et al. (1987) find that moras in Japanese, when produced in isolation, can vary quite widely in duration, yet the words in which those moras appear are much more uniform in length. Small timing adjustments appear to be made at numerous locations within the word, and not necessarily at mora boundaries. Something similar might be true in English, where syllables produced in isolation are clearly not all of the same duration (cf. Chen, 1970). On the other hand, if isochrony itself is driven by a pressure for a uniform rate of information transfer, then it could be the case that true isochrony only holds over semantically defined units, such as phrases, or entire utterances (see, e.g., Aylett and Turk, 2004; Levy and Jaeger, 2007).

<sup>6</sup> Syllables with low vowels are generally longer than those with high vowels (Peterson and Lehiste, 1960); syllables with tense vowels tend to be longer than syllables with lax vowels (Peterson and Lehiste, 1960; Sharf, 1962); stressed syllables are longer than unstressed syllables (De Jong, 2004).



### 1.3 Competition

Certain temporal trade-offs observable at the syllable level in English have been modeled within Articulatory Phonology without assuming isochrony at any level. This is appealing for “compensatory” phenomena that range widely in their degree, from extreme under-compensation, to significant over-compensation (e.g., Elert, 1965; Kristoffersen, 2000; Kavitskaya, 2002; Munhall et al., 1992; Kim and Cole, 2005). In Articulatory Phonology (AP), articulatory units of various sizes are typically modeled as harmonic oscillators with different characteristic frequencies (e.g., Browman and Goldstein, 1990). Phasing relationships between such articulatory units are derived from coupling between the different harmonic oscillators. The same result can be derived from a competing constraints model in which none of the individual constraints on preferred frequencies can be perfectly satisfied, and a “compromise” frequency for the system is adopted that is somewhere between the individual frequencies. Our model is based on the premise that apparent voicing compensation can be treated in an analogous way: as the optimal solution to a set of conflicting timing constraints, but having to do with absolute duration rather than inter-gestural coordination.<sup>7</sup>

#### 1.3.1 *Number of Elements*

One of the central results of AP is the so-called c-center effect, which explains apparent vowel duration differences between syllables with simplex versus complex onsets (Browman and Goldstein, 1988; Nam and Saltzman, 2003; Saltzman et al., 2008). Syllable organization is a function of preferred phasing relationships between successive articulatory gestures. In an English CV syllable, the vocalic gesture is initiated at the target of the consonant gesture. However, a conflict arises when multiple consonants share the same preferred phasing with respect to the following vowel. Satisfying all of them would result in complete merger, or masking of the consonantal gestures. At the same time, the consonants have different timing

---

<sup>7</sup> Browman and Goldstein (1986) themselves adopt one of the phonetic explanations for the voicing effect: because voiceless stops require extra glottal opening and closing, their preceding vowels are shorter.

preferences with respect to one another. The result is that the timing for each consonant is shifted earlier or later by an amount that allows for both the C-V and the C-C phasing to deviate minimally from their preferred values, while also preserving sufficient acoustic cues for all segments. A shift earlier for the first consonant has the effect of lengthening the syllable somewhat, while a shift later for subsequent consonants has the effect of masking more of the vowel. In the latter case, the vowel is acoustically shorter, but not articulatorily.<sup>8</sup>

This type of apparent compensation is found at a number of different unit sizes: words are shorter, the more words there are in the same utterance; stems are shorter the more affixes are attached (Lehiste, 1972);<sup>9</sup> and stressed syllables are shorter, the greater the number of following unstressed syllables (e.g., Fowler, 1981). As with segment-level compensation, stressed syllable duration consistently under-compensates, such that total duration increases (non-linearly) for each additional unstressed syllable (Lindblom and Rapp, 1971; Kim and Cole, 2005). So-called polysyllabic shortening has also been modeled as the result of competing constraints, instantiated as a coupled oscillator system in which the preferred frequency of the oscillator at the lower level of the prosodic hierarchy (e.g., the syllable level) conflicts with the preferred frequency of the oscillator at the higher level of the prosodic hierarchy (e.g., the foot), resulting in a frequency intermediate between the two (O'Dell and Nieminen, 1999).

---

<sup>8</sup> In the case of coda clusters, it has been proposed that something similar to a c-center effect could account for the apparently compensatory behavior (Fowler et al., 1986; Munhall et al., 1992). However, this seems to contradict an earlier finding that only the initial consonant of a coda cluster is coordinated with the vowel, while the remainder are only coordinated with their immediately preceding consonant (e.g., Browman and Goldstein, 1988). The addition of more consonants should thus make the syllable longer but have little to no impact on the acoustic duration of the vowel. Nevertheless, articulatory overlap, leading to acoustic masking, may help explain the duration difference between an open versus a closed syllable. This mechanism alone, however, would not be able to account for the wildly varying degree of compensation (from 13 to 100 ms) reported by Maddieson (1985).

<sup>9</sup> Lehiste (1972) finds that stems are shorter in the affixed form than in isolation (e.g., *sleep/sleepy*). Furthermore, “shortening” increases with the addition of a second affix. This effect interacts with final voicing, such that the amount of “shortening” for voiced-final stems is greater (both absolutely, and proportionally) than that for voiceless-final. A difference is also found between inherently longer and shorter vowels, with longer being more “compressible”. Words are also shorter, the more words in a given utterance, and this interacts with position, with words successively longer the closer they are to the end of the utterance (and thus the phrase boundary).

### 1.3.2 *Intrinsic Duration*

In our proposed model it is preferred durations at the segment level that drive the voicing effect. Something similar may occur in so-called prominence-based compensation. While vowels, as the more expandable segments, seem to compensate for consonant duration, consonants rarely, if ever, seem to compensate for inherent vowel duration differences.<sup>10</sup> However, compensation may occur between two syllables of inherently different durations within the same word.

Final lengthening associated with phrasal boundaries is typically strongest for the segment closest to the boundary, and extends only as far as the onset of the final syllable in most cases (Turk and Shattuck-Hufnagel, 2007; Cambier-Langeveld, 1997; Berkovits, 1993; Hofhuis et al., 1995; Campbell, 1992; Port and Cummins, 1992). However, Cambier-Langeveld (1997; 2000) show that, in Dutch, the penultimate syllable of the final word also sometimes experiences significant lengthening. This happens only when the final syllable is unstressed, or contains a schwa vowel (see also, Turk and Shattuck-Hufnagel, 2007). Katsika (2016) reports a similar finding for Greek, with the articulatory mechanisms for final lengthening appearing to shift towards a stressed penultimate syllable.

The characteristically shorter duration of unstressed vowels seems to prevent them from lengthening to a degree sufficient to satisfy the requirements of phrase-final lengthening.<sup>11</sup> The voicing effect can be described in similar terms: lengthening (also often due to a phrase-final boundary) “shifts” to earlier segments (the vowel) when the final segment (the voiced obstruent) cannot be lengthened sufficiently.<sup>12</sup>

---

<sup>10</sup> Munhall et al. (1992) find a very small difference (on the order of a few milliseconds in consonant duration following vowels of different lengths).

<sup>11</sup> Within the AP framework, prominence-based compensation would arise from the interaction between two different types of basic gestures: the  $\mu$ -gesture, that is associated with stressed syllables, and the  $\pi$ -gesture that is associated with boundary edges. Both are conceptualized as localized “clock-slowness” gestures that result in lengthening (e.g., Byrd and Saltzman, 2003; Saltzman et al., 2008). Thus the interaction, or coupling, between these two gesture types should, in principle, result in prominence-based compensation. However, as far as we are aware, this has not been explicitly modeled.

<sup>12</sup> Although Munhall et al. (1992) suggest that differences in vowel duration preceding voiced versus voiceless obstruents can be explained by differences in the phasing of the two consonants with respect to the preceding vowel, they do not actually provide any evidence in support of this view.

## 1.4 Elasticity

We model this behavior by assigning characteristic elasticities that mediate the degree to which the segment resists pressures to lengthen or shorten from its preferred duration. (see also Cambier-Langeveld (2000) and Miller (1981) on inherent segment elasticity). At short durations/fast speeds, voiced and voiceless obstruents are similar in duration (as observed in the corpus data). At longer durations/slower speeds, vowel differences mirror consonant differences, not because all syllables are the same duration, but because voiced obstruents resist lengthening much more than vowels and voiceless obstruents do. The concept of elasticity is related to the concept of spring stiffness in Articulatory Phonology. The application of a given force will cause a greater perturbation to a spring with a smaller stiffness parameter, resulting in a longer duration than a stiffer spring. However, single gestures do not necessarily correspond to phonemes, and phonemes do not comprise part of the prosodic hierarchy. What this means is that, while an oscillator-coupling model between the syllable and the foot level can be used to predict syllable duration, it is not clear how timing constraints at the syllable level should be distributed among the constituent phonemes. Furthermore, a decrease in speaking rate or lengthening at phrase-boundaries is typically treated as an utterance-level clock-slowness gesture (which may effectively alter the “stiffness” of the component gestures), but there is no differentiation between individual gestures unless they belong to different types of lower-level prosodic units, such as stressed versus unstressed syllables.

What is crucial is that relative elasticity determine the proportion of the syllable that each segment comprises, and that this proportion vary as a function of duration. Expandability should decrease with increasing syllable length, and those segments whose expandability decreases the most rapidly will account for less and less of the total syllable duration, while the more expandable segments will take on an increasingly larger proportion.<sup>13</sup> We call this the Expandability Hypothesis.

---

<sup>13</sup> This model can account for the asymmetry in apparent compensation between vowels and consonants if vowels largely have greater elasticity than consonants.

## (1) The Expandability Hypothesis

All segments have a characteristic elasticity that determines their resistance to lengthening

Resistance to lengthening increases with increasing duration for all segments

Lower elasticity equates with a more rapid increase in resistance

Relative resistance determines the distribution of duration across the syllable

Modeling voiceless obstruents as high elasticity, and voiced obstruents, as low elasticity, we will show that the Expandability Hypothesis parsimoniously accounts for the production data on the voicing effect. It predicts that vowel duration differences should only be seen when there is a complimentary difference in consonant duration, and that the size of the effect should increase with increasing duration. Crucially, the Expandability Hypothesis is also consistent with the perception data, as we demonstrate in Section 5.1.

## 1.5 Perception

It has been repeatedly demonstrated that categorical perception of phonological voicing on obstruents can be achieved by varying only preceding vowel duration (e.g., Raphael, 1972; Crowther and Mann, 1992; Klatt, 1976; Hillenbrand et al., 1984; Denes, 1955). However, the contexts that fail to yield a robust voicing effect in production (fast rate, lax vowels, unstressed vowels, etc.) also fail to show categorical perception based on vowel duration. And while listeners can, and do, make use of preceding vowel duration to identify ambiguous following stops, they make use of other cues as well. While most studies do not directly test different cues against one another, among those that do, the balance of evidence, in fact, comes down against the effectiveness of the vowel duration cue. Both Raphael (1972) and Crowther and Mann (1992) report that preceding vowel duration is stronger than F1 as a cue to voicing. However, Wardrip-Fruin (1982) demonstrates that when preceding vowel duration conflicts with either formant transition cues, or actual vocal fold vibration, the latter dominates. Hogan and Rozsypal (1980) also report that, for certain voiceless-final words, lengthening the vowel does not change

the percept to voiced, but produces no effect, or results in stimuli that sound unnatural. Revoile et al. (1982), using naturally produced stimuli, find that the identification of voiced stops is most strongly disrupted by the removal of vowel offset cues (see also O’Kane, 1978; Nittrouer, 2004), while the identification of voiceless stops is most strongly disrupted by removing the release burst. Similarly, Repp and Williams (1985) find that the addition of a release burst to otherwise ambiguous stimuli reduces voiced responses. Changes to vowel duration, on the other hand, have little effect on voicing perception in their study. Raphael (1981) concludes that vowel duration is only a weak cue to voicing for natural stimuli produced in carrier phrases, and that the effectiveness of various cues is strongly context-dependent.

It was established quite early on that the perceptual boundary between the fricatives /s/ and /z/ in final position is dependent on both consonant and vowel duration (Denes, 1955). However, compared to the number of studies that test perception based on preceding vowel duration alone, there are relatively few that manipulate, or even report, the duration of final stops. These studies, for whatever reason, also tend to be cited less frequently. However, Raphael (1981) found that swapping closure durations for naturally produced “peg” and “peck” effectively switched the voicing percept for the two tokens. Repp and Williams (1985) similarly found an effect of overall closure duration on the perception of voicing on stop-final syllables followed by a stop-initial syllable (e.g, “lab coat” vs. “lap-coat”).

## **1.6 The “Voicing Effect” in Initial and Medial Position**

If the Expandability Hypothesis is correct, then it should apply equally well to obstruents in non-final position. Historically, the behavior of stops in medial position, and initial position especially, have been described in very different terms to the behavior of word-final stops. It is not clear, however, that this distinction is theoretically warranted. While syllable position is expected to affect the realization of a given segment, we assume that voiced stops have the same inherent elasticity regardless of the position in which they occur.

Like word-final stops, medial stops are post-vocalic, as well as subject to neutralization of

both voicing and aspiration, resulting in productions that are essentially just short periods of silence (oral cavity closure). Vowel duration has also been shown to be a sufficient cue to voicing in medial position. In this literature, however, it is standard to describe the perceptual boundary in terms of the ratio of closure duration to preceding vowel duration (e.g., Port and Dalby, 1982; Port, 1979, 1981; Lisker, 1957b). The C/V ratio effectively normalizes stop duration relative to estimated speaking rate. This type of normalization presumably also applies in final position, yet because final closure duration is not typically measured, let alone systematically varied, it has become tacitly assumed that consonant duration plays no role in perception – at least for stops, and pre-pausally. However, the size of the “voicing” effect seems to be the only real difference between medial and final position. Voiceless closure durations are reported to be from 30-45 ms longer than voiced;<sup>14</sup> and pre-voiced vowel durations, 25-35 ms longer than pre-voiceless (Lisker, 1957a; Sharf, 1962; Davis and Van Summers, 1989). These smaller values are to be expected, given that word- and phrase- final lengthening do not apply,<sup>15</sup> and words are, by definition, polysyllabic.

There is a very large body of work devoted to word-initial stops, in which VOTs in pre-stressed position are typically measured. The relationship between VOT and following vowel duration, however, has been much less studied. When post-stop vowel duration is manipulated in perception experiments, it is almost always done as part of a speaking rate study in which stop duration is inversely co-varied, making it difficult to determine the relationship between consonant and vowel (e.g., Miller and Baer, 1983; Miller and Volaitis, 1989; Volaitis and Miller, 1992). However, resistance to lengthening under slowed speaking rate is also observed for voiced stops in initial position. Furthermore, in the handful of studies that vary vowel, rather than syllable, duration the results are qualitatively similar to what is found in medial and final position:

---

<sup>14</sup> These values are for labial and velar place. Coronals are typically flapped in this environment and show little to no duration differences.

<sup>15</sup> In medial position the consonant’s syllabic affiliation is ambiguous. The onset maximization principle (e.g., Clements and Keyser, 1983) places a single medial consonant in the onset of the following syllable. However, there is evidence that stress and sonority both affect syllabification, such that a medial consonant will be syllabified in the coda of a preceding syllable if that syllable is stressed, and the following syllable is not (e.g., Treiman et al., 1994; Eddington et al., 2013).

perception of voicing switches based on vowel duration, but only for longer vowels. Viswanathan et al. (2019) report a significant difference in voicing perception between vowels of 175 and 225 ms., but no difference between longer vowels, at 225 and 275 ms., and no difference between shorter vowels, at 125, 150 and 175 ms. Toscano and McMurray (2015) find a significant difference in response rate across the same boundary, between vowels of 189 ms. and those of 377 ms.

VOT ranges for stops in initial position are comparable to what we see for closure durations in final position in the corpus and in the production study (Section 3): from roughly 50-150 ms for the voiceless stop, and 10-50 ms for the voiced (Allen and Miller, 1999; Pind, 1995; Miller and Baer, 1983; Miller et al., 1986). However, post-voiced vowels were only 10-19% longer than post-voiceless, compared to 30-40% for pre-voiced versus pre-voiceless. As with medial position, we attribute some of the smaller effect size to the lack of a final lengthening effect.<sup>16</sup> However, the difference in gestural timing relationships between vowel and onset versus vowel and coda is also expected to contribute to this outcome. Additionally, onsets don't add to syllable weight, and may vary less than codas under changes in speaking rate. Using the same model, but incorporating these two changes, we are able to account for the "voicing" effect in initial position (see Appendix (C)).

## **2 A Corpus Study**

In this section we provide an in-depth analysis of the voicing effect in conversational speech, using the Buckeye Corpus (Pitt et al., 1997). Although the corpus is not balanced, it provides much more data, and a larger range of speaking rates and contexts than any single laboratory experiment. A corpus study allows us, first, to quantify the voicing effect in actual

---

<sup>16</sup> While lengthening occurs preceding, as well as following, prosodic boundaries, the effects are not the same. The consonant immediately following a prosodic boundary shows lengthening, but the vowel following that consonant typically does not (e.g., Fougerson and Keating, 1997; Cho and Keating, 2009; Kim and Cho, 2012). Byrd et al. (2005) find that coda consonants show a larger difference in duration when compared across medial and boundary position than onsets. Interestingly, the difference seems to lie in the fact that more of the coda gesture is lengthened. For both codas and onsets, the portion of the gesture closest to the boundary is lengthened the most – for codas the release portion, and for onsets, the constriction portion. However, the constriction portion for the coda consonant is also lengthened to a lesser degree, while the release for the onset is not (or at least, not consistently).



usage. Secondly, it allows us to probe more deeply into the factors that affect the realization of the effect. Conversational styles of speech are expected to exhibit considerable reduction in the realization of individual words, some component sounds of which may be entirely missing (e.g., Harris and Umeda, 1974; Johnson, 2004; Jurafsky et al., 1998). This reduction could neutralize small differences in duration that result from an underlying voicing effect. And previous studies with read scripts have shown a reduced voicing effect in comparison to single sentence or word list productions (Crystal and House, 1982, 1988). We find that there is an inconsistent effect of voicing that is dependent on model structure. When consonant duration is added to the model the voicing effect either goes away, or switches direction. Consonant duration, on the other hand, is consistently negatively correlated with vowel duration. Furthermore, both effects are found to participate in interactions with speaking rate and frequency, exhibiting the expected dependence on duration predicted by the Expandability Hypothesis.

## **2.1 The data**

The Buckeye Corpus consists of segmented and transcribed sound files. These are taken from interviews, each lasting about an hour, with 40 different speakers, all middle-class and Caucasian, who are also natives of central Ohio. Intertranscriber reliability of the phonetic symbols for stops and fricatives was reported for a sample of the Buckeye Corpus at 91.2% and 92.9%, respectively. For the unanimously transcribed subset of this sample, segmentation boundaries differed an average of 16 ms. (Pitt et al., 2005). However, Raymond et al. (2002) report a difference in segmentation agreement for shorter versus longer phones. 73% of phones that were longer than average agreed within 20% of the average length of the two phones on either side of the segment boundary, whereas only 50% of phones that were shorter than average agreed within 20%. Shorter phones were thus proportionally less consistently transcribed than longer phones. In the absence of a consistent bias in the placement of the boundary, this should not affect our analysis for the most part. However, errors in segmentation at the short end of the continuum may wash out a small voicing effect. The segmentation of the vowel and final

consonant are inherently negatively correlated; an error in which the vowel duration is longer will also produce an error in which the final obstruent is shorter (assuming no error in determining the endpoint of the obstruent). If such errors are more likely to happen with voiced obstruents, then we would expect an augmented voicing effect.

From the Buckeye Corpus we extracted all monosyllabic words of the form (C)onsonant-(V)owel-(C)onsonant ending with one of the following obstruents: voiced (d,b,g,z,ʒ,v) or voiceless (t,p,k,s,ʃ,f). CVC words were selected because they were expected to show the largest voicing effect. Complex onsets were excluded to eliminate potential variability. No nasalized or rhotacized vowels were included, to be sure that each word had exactly three underlying segments. Only tokens that were both phonemically and phonetically CVCs were included. For example, tokens of “past” realized as [pæʃ], and tokens of “allowed” realized as [lɑʊd] were both excluded. Because the transcription of the corpus is quasi-phonetic, we constructed a dictionary of citation forms to ensure that the phonological voicing category was correctly assigned to each word. Because there were no words ending in voiced dental fricatives, those ending in voiceless dental fricatives were also removed. The vowel /ɔɪ/ was also excluded for reasons of data sparsity. 20.3% of the stops in the remaining data were transcribed as glottalized (tq), which could represent a glottal stop or unreleased stop with glottalization on the vowel, but less than 1% of those were underlyingly voiced, so all such tokens were removed from analysis. Affricates were excluded due to the possibility that they might straddle a word boundary.

In Figure 1 raw vowel durations for the set of CVC word tokens used in the following analyses are plotted as a function of the voicing feature of the final obstruent. The density plot on the right suggests that there is a very small effect of voicing at the longest durations. However, the actual counts given in the left panel show that there are never more voiced than voiceless tokens at any duration. This is due to the fact that there are considerably more word tokens with (phonemically) voiceless coda obstruents (over twice as many as voiced tokens, although there are more voiced than voiceless fricative tokens. See Appendix A). Vowels preceding voiceless obstruents have a slightly longer mode than those preceding voiced obstruents, and at the longer

durations (above 175 ms.), the *relative* proportion of the pre-voiced distribution is larger than the pre-voiceless. For the most part, however, the two distributions are completely overlapped, showing no transparent voicing effect.

\*\*\*\* FIGURE 1 ABOUT HERE \*\*\*\*

## 2.2 Model Factors

If a voicing effect does exist in these data, it is masked by factors that affect vowel duration more strongly. The following factors, each of which is known to affect segment duration, are included in the statistical model of vowel duration. Because the analysis was limited to CVC words, stress and word length are not included.

- **INHERENT VOWEL CLASS:** Tense and lax vowels in English are differentiated in part by duration. /ɪ, ɛ, ʊ, ʌ/, all lax vowels, are reliably shorter than their tense counterparts (e.g., Peterson and Lehiste, 1960; Klatt, 1976; Stevens and House, 1963). /æ/, although technically lax, has much longer durations than any other lax vowel (e.g., Hillenbrand et al., 2000; Crystal and House, 1988), and is actually a diphthong in some dialects, thus it is grouped with other inherently long (tense) vowels. Reduced or absent voicing effects have been reported for both unstressed and lax vowels (Umeda, 1975; Crystal and House, 1982; De Jong, 2004). Vowel class is modeled as a factor with 2 levels: Short (ɪ, ɛ, ʊ, ʌ), and Long (all other vowels, namely, i, e, u, a, o, ɔ, æ, ɔɪ, aɔ), coded as 1, and -1, respectively.

- **VOWEL HEIGHT:** Because high vowels tend to be shorter than low vowels, this can affect the realization of the voicing effect. Vowel height is a factor with 2 levels : high (i, u, ɪ, ʊ), and non-high (all other vowels), coded as -1, and 1, respectively.

- **SPEAKING RATE DEVIATION:** The z-scored average difference between expected and observed duration was used as a proxy for rate difference (see Gahl et al., 2012; Priva, 2017). Mean segment durations by speaker were taken as the expected value for each phoneme. This was computed over all tokens (regardless of word position or phonological context) and for all words

in the Buckeye Corpus (not just the CVC words used in the analysis) The deviation was calculated for each segment within a word, and then averaged. Because this is actually a duration measure, a positive difference indicates that the individual segments within the word are generally longer than their average durations, and thus that the speaking rate is slower than average. Speaking rate deviation is modeled as a continuous variable.

- **WORD FREQUENCY:** More frequently used words generally have shorter durations than less frequently used words, and both vowels and consonants within those words are affected (e.g., Jurafsky et al., 2001; Fidelholtz, 1975; Fosler-Lussier and Morgan, 1999; Hooper, 1976; Pluymaekers et al., 2005). Function words, generally the most frequent and the most contextually predictable words, are consistently shorter than content words (Bell et al., 2009; Umeda, 1975). Because the difference in frequency between content and function words is several orders of magnitude, Zipf scores,  $\log_{10}(\text{Frequency})$ , were used. Word frequencies were supplied as counts per million from the SUBTLEX corpus (Van Heuven et al., 2014). Log-frequency is modeled as a continuous variable.

- **PHRASAL POSITION:** Prosodic boundaries have the effect of lengthening adjacent segments. The greater the number of nested phrases marked by the boundary, the greater the degree of lengthening, and the further its spread (Oller, 1973; Wightman et al., 1992; Fougeron and Keating, 1997; Byrd and Saltzman, 2003). Because the Buckeye Corpus is not annotated for syntactic boundaries, tokens were classified only as pre-pausal or non-pre-pausal, based on the end of a transcribed utterance. Pre-pausal position is expected to show the largest lengthening effects (see, e.g., Crystal and House, 1988; Klatt, 1975). The following tags were used to identify a boundary: SIL (silence), E\_TRANS (end of phonetic transcription), IVER (interviewer speaking), VOCNOISE (non-speech sound such as a cough, or laugh). Position is modeled as a factor with 2 levels: phrase-final and non-phrase-final, coded as 1 and -1, respectively.

- **PHONETIC VOICING:** Phonetically voiced segments exhibit acoustic evidence of voicing, as transcribed by corpus annotators. Phonetic voicing is modeled as a 2 level factor: voiced, and voiceless, coded as 1 and -1, respectively.

- **PHONEMIC VOICING:** Phonemic voicing refers to the category of the phoneme in the citation form of the word. Phonemic voicing is modeled as a 2 level factor: voiced, and voiceless, coded as 1 and -1, respectively.

- **OBSTRUENT TYPE :** A 2-level factor: stop, or fricative, coded as 1 and -1, respectively.

- **CONSONANT DURATION:** The duration of the final consonant as measured by corpus annotators. This is a continuous variable.

Both phonetic and phonemic voicing were included in the model because it was not known if there might be an effect of actual voicing above and beyond the effect of phonological voicing. We soon found that phonetic voicing did not differ appreciably from phonemic voicing, and it was dropped from the analyses. For the remainder of the paper, “voicing” will refer to phonological voicing. Although vowel duration based on vowel quality is not actually binary, data sparsity for certain low-frequency vowels makes using vowel quality itself problematic as a finer-grained determiner of inherent duration<sup>17</sup>.

We use a measure of speaking rate that is based on the duration by which each word differs from the average of the mean values of its individual phonemes. This measure is based on similar metrics in which an expected duration is compared to an observed duration (e.g., Priva, 2010; Gahl et al., 2012; Gahl and Baayen, 2022). Such measures are used in order to make estimates of rate as independent of segment duration as possible. In principle, the two will be linked because it is not possible to avoid a confound when using acoustic data that has not been generated in explicitly rate-controlled contexts. Ambiguity in attributing duration differences among different factors is partially resolved by calculating deviation over the entire word. For example, a duration deviation in the vowel preceding a voiced coda may be attributable to the voicing effect, or to speaking rate. If the latter, then the onset and coda of the word in question should also be longer than average. If not, then the calculated speaking rate deviation will be smaller, leaving all, or most of the variance, to the voicing effect.

---

<sup>17</sup> Note that Tanner et al. (2019) include a random intercept for vowel quality, which we adopt in attempting to replicate their results.

## 2.3 Methods

All statistics were performed using the lme4 package in R. Linear mixed effects models were run using the function lmer, fit by REML. T-tests used Satterthwaite's method, and the lmerTest function was used to obtain estimated p-values. All continuous numerical variables were log-transformed and mean-centered to approximate a normal distribution with a mean of zero. Following Tanner et al. (2019), we normalize by dividing by two standard deviations. Random intercepts for word and speaker were included in all models. Place of articulation of the final obstruent, although known to affect consonant duration, was too small of an effect to significantly improve model fit, and was therefore left out of the final model. Due to the asymmetric distribution of the data, it was not possible to use paired data in analyzing the voicing effect. All factors were sum-coded so that each individual factor was assessed at the mean value of all other factors. Three-way interactions were avoided for reasons of interpretability as well as model convergence.

## 2.4 Results

Tanner et al. (2019) report a voicing effect for phrase-final monosyllabic words in the Buckeye Corpus. Therefore, we begin by attempting to reproduce their analysis as closely as possible. Nevertheless, there remain differences in how the samples were selected. Tanner et al. (2019) defined pre-pausal as preceding a silence of at least 150 ms of silence, although they do not specify how they determined silent intervals. Our measure was more conservative, using only tokens with segmentation that indicated that the talker had stopped speaking long enough for the transcription to include an interval with silence, noise, or change in talker. This measure corresponds more closely to utterance-final position. Tanner et al. (2019) also excluded vowel tokens under 50 ms. We kept in tokens under 50 ms because that is not unusual for casual speech (there were 167 such tokens across stops and fricatives in utterance-final position). We also did not include stops that had been transcribed as flapped or glottalized, which further reduced the sample size, but gave us more homogeneous data. Tanner et al. (2019) included tokens with

complex onsets and codas, while we analyzed only CVC forms. It is not clear whether they included reduced forms of polysyllabic words, or corrected for phonetic transcriptions of voiceless obstruents when the underlying value was voiced.

Our final token count was 3200, while Tanner et al. (2019) analyzed 5500 tokens. Their model includes interactions between voicing and frequency, voicing and vowel type, voicing and obstruent type, and voicing and word class. They included random intercepts for speaker, word, and vowel quality, as well as random slopes for speaker by frequency, vowel type, obstruent type, word class, and by the interaction of voicing and obstruent type. Random slopes were also included for word by both speaking rate measures that they used. See 2 for the full model specification. This model overfits our data and does not converge. Therefore, we simplified the model to what is shown in (3).

- (2) Vowel Duration ~ Local Speaking Rate + Global Speaking Rate + Frequency + Vowel Class + Vowel Height + Voicing + Stop Type + Word Class + Voicing:(Frequency + Vowel Class + V Height + Stop Type + Word Class) + (Frequency + Vowel Class + Vowel Height + Stop Type + Voicing + Word Class + Voicing: Stop Type | Speaker) + (Local Speaking Rate + Global Speaking Rate | Word) + (1 | Vowel Quality)
- (3) Vowel Duration ~ Vowel Height + Vowel Class + Voicing + Speaking Rate Deviation + Obstruent Type + Frequency + Voicing:(V Height + Speaking Rate Deviation + Obstruent Type + Frequency + V Class) + (Voicing + Obstruent Type + Frequency|Speaker) + (Speaking Rate Deviation|Word) + (1|Vowel Quality)

For each variable, the average value of its levels (if a factor), or of its range of values (if a continuous numerical variable) was the baseline for analysis. This allows us to conceptualize the results in a way that is similar to ANOVA, where each effect is an adjustment to the average value for the model. For example, the effect of Vowel Class is determined by whether the average duration of the class of Long vowels is significantly different from the global vowel duration average, calculated over both Long and Short vowels.

As expected, there was a significant main effect of speaking rate. See Table 1. Longer vowel durations were found at slower than average speaking rates. Word frequency also had the expected negative effect on vowel duration, such that words with higher than average frequency had shorter vowel durations. As predicted, long vowels were also longer than the average of long and short vowels, and high vowels were shorter than the average of high and low vowels. Phonological voicing, however, did not have a significant effect on vowel duration. Nevertheless, in the interactions we see a positive adjustment to voicing at slower speaking rates, and a negative adjustment in the effect in higher frequency words. These two terms go in the expected direction: an emergent voicing effect for tokens of longer duration, whether low-frequency, or spoken slowly.

\*\*\*\*\* TABLE 1 ABOUT HERE\*\*\*\*\*

It is not entirely clear why we fail to find a voicing effect where Tanner et al. (2019) did find one. It could be because we included very short tokens, or because we had a smaller sample of data. However, removing tokens under 50 ms did not appreciably change the results, and a model of non-phrase-final tokens did not produce a significant voicing effect either. The local speaking rate measure used by Tanner et al. (2019) is a count of the number of segments within an inter-pause interval containing the target word, therefore it is much less local than the one we use. This type of measure also does poorly in estimating rate for pre-pausal tokens that are subject to final lengthening, and at the other end of the continuum, with the very shortest tokens. These two ranges show non-linear behavior for this type of speaking rate measure, which is presumably why linear regression models do not capture them well (see Gahl, 2009). If less of the over-all variance is attributed to speaking rate, this might account for the difference between our results.

Approaching the analysis from the other direction, we ran a highly simplified model with no interactions, but keeping random intercepts and slopes. See Table 2. All model formulae are included in Appendix D.

\*\*\*\*\* TABLE 2 ABOUT HERE\*\*\*\*\*

A significant effect of voicing emerged for both phrase-final tokens and phrase-medial



tokens, run separately. Obstruent type also became significant, in the negative direction, indicating that vowels were longer preceding fricatives than stops. Thus, reducing the model itself, rather than the sample, seems to result in the expected effect of voicing.

Out next step was to add consonant duration to the model. The interaction terms had to be modified slightly so as to converge. The result is a (negative) significant effect of consonant duration, but voicing now becomes significant in the wrong direction. The interaction between speaking rate and voicing is still significant, while the interaction between voicing and frequency is marginal. There is also a significant interaction between voicing and consonant duration in the negative direction. See Table 3. To determine the degree of correlation in our continuous variables we used the vif function from the faraway package for R. Consonant duration, speaking rate and frequency, had respective variance inflation scores of 1.62, 1.77, and 1.14. VIF scores below 5 are usually considered to be unproblematic.

\*\*\*\*\* TABLE 3 ABOUT HERE\*\*\*\*\*

For the simplified model (without interactions), adding consonant duration has the effect of rendering the voicing effect non-significant. See Table 4. Meanwhile, both the simplified model and the interaction model with consonant duration, but with voicing removed, produce consistent effects: consonant duration is negatively correlated with vowel duration, and all other terms are significant in the expected directions. We suggest that the variable results from the voicing factor - non-significant, significant, and significant in the wrong direction – indicate the fragility of the voicing effect. In contrast, a significant amount of variance in vowel duration can be robustly captured by consonant duration. While both vowel and consonant duration increase with decreasing speaking rate, longer consonant durations also predict shorter vowel durations. Our final model was run on phrase-medial tokens to corroborate these results. See Table 5. The VIF scores for the non-phrase-final sample, for consonant duration, speaking rate, and frequency, respectively, are 1.37, 1.53, and 1.14.

\*\*\*\*\* TABLE 4 ABOUT HERE\*\*\*\*\*

\*\*\*\*\* TABLE 5 ABOUT HERE\*\*\*\*\*

## 2.5 Obstruent duration

In laboratory settings, voiced obstruents are consistently found to be shorter than voiceless obstruents (e.g., Klatt, 1976; Umeda, 1975; Miller and Volaitis, 1989; Chen, 1970; Luce and Charles-Luce, 1985). In conversational speech, however, the duration distributions are almost completely overlapped, as can be seen in Figure 2, which shows the raw consonant durations as both counts and probability densities. These distributions look strikingly similar to the vowel duration distributions from Figure 1.

\*\*\*\*\*FIGURE 2 ABOUT HERE\*\*\*

A regression model with obstruent duration as the dependent variable shows that voicing is significantly negatively correlated for the non-phrase-final data. See Table 6. There are also significant interactions of voicing with rate and with frequency, as was seen in the vowel model. The negative effect of voicing is increased at slow rates, but decreased for high frequencies. These interactions are consistent with the observation that consonant duration differences are larger for longer tokens. Although these results, in and of themselves, cannot prove the hypothesis that vowel duration differences arise from consonant duration differences, they are consistent with that hypothesis.

\*\*\*\*\*TABLE 6 ABOUT HERE\*\*\*

## 2.6 Summary & Discussion Of Corpus Results

The corpus results show that, in conversational speech, we do not consistently see the expected effect of voicing on vowel duration, failing to replicate the finding of Tanner et al. (2019). It has generally been found that the voicing effect is reduced or absent in higher-frequency words, at faster rates, for inherently shorter vowels, phrase and word medially, for unstressed vowels, and for vowels in polysyllabic words – all environments that exert a shortening effect to some degree. We also see a dependence of the voicing factor on absolute durations in interaction terms with speaking rate and frequency. Although vowel duration and consonant duration both increase as absolute word duration increases, vowel duration is also

inversely correlated with consonant duration – consistently across all models.

While the effect of voicing is only present in interactions, or when consonant duration is not included in the model, laboratory production studies find that pre-voiced vowels can be up to 50% longer than pre-voiceless vowels (Peterson and Lehiste, 1960; Mack, 1982; House, 1961; Luce and Charles-Luce, 1985; Umeda, 1975). Similarly, voiceless stop closure durations can be from 25% to 50% longer than voiced stop closures (Chen, 1970; Luce, 1986). These discrepancies can be explained by the large difference in absolute durations between the corpus and the laboratory. Vowel durations in those studies were reported in the range of 175 to 300 milliseconds (Peterson and Lehiste, 1960; Mack, 1982; House, 1961; Luce and Charles-Luce, 1985; Umeda, 1975), with voiceless stop closures ranging from 95-140 milliseconds (Luce and Charles-Luce, 1985; Chen, 1970). For the vowel tokens in the Buckeye Corpus, on the other hand, durations this long are rare. Among the set of CVC words ending in voiced obstruents, less than 7% reach durations of 200 ms or above. Even restricting the sample to just characteristically longer vowels, only 13% of such tokens fall in this range. Median vowel duration over the complete set of CVC words used in this study is only 83 ms. Median vowel duration for just the voiced tokens is actually lower than that, at 75 ms. Similarly, only 9% of CVC-final voiceless stops reach durations of 100 ms or above in the Buckeye Corpus, while the median closure duration is 46 milliseconds.

### **3 A Production Study**

We take the corpus results, in conjunction with the production literature as a whole, to provide strong preliminary support for the Expandability Hypothesis. However, because paired data are not available in the corpus,<sup>18</sup>our predictions must be confirmed in a setting where sources of variation can be controlled for. In this section we report the results of a production experiment in which we asked native English speakers to repeat a series of CVC minimal pairs at varying rates. This allows us to directly compare the lengthening behavior of final voiced obstruents to

---

<sup>18</sup> And may not be readily available to listeners in any case, if the distributional properties are similar across other spoken corpora.

that of final voiceless obstruents. Additionally, the difference in obstruent duration can be compared to the difference in vowel duration as a function of speaking rate.

### **3.1 Methodology**

#### **3.1.1 Participants**

All participants were undergraduate students at The Ohio State University who were given course credit for completing the experiment. A total of 45 participants were run: 25 were female, and 20 were male. The average age of participants was 21. Of this group, 11 were excluded from analysis for the following reasons: they reported hearing issues (3); they reported learning a language other than English before the age of 7 (7); they did not learn English until after the age of 7 (1).

#### **3.1.2 Procedure**

Participants were seated in front of a computer monitor inside a sound-attenuated booth. Continuous audio was recorded from a desktop microphone using the sound editing software Audacity.<sup>19</sup> Participants were instructed that they would be asked to speak into the microphone in response to prompts on the computer screen. The entire experiment took less than an hour to complete.

The experiment began with a practice block to acclimate participants to the experimental task, and the different repetition rates involved. Prior to the start of the practice block, participants were given the following instructions:

*A + sign will appear on the screen. It will be black to begin with, then will change to red, and keep alternating. Your job is to repeat the word on the screen every time + changes color. Try to use the entire time that the + does NOT change color to say the word. Keep going until the flashing stops. Press any key when you are ready to practice with the word “lab”.*

---

<sup>19</sup> Available at <http://audacity.sourceforge.net>.

For the first trial, participants saw the following text: “*Here’s the fastest speed*”. The word “lab” appeared 1.5 seconds later. The word stayed on the screen as the “+” immediately appeared and began to change color. Color changes occurred 8 times. At the end of the 8 cycles, a new trial began. For each new trial, participants were alerted to the change with the following text: “*A little slower*”. The same word then appeared 2 seconds later. There were 5 different rates, corresponding to the time it took for the plus sign to change from black to red: 350, 550, 750, 950, and 1150 ms. The slowest and fastest rates were chosen to be as extreme as possible while still being within the ability of participants to match.<sup>20</sup>

At the end of the practice session participants were told that they could begin the experiment whenever they were ready. The experimental trials were identical to the practice except that the rates went in order from slowest to fastest. Participants were presented with the following text: “*You will begin with the SLOWEST speed, and the flashing will become faster*”. Subsequently, each rate change was signaled with: “*The speaking rate will now speed up a bit*”. Trials were blocked by word, such that participants experienced all rates before beginning with a new word. At the end of a given block, participants were alerted that “*The next item will now appear on the screen*”, with a pause of 2 seconds before the word appeared. Word order was randomized across participants, but the order of rate presentation was fixed. Each word/rate pair was presented once.

The minimal pairs reported in this paper were chosen to vary across vowel quality (o or i), consonant manner (stop or fricative), and final consonant place (coronal or labial): feet/feed, thief/thieve, lobe/lope, and doze/dose. Differences in part of speech and morphological complexity were largely unavoidable in constructing CVC minimal pairs, but those factors are not expected to show interactions with speaking rate. No effort was made to balance word frequency, beyond the avoidance of archaic forms, for the same reason. While higher frequency words would

---

<sup>20</sup> Note that the fastest change time, 350 ms, is quite long in terms of vowel duration alone, as measured in the Buckeye Corpus. This presumably reflects the fact that coarticulation and reduction, along with prosodic organization, allow for individual segments to be much shorter in normal speech than in a laboratory word-repetition task.

be expected to be somewhat shorter across the board, there was no reason to believe that speaking rate would affect the individual segments differently.

### **3.2 Data Selection and Annotation**

Each participant produced approximately 8 tokens of each word at each rate. To avoid edge effects, and fluctuations in rate, a single representative token from the center of the group was selected and measured. Because each token was surrounded by other tokens at the same repetition rate, it was possible to segment both the closure and the release interval for each stop. However, at the fastest rates, final stops did not always have a clear release. In those cases, the end of the stop was set to the end of the voicing bar (for voiced stops), or the point at which the amplitude first dropped to background levels, indicating the end of the vowel and the beginning of the stop. Background level was estimated by the amount of noise visible during the gaps between successive words. The most ambiguous cases involved the segmentation of the sonorant /l/ from the following /o/ vowel, given a large degree of coarticulation. At faster rates, the point at which the release of the final /d/ became the initial fricative of the following token of “feed” could also be hard to determine. This was also true of the final /v/ and the initial /θ/ in “thieve” sequences. Measurement variability is likely to be highest in those contexts. Note, however, that any measurement errors for these kinds of tokens will only introduce error for one of the measured variables. For “lo”, the vowel duration will be affected by where the segment boundary is placed, but not the final p/b. For “feedfeed” and “thievethieve” the coda duration will be affected by where the segment boundary is placed, but not the vowel duration. Furthermore, any possible annotator bias in segmenting the sequence “lo”, for example, would have a minimal impact on the results, firstly because each participant was assigned to a single annotator, meaning that any effect could be absorbed into a random effect by speaker, and secondly, because both the voiced and voiceless minimal pairs would be segmented in the same way, such that the voicing effect (the difference in vowel durations) would not be affected by any bias. There is a possibility of resyllabification for the fastest word repetition rates, but this is only likely for p#l and b#l

sequences in the lope/lobe pair. If such resyllabification occurred, we might expect the stop to be shorter, with reduced aspiration. In fact, this might explain the abrupt drop in VOT between rates 4 and 5 for this word pair (See Fig. 3).

The data for the first two word pairs (feet/feed, thief/thieve) were randomly assigned to three undergraduate research assistants or annotation. One of the authors and two of the RAs then re-measured a subset of the data produced by the other two annotators. Discrepancies between any two raters were discussed as a group to establish shared criteria for ambiguous tokens. The two RAs then individually reviewed their previous measurements and made adjustments where their original segmentation did not meet the discussed criteria. The same two RAs each also re-measured half the data of the third RA who had left the lab at that point. The second set of words (lobe/lope, doze/dose) were measured later, by an additional two RAs. Measurement verification was conducted in the same way. It was stressed that the most important criterion was consistency. As a final check, 2% of all tokens from each annotator were re-measured by the first author, selected in pairs in order to assess the discrepancy in the measured voicing effect. In terms of absolute durations, vowel measurements differed by an average of 12 ms, and total consonant durations difference by an average of 21 ms. The difference in vowel duration between voiced and voiceless minimal pairs differed by 13 ms, and the difference in total consonant duration by 30 ms. However, because the durations were sometimes longer than the first author's measurements, and sometimes shorter, the actual effect of discrepancies in this sample of data were much smaller: .06 ms shorter for vowel duration; 12 ms shorter for total consonant duration, a vowel duration difference that was 3.6 ms smaller, and a consonant duration difference that was 15 ms larger.

Occasionally the voiced stops and fricatives at the slower repetition rates were produced with a final epenthetic schwa. There were 29 such tokens. Any words with final schwa were removed from the analysis. In many cases, participants produced dose and doze tokens that were difficult to disambiguate. Two such participants were removed due to their productions of final s and z being practically identical. Five participants were removed for either failing to vary their

speaking rate significantly across trials, or varying only inter-word pause duration rather than word duration. An additional participant was removed due to adopting a sing-song (high-low) prosody to the word repetition. The results from the 26 remaining participants, and the 5049 measured tokens, are given below. Praat (Boersma and Weenink, 2009) was used for segmentation and annotation.

### 3.3 Results

In Fig. 3 final stop durations are plotted as a function of repetition rate (shown as a number between 1 and 5, where 5 is the fastest rate, and 1 the slowest). Voiced and voiceless tokens are plotted separately, and three different duration measures are given: closure (black), VOT (light gray), and the sum of the two (TDur: dark gray). Closure duration for final voiced stops varied relatively little across repetition rates. However, most stops were also produced with a period of aspiration (VOT). Voiced stops show a clear increase in total duration as rate decreases, but one that appears to plateau at the slowest rates. For voiceless stops, closure duration increases steadily, patterning very closely with VOT. Because both duration measures show dependence on rate, total duration was used as the dependent variable for testing the Expandability Hypothesis.

\*\*\*FIGURE 3 ABOUT HERE\*\*\*

Figure 4 provides duration data for the full set of words, both vowel duration (triangles), and total obstruent duration (filled circles). Visual inspection shows that larger vowel durations were reached by the voiced member of each minimal pair, while larger obstruent durations were reached by the voiceless member. There is also a larger difference between consonant and vowel durations for voiced-final tokens across all repetition rates, and that difference increases with decreasing repetition rate.

\*\*\*FIGURE 4 ABOUT HERE\*\*\*

A linear mixed-effects model was fit to the vowel duration data as a function of repetition rate and consonant duration. Consonant duration was treated as a continuous variable, and



repetition rate, as an ordinal variable. Random intercepts for participant and word were included. Random slopes were excluded as they caused the model to fail to converge, or led to singularity. As expected, a significant (linear and quadratic) effect of speaking rate was found (vowels were longer at slower speaking rates). There was also a main effect of consonant duration; vowels were longer when the coda consonant was shorter. The interaction between rate and consonant duration also reached significance; the negative effect of consonant duration was strongest at the slowest rates. See Table 7. Only significant effects are shown. Adding voicing to this model did not improve fit.

\*\*\*TABLE 7 ABOUT HERE\*\*\*

A separate model of vowel duration as a function of rate and voicing behaves very similarly. As before, random intercepts were used for participant and word, but random slopes were excluded as they caused the model to fail to converge, or led to singularity. Main effects of (linear) rate and voicing are found (reference level is Voiceless), as well as an interaction between voicing and rate such that the positive effect of voicing is strongest at slower rates. However, adding consonant duration to this model *does* significantly improve model fit. In fact, adding consonant duration renders voicing only significant in interactions. Voicing interacts significantly with rate (both linear and quadratic terms) and consonant duration; and consonant duration interacts significantly with rate (both linear and quadratic terms). There is an additional three way interaction between consonant duration, voicing, and linear rate that also reaches significance. The interaction of both voicing and consonant duration with rate is expected based on the earlier observation that durational differences get larger with slower speaking rate (a negative effect of consonant duration is enhanced), and the voicing distinction becomes more strongly associated with length differences as speaking rate decreases (a positive effect of voicing is enhanced). The 3-way interaction shows that there is an additional positive effect for voicing at slower rates, and for longer consonants. This is likely because consonant duration, even though it is negatively correlated with vowel duration, is positively correlated with speaking rate, as is vowel duration.

Slower rates mean longer consonant durations, and larger duration differences, producing a larger voicing effect. See Table 8. Only significant effects (other than voicing) are shown.

\*\*\*\*TABLE 8 ABOUT HERE\*\*\*\*

The model of consonant duration as a function of voicing confirms the interpretation that the “voicing” effect is driven by consonant duration. See Table 9. Random intercepts for participant and word were included. Random slopes were excluded as they caused the model to fail to converge, or led to singularity. A fully crossed rate, voicing, and manner model produced significant main effects of speaking rate (linear), voicing, and manner. Fricatives were significantly longer than stops (reference level is Stops). A significant interaction between rate (linear) and voicing was also found, indicating, as expected, that differences in duration between voiced and voiceless consonants increased with decreasing repetition rate. An interaction between manner and rate (linear) also reached significance: the difference in duration between fricatives and stops was even larger at slower rates. Only significant interactions are shown.

\*\*\*\*TABLE 9 ABOUT HERE\*\*\*\*

A final analysis of the paired duration differences confirms the negative correlation between the *difference* in duration of voiceless and voiced consonants, and the *difference* in duration of their preceding vowels. Random intercepts for participant and word were included. Random slopes were excluded as they caused the model to fail to converge, or led to singularity. Adding manner to the model also resulted in singularity. The final model of  $\Delta V$  ( $=V_{VL} - V_{VD}$ ) included rate and consonant duration difference ( $\Delta C = C_{VL} - C_{VD}$ ) and their interactions. A significant main effect of rate (quadratic) and  $\Delta C$  were found, and significant interactions between  $\Delta C$  and rate, for both the linear and the quadratic terms. Thus the voicing effect ( $\Delta V$ ) is shown to be larger for larger negative values of  $\Delta C$ , which are enhanced at the slowest speeds. See Table 10. Only significant factors are shown.

\*\*\*\*TABLE 10 ABOUT HERE\*\*\*\*

### 3.4 Discussion

These results strongly support the Expandability Hypothesis. Firstly, we confirm the predicted difference in lengthening between voiced and voiceless consonants in coda position, paralleling what has been repeatedly found for consonants in initial and medial position (Port, 1976, 1981; Miller and Baer, 1983; Miller and Volaitis, 1989; Volaitis and Miller, 1992). There is a difference in consonant durations at all rates,<sup>21</sup> and there is also a large difference in the slopes of the duration curves. The difference in consonant duration increases with decreasing rate, as does the vowel duration difference. Pairing consonant duration differences with vowel duration differences at each speaking rate shows that the strength of the voicing effect is significantly correlated with the size of the consonant duration difference. The significant interaction between rate and consonant duration (vowels), and between rate and voicing (consonants), is precisely what is predicted if vowel duration differences derive from consonant duration differences, rather than depending on phonetic voicing, or an abstract phonological voicing feature. In fact, absolute vowel duration differences and consonant duration differences are very close. Rhyme duration differences were significantly different between voiceless and voiced, but not large, at 26 ms. For final stops, the rhyme duration difference was only 2.8 ms. These results probably over-estimate the degree to which vowel and consonant duration are traded off, given that the experimental task is highly unnatural, and likely to bias more towards uniform syllable duration than natural speech contexts.

## 4 The Expandability Hypothesis: Modeling the Corpus Data

The corpus and production study results, combined with the previous research summarized in Section 1, strongly suggest that final obstruent duration trades off against preceding vowel duration, and that the size of the resulting voicing effect depends on absolute

---

<sup>21</sup> Note that the shortest vowel durations in this study are between 150 and 200 ms, already in the upper range of values found in the conversational speech of the Buckeye Corpus.

duration. To account for both of these properties, in addition to the fact that apparent compensation is not “perfect” (cf., Chen, 1970; Keating, 1985; Port and Dalby, 1982), we propose a competition-based model where trade-offs in duration arise, not from isochrony, but from pressures to meet certain duration targets, none of which can be fully satisfied. It is instructive, however, to first consider a perfect compensation model, and the range of outcomes it supports. Although his goal was not specifically to model the voicing effect, Campbell (1992) provides such a model, making use of an elasticity parameter that we will adapt to our own model.

#### 4.1 A pure compensation model

In Campbell’s model, duration is specified at the syllable level, and distributed over the segments within the syllable according to their relative elasticity. The hypothesis is that patterns of variation observed at the segment level can be derived from just two parameters: the inherent elasticity of the segment (which is fixed), and what we will call the expansion coefficient ( $\epsilon$ ), which varies as a function of target syllable duration (see also Campbell and Isard, 1991 and Campbell, 1990). The function for calculating the expansion coefficient for a given syllable,  $\sigma_k$ , is given in Equation (4). The solution,  $\epsilon_k(\sigma_T)$ , is the value that, when distributed to each segment in the syllable according to their specific elasticities ( $\kappa_i$ ), will result in the necessary total duration change from the underlying syllable duration ( $\bar{\sigma}_k$ ), to the target syllable duration ( $\sigma_T$ ).

$$\epsilon_k(\sigma_T) = \frac{\sigma_T - \bar{\sigma}_k}{\sum_i \kappa_i} \quad (4)$$

In Campbell’s model, compensation effects derive from the dependence of the expansion coefficient on total elasticity. The more segments within a syllable, the smaller  $\epsilon$ , and the less any given segment is expanded. Higher-elasticity segments within the syllable produce the same effect. Conversely, segments with lower elasticity force more lengthening to take place over higher-elasticity segments within the same syllable. Because voiced obstruents have lower elasticity than voiceless, a larger expansion coefficient is required for the syllable closed by the voiced obstruent to reach the same target duration as the syllable closed by the voiceless

obstruent. The larger expansion coefficient, in turn, results in a longer vowel. Campbell notes that his model “...appears to account quite simply, *though in fact not completely*, for the lengthening that has been observed in vowels of English before voiced consonants.” (Campbell, 1992, p. 218, emphasis ours). An underlying difference in mean duration between voiced and voiceless obstruents comprises part of this voicing effect. However, the relative contribution of this fixed value decreases as syllable length increases. The difference in duration due to the differing elasticities of the two segments, on the other hand, increases as syllable length increases. Expansion is a linear function of  $\sigma_T$ , therefore the difference in expansion,  $(\varepsilon_{vd} - \varepsilon_{vl})$ , also increases linearly. See Appendix (B) for more details. Campbell uses standard deviation, and mean duration (both values estimated from the British English corpus SCRIBE), as proxies for segment elasticity, and underlying duration, respectively, based on the observation that longer segments generally show more variability in their duration (e.g., Lehiste, 1972).

Because Campbell does not limit the degree to which voiced obstruents can lengthen, the model under-estimates the voicing effect at the longest durations. Without additional mechanisms, the model also fails to account for the shortest end of the distribution, predicting, in fact, that the voicing effect should reverse under compression. Both of these mismatches are due to the use of a linear expansion function in regions of the duration space that do not behave linearly. Although Campbell (1992) does not explicitly require all syllables to be the same length, the fact that target syllable durations are strictly enforced means that any two syllables can be set to the same duration. In which case, the difference in vowel duration must be equivalent to the difference in obstruent duration between any two  $VC_1$ ,  $VC_2$  minimal pairs. Campbell additionally predicts that vowels in open syllables will be longer than vowels in closed syllables, and vowels in closed syllables with complex codas will be shorter than vowels in syllables with simplex codas, in both cases, by exactly the duration of the coda consonant.

## 4.2 A Competing Constraints Model of the Voicing Effect

In this section we model the voicing effect as the outcome of a competition between conflicting duration targets at the segment and syllable level. The primary differences from the model of Campbell (1992) are the following: expandability is not constant but varies as a function of duration, target syllable duration is treated as a constraint itself, and all constraints are gradiently “violable”, such that no given constraint need be perfectly satisfied (allowing for undershoot and overshoot). These modifications allow us to capture non-perfect “compensation”, the lack of a significant voicing effect for conversational speech, the general dependence of the voicing effect on duration, and the behavior of voiced obstruents under lengthening (see also Section 3). The results reported here are for VC syllables. See Appendix (C) for the treatment of CV syllables, and the “voicing” effect in onset position.

Constraints in the competition model are all implemented as Normal probability distributions. Each distribution assigns the highest probability to its preferred duration (the mean of the distribution), and smoothly decreasing probabilities for durations both longer and shorter than that mean. The variance of the distribution controls how quickly the probability decreases.<sup>22</sup> The smaller the variance, the more rapid the decrease, and the greater the resistance to deviations from the mean. Its variance thus acts effectively as a weighting factor for each constraint. This means that, all else being equal, a segment with a broader probability distribution will be lengthened or shortened more than a segment with a narrower probability distribution. Variance thus also maps to segment elasticity. Constraint “competition” in this model is realized through maximization of the joint probability function over all constraints. This function exhibits the desired behavior: one constraint may be “violated” to a greater degree (decrease in probability) if this allows another, more highly weighted constraint to be less “violated” (*greater* increase in probability).

The three segment-level constraints for the voicing effect model are shown graphically in

---

<sup>22</sup> Note that distribution variance is not a measure of actual duration variance. The latter is determined by the interaction of all constraints.

Fig. 5. Voiced and voiceless obstruent constraints are given the same mean value in these simulations, differing only in their variance. Target syllable duration is treated as a random variable, an external specification for a specific speaking rate or duration. Thus, each pair of values, (V, C), derived by the model are conditioned on a particular target syllable duration. In addition, two inter-segmental constraints specify preferred values for the  $\frac{C}{V}$  duration ratio and the  $\frac{V}{\sigma}$  duration ratio, respectively. The  $\frac{V}{\sigma}$  duration constraint forces vowel duration to lengthen with lengthening target syllable duration, while the  $\frac{C}{V}$  duration constraint requires consonant duration to do the same. Together they enforce monotonic behavior for both segments, such that they never shorten with an increase in target syllable duration, or lengthen with a decrease in target duration.

\*\*\*FIGURE 5 ABOUT HERE\*\*\*

If target syllable duration is strictly enforced, then this model will also exhibit perfect compensation, meaning that a given voiced syllable and a given voiceless syllable can always be set to the same duration. To avoid this, we introduce a final, violable, constraint for matching the target syllable duration. This constraint has a probability distribution centered at zero, over the variable  $\frac{\sigma_T - \sigma}{\sigma_T}$ : the normalized difference between actual and target durations. In the absence of an imposed target duration, all segments would be realized at their preferred absolute durations. In actuality, however, duration will always depend on how quickly one is speaking. In the model, this is the result of the competing pressure to reach the target syllable duration and the resistance of the individual segments to expansion or compression.

For a given target syllable duration, the model conducts a brute force search for the durations of the coda consonant (D or T) and vowel (V) that result in the highest joint probability over the entire set of constraints.<sup>23</sup> Although the model simply tries all possible combinations of values, the search space is restricted within a range where the maximum possible vowel duration

---

<sup>23</sup> Following Browman and Goldstein (1986) inter alia, we assume that there is a preferred timing relationship for a VC syllable which governs the degree of overlap between the articulatory gestures corresponding to the nucleus, and those corresponding to the coda. This parameter affects the apparent acoustic duration of the vowel, i.e., the portion that is not masked by the following consonant. Although we assume that modifications to this phasing relationship are possible, it does not vary in the current model. In all cases, there is no overlap between the two segments, such that the acoustic syllable duration is given by the sum of the vowel and consonant durations.

is set to the maximum syllable duration, and the minimum vowel duration is set to half the maximum syllable duration (such that consonant duration can never be larger than vowel duration). A fixed step size of 1 ms for both consonant and vowel is used to search this space. Each variable is assumed to be independent, therefore the joint probability is given as the product of the individual probabilities. See Appendix (C) for further details of the model.

Figure 6 shows the result of running the model for the set of target syllable durations ranging from 30 to 700 ms, sampled at 20 ms intervals (x-axis). On the y-axis, vowel duration, consonant duration, and syllable duration (V+C) are plotted for both voiced (black) and voiceless (gray) syllables. Each point on the graph corresponds to a VC word produced at the specified duration/rate. For example, for a target syllable duration of 330 ms, and a voiced-final syllable, the optimal vowel duration is 249 ms, and the optimal voiced obstruent duration is 62 ms. These points are shown as filled circles in Figure 5. Note that actual syllable duration is less than the target, at 311 ms. For a voiceless-final syllable, on the other hand, the optimal vowel duration is 238 ms, and the optimal voiceless obstruent duration is 80 ms (for a syllable duration of 318 ms). These points are shown as unfilled circles in Figure 5. The vowel, like each obstruent type, prefers the mean duration of its probability distribution (100 ms in this case). It is forced to lengthen due to the pressures of the other constraints. Because the voiced obstruent constraint is more highly weighted than the voiceless obstruent constraint (smaller variance), it does not shift as far from its preferred duration (at 50 ms). Therefore, the vowel is forced to lengthen more when it co-occurs with a voiced obstruent than with a voiceless obstruent.

At shorter target syllable durations, voiced and voiceless obstruents (black and gray solid lines, respectively) are more or less identical in duration; preceding vowel durations (black and gray dashed lines) are also identical within the same range. As target syllable duration continues to increase, the consonant durations start to diverge. Because of its much smaller variance, the voiced obstruent resists lengthening more strongly than the voiceless, and that resistance also grows faster, leading to smaller and smaller increases in duration. As a result, either the vowel must lengthen more, or the divergence from the target syllable duration must increase, or both. As



the obstruent duration difference continues to increase, the vowel duration difference will also continue to increase, meaning that the magnitude of the voicing effect will increase with increasing duration.

\*\*\*FIGURE 6 ABOUT HERE\*\*\*

As we have seen, because the constraint to match target duration competes with other constraints, a given syllable does not always match the target exactly. And because coda elasticity affects the outcome, voiced and voiceless syllables at the same target syllable duration do not necessarily have the same actual syllable durations. Thus this model allows for under- and over-compensation. This occurs only when imperfect matching would increase over-all probability. In the expansion regime, where vowel and consonant durations both increase past their preferred means (indicated by the filled circles in Fig. 6), both voiced and voiceless syllables are shorter than they would be if target syllable duration were strictly enforced, and the degree of under-compensation increases with increasing duration. This is the difference between the gray line ( $\sigma = \sigma_T$ ), and the two dotted lines that indicate actual syllable duration in Fig. 6. The increasing deviation from target is partially due to the fact that variance is expressed as a proportion: a larger deviation is tolerated for a longer syllable.

Voiced syllables are also systematically shorter than voiceless.<sup>24</sup> This is because even the highly-expandable vowel has a preference for its mean duration. A balance is struck between the length of the vowel and the amount of deviation from the target. In the compression regime, where segments shorten past their preferred durations, syllable durations are slightly longer than the target. For consonant durations below the mean, voiceless obstruents become slightly shorter than voiced. This occurs because elasticity is bi-directional; voiceless obstruents are both more expandable and more compressible than voiced stops.

Using this model, we simulated the corpus data by sampling 1000 points from a Normal distribution of target syllable durations,<sup>25</sup> durations that fall mostly in the range where there is a

---

<sup>24</sup> This is not necessarily the case under different parameter values.

<sup>25</sup> All durations less than 30 ms were set to 0.

negligible difference in consonant duration. This sample is depicted by the light blue vertical bars in Fig. 6. Each point from the sample represents a VC word produced at a given target duration/rate. For each of these 1000 points, the vowel and consonant durations that optimized joint probability were calculated, for voiced and voiceless syllables separately. Because words within each distribution were treated as identical (VD or VT), all words with the same target syllable duration had the same segment durations. The resulting data are plotted as probability density functions in Figure 7, for comparison with the corpus data in figures 1 and 2. Since only one type of vowel (fixed variance) was used, the simulation results do not entirely map to the corpus. However, they demonstrate that vowels of different inherent length are not necessary to derive the observed duration dependence of the voicing effect.

\*\*\*FIGURE 7 ABOUT HERE\*\*\*

As a proof of concept, the model does quite well at capturing the critical behaviors that motivated our re-analysis of the voicing effect in English, and without a directly compensatory mechanism. Languages other than English can be modeled by changing the relative variances of the obstruent probability distributions. A smaller difference in elasticity leads to a smaller voicing effect. The model can also capture the interaction between the voicing effect and vowel quality, using a lower elasticity parameter for inherently shorter vowels. Reducing the variance of the vowel probability distribution, but keeping all other parameters the same, results in a smaller voicing effect, and shorter syllables over-all. The voiced obstruents become slightly longer under these conditions, but the largest change is in how closely the target syllable duration is approximated. In this model, the reduced variance of the vowel results in shorter vowel durations, also causing greater undershoot at the syllable level. Qualitatively, this behavior is consistent with the finding that the voicing effect is significantly reduced in preceding vowels that are inherently short (Umeda, 1975; Crystal and House, 1982; De Jong, 2004). Note that the difference in duration between the obstruents themselves can, in principle, still grow quite large. See Appendix (C). Because very few studies on the voicing effect report final obstruent durations, it remains to

be seen whether this prediction is borne out.

Our model can also be used to model the behavior of CV syllables that differ with respect to the voicing of an onset consonant. In this model it is necessary to make the assumption that it is actually the rhyme, rather than the entire syllable that is the relevant prosodic unit. Additionally, a different phasing relationship is used to capture the onset-nucleus timing: one in which the consonant completely overlaps with the vowel (see, e.g., Browman and Goldstein, 1988). Because target rhyme duration does not affect the onset consonant, there is no trade-off of consonant duration with vowel duration. However, longer onset durations mask more of the vowel. As the durations of the voiced and voiceless obstruents diverge, vowels phased with voiceless obstruents become acoustically shorter than vowels phased with voiced obstruents. The behavior of this model is very similar to the short vowel model: a smaller voicing effect, but a similarly large divergence in obstruent duration. These two features are consistent with the known data. Voiced obstruents lengthen very little in speaking rate studies of CV syllables, while voiceless obstruents lengthen considerably, and vowel duration differences, when observed, are quite small (Allen and Miller, 1999; Pind, 1995; Miller and Baer, 1983; Miller et al., 1986). See Appendix (C) for more details of this model.

The competing constraints model does not differentiate between sources of lengthening, modeling only what occurs at the segment level to meet specified targets at some higher prosodic level, whether rhyme, syllable, word or foot. For very slow speaking rates, of the kind encountered in laboratory speech, a robust voicing effect can be observed. Similarly, pre-pausal lengthening can also produce a significant voicing effect. A particularly large final lengthening effect in English (e.g., Delattre, 1966), we conjecture, may be largely responsible for the particularly large voicing effect in this language.

## **5 Further Tests of The Expandability Hypothesis**

In the previous sections we have shown that vowel duration is better predicted by coda duration than by coda voicing. The implication being that the correlation between obstruent

duration and voicing is the source of the apparent voicing effect. It has also been demonstrated that a model of competing durational constraints can qualitatively capture the duration trade-offs between consonant and vowel duration. However, the Expandability Hypothesis, in and of itself, does not explain the ability of listeners to reliably use vowel duration to predict post-vocalic obstruent voicing. In the next section we will show that not only is the Expandability Hypothesis consistent with the perception literature, it is confirmed by certain results. For the remainder of the paper we will focus on word-final stops because there are often very limited cues to stops in final position, and it is primarily for stops that preceding vowel duration has been characterized as a contrastive cue.

## **5.1 Perception of voicing in final position**

A review of the perception literature in Section 1.5 has shown that other cues to the voicing contrast are likely to be stronger than preceding vowel duration, and categorical perception results may only be possible with highly impoverished stimuli. Meanwhile, categorical perception results have been obtained by varying obstruent duration alone. Based on these results, we hypothesize that listeners are using stop duration itself as the cue to voicing when final stops are both voiceless and unaspirated. Vowel duration factors into the classification decision insofar as it provides information about stop duration indirectly, as a measure of speaking rate.<sup>26</sup> In essence, the listener's task is to decide whether what they are hearing is a voiced stop spoken slowly or a voiceless stop spoken quickly. Shorter vowel durations, which comprise the majority of the corpus data, correspond to speaking rates at which voiced and voiceless stop durations are not significantly different from each other. In this range, vowel duration is ineffective as a cue to voicing. Only as speaking rate slows to the point where the voiced and voiceless expansion trajectories begin to diverge, does vowel duration become predictive.

The competition model of Section 4 is used to illustrate this hypothesis. See Fig. 8. The duration of the voiceless stop (gray solid line) gradually diverges from the duration of its voiced

---

<sup>26</sup> It is common practice to use stressed vowel duration as a proxy for local speaking rate (e.g., Crystal and House, 1982; Summerfield, 1981; Port and Dalby, 1982).

counterpart (black solid line), as the syllable is lengthened. This divergence is mirrored in the preceding vowel duration (gray dashed line – preceding voiceless stop; black dashed line – preceding voiced stop). If the listener is exposed to a relatively short vowel (Fig. 7a: upper horizontal dotted line), their expectation for the duration of the upcoming stop will be roughly the same regardless of whether it is voiced or voiceless (vertical difference between the lower open circles). An observed stop duration (lower dotted line) that falls close enough to both expected values is assumed to be acceptable for either member of the pair, and will not be sufficient to distinguish between the two in the absence of other cues.

\*\*\*FIGURE 8 ABOUT HERE\*\*\*

For a longer vowel, on the other hand, there is a larger difference in the expected durations of the voiced and voiceless stops. See Figure 7b. The same observed stop duration (lower dotted line) now falls significantly below both expected values. In a two-alternative forced choice task we predict that this stimulus should sound more like a voiced than a voiceless stop. In general, an ambiguous final stop of fixed duration becomes less ambiguous as vowel duration increases (speaking rate decreases). We assume that the category cross-over point from voiceless to voiced falls where the stimulus is significantly shorter than expected for a voiceless stop at that rate. After that, the likelihood of a voiced stop continues to increase (cf. Massaro and Cohen, 1983).

The foregoing can thus explain the increase in voiced responses with increasing vowel duration. However, given that we hypothesize that shorter vowels should not provide any cues to the voicing contrast, we would expect, all else being equal, that listeners would be at chance in identifying tokens in the short half of the continuum. Here it is the nature of the actual experimental stimuli that may bias perception strongly towards the voiceless stop. In the first place, ambiguous tokens are, by definition, phonetically voiceless. Depending on how exactly such stimuli were created, they may retain other cues to the original speech token from which they were generated, such as an F1 offset that is more consistent with a voiceless, than a voiced, stop. The synthetic stimuli used in Denes (1955), for example, were based on originally voiceless

tokens. Whereas Repp and Williams (1985), using naturally produced stimuli, found a large perceptual difference between continua generated from an originally voiced stop (lab), versus an originally voiceless stop (lap). Voiced responses were about 40% higher for the former across all but the two longest vowel durations.<sup>27</sup>

We therefore posit that the categorical perception results are due, firstly, to a default voiceless percept, based on residual cues that are more consistent with the voiceless member of the contrast, and secondly, to unusually long vowel durations. At the longest vowel durations (vanishingly rare in the speech corpus), we posit that the expected duration of a voiceless stop is so long that its likelihood approaches zero. For such extreme tokens, selection/perception of the voiced alternative may occur prior to actually hearing the final segment. However, it appears that the addition of a period of strong aspiration at the end of the stop is sufficient to switch the percept to voiceless.<sup>28</sup> Listeners may also be able to reliably select the voiced member of a minimal pair when final stops are entirely removed. We suspect that this is only possible in an explicit comparison task where listeners must label one token as voiced, and one as voiceless. In such a task it is likely that listeners assume a uniform speech rate, leading them to attribute a somewhat longer vowel duration to the effect of a following voiced stop.

Additional support for this account of voicing perception comes from studies of the voicing contrast in initial position. It has been consistently found that the perceptual VOT boundary is longer than the boundary estimated from production data (e.g., Miller et al., 1986; Miller and Volaitis, 1989; Volaitis and Miller, 1992). However, the two boundaries coincide when naturally produced, unedited stimuli are used in the perception task. Nagao and de Jong (2007) suggest that the mismatch may arise from the fact that the stimuli typically used in perception experiments are artificially impoverished. In other words, the edited tokens are so ambiguous that they can only be confidently classified at very long VOT, or very slow speaking rates. The

---

<sup>27</sup> Although Raphael (1972) tested both “voiced” and “voiceless” final synthetic stimuli, the only difference was that the voiceless lacked any F1 values at all during the transition period. All tokens in both experiments lacked a voicing bar, and contained no release bursts.

<sup>28</sup> This was established anecdotally when the spliced stimuli were played for various audiences.

consistency in the reported perceptual cross-over point across experiments on word-final stops may be explained by the same artificiality. For voiceless closures with no audible release, the duration of the coda stop is indeterminate. Listeners may therefore assume a duration that is plausible given their language experience and consistent with experimental variables such as the inter-stimulus interval. It is therefore likely to be relatively stable across experiments involving native speakers of English.

## 5.2 Predictions

Our explanation of the perception results generates at least one testable hypothesis. We predict that a change in the perception of voicing should lead to a change in the perception of speaking rate. During the course of vowel production, it is assumed that a hypothesis about both speaking rate and following segment duration is generated by the listener. In the absence of any information about the duration of the following stop (silent and unreleased), we posit that listeners will infer a duration that is consistent with those hypotheses. For a particularly long vowel, an expectation for a following phonologically voiced stop should lead listeners to infer the expected duration for a voiced obstruent, and the speaking rate associated with that duration (as depicted in Figure 7b: the intercepts of the leftmost vertical line with the voiced obstruent duration curve and the x-axis, respectively). However, if listeners subsequently experience unambiguous release or aspiration cues, then we hypothesize that there should be a noticeable correction to both the perceived stop class and the perceived speaking rate. The voiceless stop should indicate that the speaking rate is actually slower than previously supposed (represented by the x-intercept of the rightmost vertical line in Figure 7b).<sup>29</sup> Sanker (2019) has shown that the judgment of whether a vowel is “long” or “short” depends not only on the duration of the vowel, but on whether it is followed by a voiced or a voiceless obstruent. For vowels preceding voiced obstruents, longer durations are required to elicit a “long” response. Although she did not report

---

<sup>29</sup> The expected voiceless obstruent duration for that vowel duration is also expected to be longer. However, because speaking rate perception likely depends more on vowel duration than consonant duration, a change in the percept of voicing alone, without a change in the actual obstruent duration may be sufficient to trigger a change in the perception of speaking rate.

obstruent duration, we interpret her results as deriving from the expectation for a specific vowel duration given the unambiguous obstruent duration and its voicing. Vowels shorter than this expected value would be perceived as “short”, and vowels longer than this value would be perceived as “long”. The results of Fowler (1992) may indicate something similar. Also using a long/short judgment task, but VCVC disyllables, she finds that the longer the following closure duration, the larger the percentage of “long” vowel responses. As Sanker (2019) points out, this counter-intuitive result can be explained if lengthening the closure results in a change from a voiced to a voiceless percept. The longer the closure, the more likely the stop is to be classified as phonologically voiceless; and a vowel that was of the expected duration when the following stop was classified as voiced will be longer than expected for a following voiceless obstruent.<sup>30</sup>

The Expandability Hypothesis also predicts that it should be possible to find apparent compensation with segments other than immediately preceding or following vowels, as long as they are more expandable than voiced obstruents. This is corroborated to a certain extent. Weismer (1979) and Choi et al. (2016) have found that the VOT is longer for voiceless stops *word-initially* in CVC words when the final stop is voiced, than when it is voiceless.<sup>31</sup> A difference in nasal duration preceding voiced versus voiceless stops has also been found both for monosyllabic words of the form “dens/dense” (Raphael et al., 1975; Port and Cummins, 1992; Beddor, 2009), and polysyllabic words of the form “cantor/candor” (Vatikiotis-Bateson, 1984). Furthermore, Raphael et al. (1975) find that both vowel and nasal duration affect perception of voicing on final stops. In an eye-tracking study by Beddor et al. (2013), participants heard CVND words (such as “bend”), CVNT words (such as “bent”), and  $C\tilde{V}C$  words ([bẽd] vs [bẽt]), in which the nasal was missing but the vowel was nasalized. They found that, for  $C\tilde{V}C$  tokens, participants were overall more likely to fixate on the image corresponding to the CVNT word than the CVND

---

<sup>30</sup> Longer vowels also increase percentage “long” responses for closures in Fowler (1992). We speculate that this might be the result of re-syllabification. Vowels in open syllables tend to be longer than vowels in closed syllables. Therefore, the longer the vowel, the higher the likelihood that there is a syllable boundary between  $V_1$  and  $C_2$ . Since syllable-final lengthening is expected to apply to coda consonants, an expected duration for  $C_2$  as coda will be longer than expected for  $C_2$  as onset.

<sup>31</sup> Although we treat onsets as external to timing considerations, this is clearly a simplification. Onsets do contribute something to syllable duration, even if they play a smaller role in prosodic phenomena than codas.



word. They interpret this result as deriving from listener expectation that the nasal gesture will be coordinated differently in the two contexts: initiating earlier before a voiceless stop, and later before a voiced stop. However, no explanation is offered as to why the phasing relationship should be different in the two contexts. This difference, however, can be accounted for under the Expandability Hypothesis if the competition at the word (or syllable) level affect both the duration of individual gestures, as well as their phasing, as occurs under changes in speaking rate (e.g., Stetson, 1928; Hardcastle, 1985), and other types of prosodic lengthening (e.g., Byrd and Saltzman, 1998; Byrd et al., 2000). A shorter voiced stop would thus correlate with both longer tautosyllabic segments, as well as a preceding VN sequence that is less coarticulated. Less coarticulation, in turn, would result in less vowel nasalization. Thus, a highly nasalized vowel is more likely to occur preceding a [t] than a [d].

An additional corollary of our account of the voicing effect is that actual voicing, or any feature other than length, is not required for a “voicing” effect to arise. In fact, active phonetic voicing cannot be a requirement when the strongest effect is seen in English pre-pausally, where final voiced obstruents are likely to undergo devoicing. Sharf (1964) explicitly found that vowel duration differences were approximately the same whether words were produced normally or whispered, thus demonstrating that the effect was independent of actual vocal fold vibration.<sup>32</sup>

Given our hypothesis, however, it should be possible to find a “voicing” effect involving segments that have low elasticity for a reason not related to historic voicing. In principle, any apparent temporal compensation phenomenon could be modeled using the competing constraints framework (see Section 4.2). For example, characteristic differences in stop duration across place of articulation are fairly robust, and have also been associated with following vowel duration differences. However, the differences are very small, on the order of 10 ms for following vowels and 5 ms for the consonants themselves (Fant, 1973; House, 1961; Luce and Charles-Luce, 1985; Elert, 1965). Vowel duration differences preceding stops at various places of articulation, are

---

<sup>32</sup> This result is usually cited as evidence for the “linguistically determined” (i.e., phonologized) nature of the vowel duration cue.

similarly small, on the order of 5-15 ms. (Elert, 1965; Peterson and Lehiste, 1960). It is not clear whether these differences in consonant duration should be attributed to differences in characteristic duration, or differences in elasticity, or both.

Beguš (2017) finds that stop duration correlates negatively with preceding vowel duration not just for voiced and voiceless stops in Georgian, but for ejectives as well, with ejectives intermediate between voiced and voiceless stops in terms of both consonant duration and preceding vowel duration. Durvasula and Luo (2012) report vowel durations in Hindi (with a 4-way contrast) as longest preceding voiced aspirated stops, then voiced unaspirated, then voiceless aspirated, and finally voiceless unaspirated. Total stop durations were not reported, but if closure duration is the relevant measure, then these results are largely consistent with the Expandability Hypothesis: closure durations were shortest for voiced aspirated, and longest for voiceless unaspirated. Voiceless aspirated were slightly longer than than voiced unaspirated.<sup>33</sup>

All else being equal, we might predict that an appreciable difference in consonant duration should lead to a complementary difference in preceding vowel duration in monosyllabic words. However, it may prove difficult to isolate elasticity-based effects from other factors that affect syllable duration. For example, vowels in monosyllables closed by nasals have been found to be as long, or longer, than vowels in monosyllables closed by voiced obstruents in English (e.g., Peterson and Lehiste, 1960; Umeda, 1975; Crystal and House, 1988; House and Fairbanks, 1953), which is the opposite of what one would expect for a sonorous segment like a nasal. However, both Crystal and House (1988) and Klatt (1975) report nasal durations that are comparable to those for voiced stops. Thus, it may be the case that nasals (and possibly other sonorants) are not as elastic as might have been expected. Another possibility is that the phasing relationship between vowel and coda may be different in the case where the two gestures can overlap significantly without masking. Thus vowels may be measured as longer, and nasals, as shorter, if

---

<sup>33</sup> Durvasula and Luo (2012) reject this correlation based on continuous measures of vowel and consonant duration. Their analysis shows that the two measures are positively, not negatively, correlated. However, given that speaking rate controls syllable duration, both measures should be correlated with speaking rate in the same way, and thus positively correlated with one another. Thus, variation in speaking rate should lead to a positive correlation. This has also been suggested by Beguš (2017).

there is significantly more coarticulation than occurs with other consonants. If this is correct, then the vowel should be acoustically highly nasalized when the nasal is short, reflecting the true length of the nasal. Note that this would be consistent with the results for  $C\tilde{V}C$  words in Beddor et al. (2013).

It has also been found that vowels preceding voiced fricatives are longer than vowels preceding voiced stops, while vowels preceding voiceless fricatives are somewhat longer than those preceding voiceless stops (Umeda, 1975; Peterson and Lehiste, 1960).<sup>34</sup> Furthermore, the voicing effect has been reported to be larger for fricatives than for stops (e.g., House and Fairbanks, 1953; House, 1961). Although our production experiment was not designed to explicitly test fricatives against stops, our results are in line with these findings. In our data, vowel durations were longest before voiced fricatives, and a larger voicing effect was found for fricatives than stops (91 ms  $\Delta V$ , versus 56 ms). However, the Expandability Hypothesis predicts that preceding vowel durations should be similar for voiced stops and fricatives, given that voiced fricatives were only 15 ms longer than voiced stops on average. It also predicts that vowels should be shorter before voiceless fricatives than voiceless stops, given that voiceless fricatives were about 29 ms longer than voiceless stops. It may be that different consonants have different preferred C/V durations, thus exhibiting steeper or shallower curves for duration differences as a function of speaking rate. There is a suggestion of this in Figure 9, where there appear to be differences in the slopes of both  $\Delta C$  and  $\Delta V$  between the stop minimal pairs and the fricative minimal pairs. Nevertheless, this discrepancy is a potential problem for the Expandability Hypothesis, and requires further investigation.

\*\*\*FIGURE 9 ABOUT HERE\*\*\*

---

<sup>34</sup> Umeda (1975) also finds that vowel duration preceding nasals is sometimes longer than before voiced stops, sometimes shorter, depending on the vowel. While vowels before voiceless fricatives tend to be intermediate in duration between voiceless stops and nasals, low vowels are actually longer before voiceless fricatives than before nasals. In a production experiment with Russian speakers, Kavitskaya (2002) finds that the *difference* in vowel duration between open and closed syllables is smallest before voiced fricatives, consistent with the other two studies. However, she also finds that voiceless stops have the next smallest difference, followed by voiceless fricatives, voiced stops, and nasals, with liquids showing the largest difference. The apparently variable behavior of nasals, voiceless stops and voiceless fricatives suggests that a number of interacting factors affect nucleus duration.

The Expandability Hypothesis as developed here was designed for consistency with an already very large experimental literature, thus many of its predictions are actually postdictions. Nevertheless, we have offered a number of speculations that can, in principle, be tested. Among these are the hypothesis that onsets are largely excluded from rate/duration targets, that longer vowels in CVN words are highly nasalized, and that less nasalization in VNC sequences is correlated with longer VN durations. The competing constraints model also offers the hypothesis that significant differences in obstruent duration can occur without apparent compensation on vowels that are inherently short, or vowels that follow consonants at the beginning of the syllable (see Appendix C). Additional predictions about differences in effect size across final, medial, and initial position cannot be entirely determined by comparing across heterogeneous studies, but require carefully controlled experimentation to assess. More detailed information about gestural coordination between vowels and specific following consonants is also needed to fine-tune model predictions.

## **6 Summary & Conclusions**

In much modern work, the voicing effect tends to be described in simplified terms, as a regular, quasi-universal, phonetically-driven phenomenon. In English, preceding vowel duration is often said to play a contrastive role for word-final stops. A small subset set of the literature – work that demonstrates either strong categorical perception (e.g., Raphael, 1972) or large vowel duration differences in production (e.g., Mack, 1982) – is most frequently cited. Such studies are primarily conducted using monosyllabic single-word stimuli in a laboratory setting, where speaking rate is much slower than for normally produced speech (production), and cues to stop identity are significantly impoverished, if not missing altogether (perception).

Yet it has been known for some time that vowel duration differences can be quite small in continuous speech, in polysyllabic words, across a syllable boundary, and phrase-medially (e.g., Umeda, 1975). Additionally, lax, unstressed, or otherwise inherently short vowels show little to

no voicing effect even in laboratory speech (e.g., Peterson and Lehiste, 1960). We confirmed both these results using the Buckeye Speech Corpus, finding no over-all effect of final-obstruent voicing on vowel duration when interaction terms were included, and those interaction terms suggested that voicing predicted vowel duration only when it was highly correlated with obstruent duration. We confirmed that there was a significant effect of voicing on consonant duration. And we noted that, like vowels, the difference in average durations between voiced and voiceless obstruents was larger at the higher end of the duration distribution.

In production studies that manipulate speaking rate it has been shown that voiceless obstruents, in both word-initial pre-stressed (VOT, e.g., Miller and Volaitis, 1989), and word-medial post-stress (closure duration, e.g., Port, 1976) position, are longer than voiced, with that difference increasing as speaking rate decreases. We extended that finding to coda position, demonstrating that the difference in vowel duration increased in step with the inverse duration difference for obstruents.<sup>35</sup> Using paired data, we were able to show that the magnitude of the voicing effect depended on obstruent duration across the board, while voicing was only significant at the slower rates (i.e., when it was significantly correlated with duration).

Aspiration, and actual voicing, have also been shown to be stronger cues to “voicing” than preceding vowel duration (Wardrip-Fruin, 1982; Repp and Williams, 1985). Furthermore, depending on the type of stimuli, vowel duration may have no effect on phoneme identification at all (Revoile et al., 1982). Obstruent duration itself has been shown to affect voicing perception in final position (Denes, 1955; Raphael, 1981; Repp and Williams, 1985), just as it does in word-medial position (Port and Dalby, 1982). This body of results argues against preceding vowel duration as a primary cue to the voiced/voiceless contrast in English. Indeed, it strongly suggests that vowel duration directly affects the perception of obstruent duration, not voicing itself, and is only predictive of voicing class within a certain upper range of durations. We have

---

<sup>35</sup> In a similar study, Ko (2018) found that duration differences between voiced and voiceless obstruents, and between their preceding vowels, both increased with decreasing speaking rate. However, of the three speaking rates, the “normal” and “fast” conditions were largely the same, and duration differences were not analyzed as paired (voiced, voiceless) data.

offered a proposal that accommodates this full set of results, as well as additional related findings. Namely, that the voicing effect in English is the result of the inherently low elasticity of voiced obstruents, and that segment durations, in general, are determined by the components of the Expandability Hypothesis, reproduced below.

(5) The Expandability Hypothesis

All segments have a characteristic elasticity that determines their resistance to lengthening

Resistance to lengthening increases with increasing duration for all segments

Lower elasticity equates with a more rapid increase in resistance

Relative resistance determines the distribution of duration across the syllable

The inverse correlation between obstruent duration and vowel duration, and its dependence on speaking rate, are attributed to a type of compensatory effect (see also Massaro and Cohen, 1983; Campbell, 1992), but not one based on syllable isochrony. Our competing constraints model of segment timing allows for “imperfect compensation”, which appears to be the rule in language generally, rather than the exception (e.g., Browman and Goldstein, 1988; Krivokapić, 2020).

This model provides a proof of concept for deriving the voicing effect from a set of general-purpose timing constraints. Our account also covers much more empirical ground than explanations of the voicing effect that are based on actual vocal fold vibration, or articulatory effort. We are able to unify the treatment of the contrast across word and syllable position, and draw connections between effects based on differences of consonant elasticity, and those based on differences of vowel elasticity. Our explanation for the voicing effect also has ramifications for theories of contrastive features.

## 6.1 Contrast and Allophony

Throughout this paper the relevant obstruent contrast in American English has been referred to as one of voicing. This is in spite of the fact that it is precisely because phonetic voicing is often absent from “voiced” stops that preceding vowel duration can be discussed as a

possible cue to contrast. Clearly, the presence or absence of vocal fold vibration is not always necessary, or even sufficient, for phoneme identification. In order for the contrast to be described as one of voicing, it is necessary to treat the phonological voicing feature as distinct from the phonetic feature of the same name. The first is transformed to the second via a series of allophonic rules. For example, in absolute initial position the */-voice/* stop becomes *[+spread glottis]*, while the */+voice/* stop may become *[-voice]*. In final position, a */-long/* vowel preceding a */+voice/* stop becomes *[+long]*.

However, we have seen that apparent vowel lengthening varies considerably as a function of speaking rate, sentence and word position, stress, and other factors. Most importantly, longer vowels correlate with shorter consonants, and voiced obstruents tend, cross-linguistically, to be shorter than their voiceless counterparts. The apparent physiological difficulty of maintaining the necessary conditions for voicing over extended closure periods has been proposed as an explanation for this tendency (e.g. Ohala, 1983, 2011). Nevertheless, it is possible, by virtue of greater articulatory effort, to maintain voicing if desirable, at least up to a point. Partial, or total, devoicing is also a possible outcome. Thus, whether or not voiced obstruents are actually shorter than voiceless obstruents is language specific. The fact that “voiced” stops in English are now frequently devoiced means that the observed duration differences are no longer the direct result of physiological constraints, but of what has become an underlyingly specified property of the segment. The fact that the difference in behavior between voiced and voiceless obstruents is only observable at long durations means that the specification is not for absolute duration, but for something that quantifies resistance to lengthening; the lower elasticity of the voiced segment is only observed when the two segments are pushed well beyond their preferred durations. Our claim is that vowel duration differences emerge directly from these elasticity differences. Therefore, we also conclude that vowel duration is *not* a feature that is specified, either at the phonological or phonetic level.

Although categorical perception effects have been demonstrated for the vowel duration cue, this is not particularly noteworthy, given that the number of acoustic cues to the contrast that

listeners are able to exploit has been shown to be quite large. Duration and intensity of voicing, aspiration, and F0 contour, length of vowel formant transitions with respect to steady state duration (Fitch, 1981), F1 offset frequency (Crowther and Mann, 1992), speed of jaw lowering, and jaw offset position (Van Summers, 1987) all differ consistently between the two stop types in final position. In medial post-stress position, consistent differences have also been found in the timing of vocalic voice offset, and the signal decay time (Lisker, 1986), which should apply to final position as well. Furthermore, it is well known that cues can be “traded off” with one another. That is, while a long enough closure duration can cue a “voiceless” stop on its own, a shorter closure in tandem with a shortened vowel can also do so (e.g., Kohler, 1979, 1984; Fitch, 1981; Lisker, 1986; Van Summers, 1987; Bailey and Summerfield, 1980; Klatt, 1976; Malécot, 1968). Yet absolute vowel duration, and not closure duration or formant transition information, is frequently characterized as a phonological “voicing” feature, even though the latter two cues have been shown to influence perception to the same, or an even greater, degree. This may be due, in large part, to the privileging of ‘prominent’ contexts in phonological theory.

## **6.2 Prominence**

While the phonetic realization of underlyingly contrastive features is assumed to vary by context, the most prominent environment, usually initial pre-stress position, is assumed to most faithfully reflect those features. Not only that, but features are said to be enhanced, or more strongly signaled, in such contexts (e.g., Kingston and Diehl, 1994). Conversely, observed enhancement is taken to indicate features that are “controlled”, or underlyingly specified, as opposed to being supplied by context-sensitive rules (e.g., Ohala, 1981). Enhancement can be realized as an increase in acoustic amplitude, an increase in size of articulatory gestures, and/or an increase in gestural, and thus, segmental, duration. In addition to making individual features more salient, enhancement is also assumed to be a mechanism for increasing discriminability between the members of a phonemic contrast (e.g., De Jong, 1995; Cho, 2016; Cho and Jun, 2000). For the above reasons, slower than normal speaking rate is considered to be an enhancement mechanism



that should lead to lengthening, *but only of contrastively specified features* (e.g., Solé, 2007).

Underspecification theory applied to laryngeal contrasts typically makes use of the following privative features: [spread glottis], [voice], and [constricted glottis] (e.g., Kim, 1970; Iverson and Salmons, 1995). This system yields three possible two-way contrast systems, one for each of the features, with the second member always unspecified. The phonetically voiceless stops in French and Thai fail to lengthen significantly with decreased speaking rate, and are therefore taken to be unspecified for laryngeal features, while the phonetically short lag/voiced stops in English are the unspecified member of the contrast<sup>36</sup> (Kessinger and Blumstein, 1997; Beckman et al., 2013).

In the same vein, an observed interaction between a given phonetic cue, and any variable that affects duration, is taken to indicate that the cue is an inherent part of the contrast. It has been argued that vowel duration is purposefully manipulated by speakers to enhance the laryngeal contrast of the following obstruent, based on the following set of results: that the effect of stress is smaller for pre-voiceless than pre-voiced vowels in English (De Jong, 2004); that /-long/ pre-voiced vowels lengthen less than they would otherwise, in order to avoid overlapping with /+long/ pre-voiceless vowels, and preserve an existing long versus short vowel distinction in German (Braunschweiler, 1997); that vowel duration differences preceding voiced versus voiceless segments are greater for long vowels than for short vowels in English (Peterson and Lehiste, 1960); that the difference in duration between the stressed vowel in a monosyllabic word and the same vowel in a bisyllabic word is larger (by percentage) for syllables closed by voiced stops than those closed by voiceless stops (Van Summers, 1987; De Jong, 1991; Crowther and Mann, 1992; Raphael, 1975; Smith, 2002; Klatt, 1973); that the vowel shortening effect of affixation is greater (both absolutely, and proportionally) for a voiced-final stem than for a voiceless final (Lehiste, 1972).

In this paper, however, we have conceptualized stress, prosodic boundary marking, and speaking rate simply as external forces which, among others, can act to lengthen segments. Under

---

<sup>36</sup> English is usually described as a [spread glottis]/Ø system, although this is not uncontroversial.

our account, all segments are subject to such lengthening and shortening pressures. How much lengthening or shortening actually occurs for individual segments, however, is governed by the interactions of all such constraints, some of which are more highly weighted than others. The apparently asymmetric effects on voiced versus voiceless syllables do not need to be explained as the result of speaker effort to avoid phonetic ambiguity, or to maintain a specific range of phonetic values. They follow directly from these two premises: that the voicing effect derives from differences in segment elasticity; and that the resulting differences in duration increase with increasing duration. Characterizing the voicing effect as a consequence of on-line timing adjustments (to which multiple factors can contribute) is therefore more parsimonious, and more explanatorily adequate, than the hypothesis that there is both a grammatical rule of vowel lengthening, and a set of deliberate adjustments made to preserve the output of that rule. Note that this analysis requires elasticity to be underlyingly specified. This is not the same, however, as an underlying specification for an abstract voice feature. In the first place, all segments are assumed to have their own characteristic elasticity. Furthermore, a specification of this kind is necessary for independent reasons: to account for the differing degrees to which segments respond to changes in speaking rate. Finally, relative duration values within a word cannot be derived from a /voice/ feature, or even a long/short duration feature, as they depend on potentially complex interactions between all the segments within a word.

If ‘prominent’ environments do not actually enhance contrastive features, then the realization of the features in such contexts should not necessarily be taken as underlying. Doing so, in fact, requires potentially extensive transformations to arrive at the more frequent, non-prominent contexts of normal speech. If we reverse this relation, then very slow hyper-articulated speech is the exception, rather than the rule, and intense aspiration and especially long durations are derived from features that are more typical of the contrast in general. Large differences in preceding vowel duration are, almost exclusively, the product of atypical speech and therefore, in our view, should be considered the least central to the “voicing” contrast, not the most. This flipped view of contrast offers an intriguing avenue for future research.

## 7 References

- Abdelli-Beruh, N. 2004. The stop voicing contrast in French sentences: Contextual sensitivity of vowel duration, closure duration, voice onset time, stop release and closure voicing. *Phonetica* 61:201–219.
- Allen, J. Sean, and Joanne L. Miller. 1999. Effects of syllable-initial voicing and speaking rate on the temporal characteristics of monosyllabic words. *The Journal of the Acoustical Society of America* 106:2031–2039.
- Aylett, Matthew, and Alice Turk. 2004. The Smooth Signal Redundancy Hypothesis: A Functional Explanation for Relationships between Redundancy, Prosodic Prominence, and Duration in Spontaneous Speech. *Language and Speech* 47:31–56.
- Bailey, Peter J., and Quentin Summerfield. 1980. Information in speech: Observations on the perception of [s]-stop clusters. *Journal of Experimental Psychology: Human Perception and Performance* 6:536 – 563.
- Beckman, Jill, Michael Jessen, and Catherine Ringen. 2013. Empirical evidence for laryngeal features: Aspirating vs. true voice languages. *Journal of Linguistics* 49:259–284.
- Beddor, Patrice Speeter. 2009. A Coarticulatory Path to Sound Change. *Language* 85:785–821.
- Beddor, Patrice Speeter, Kevin B. McGowan, Julie E. Boland, Andries W. Coetzee, and Anthony Brasher. 2013. The time course of perception of coarticulation. *The Journal of the Acoustical Society of America* 133:2350–2366.
- Beguš, Gašper. 2017. Effects of ejective stops on preceding vowel duration. *The Journal of the Acoustical Society of America* 142:2168–2184.
- Belasco, Simon. 1958. Variations in Vowel Duration: Phonemically or Phonetically Conditioned? *The Journal of the Acoustical Society of America* 30:1049–1050.

- Bell, Alan, Jason M. Brenier, Michelle Gregory, Cynthia Girand, and Dan Jurafsky. 2009. Predictability effects on durations of content and function words in conversational English. *Journal of Memory and Language* 60:92–111.
- Berkovits, Rochele. 1993. Utterance-final lengthening and the duration of final-stop closures. *Journal of Phonetics* 21(4):479 – 489.
- Boersma, Paul, and David Weenink. 2009. Praat: Doing phonetics by computer (Version 6.0.36). computer program. URL <http://www.fon.hum.uva.nl/praat/>.
- Braunschweiler, Norbert. 1997. Integrated cues of voicing and vowel length in German: A production study. *Language and Speech* 40:353–376.
- Browman, Catherine P., and Louis M. Goldstein. 1986. Towards an articulatory phonology. *Phonology* 3:219–252.
- Browman, Catherine P., and Louis M. Goldstein. 1988. Some notes on syllable structure in articulatory phonology. *Phonetica* 45:140–155.
- Browman, Catherine P., and Louis M. Goldstein. 1990. Tiers in articulatory phonology, with some implications for casual speech. In *Papers in laboratory phonology I: Between the grammar and the physics of speech*, ed. John Kingston and Mary E. Beckman, 341–376. Cambridge: Cambridge University Press.
- Byrd, Dani, Abigail R. Kaun, Shrikanth Narayanan, and Elliot Saltzman. 2000. Phrasal signatures in articulation. In *Papers in laboratory phonology v*, ed. M. B. Broe and J. B. Pierrehumbert, 70–87. Cambridge: Cambridge University Press.
- Byrd, Dani, Sungbok Lee, Daylen Riggs, and Jason Adams. 2005. Interacting effects of syllable and phrase position on consonant articulation. *The Journal of the Acoustical Society of America* 118:3860–3873.

- Byrd, Dani, and Elliot Saltzman. 1998. Intra-gestural dynamics of multiple prosodic boundaries. *Journal of Phonetics* 26:173–199.
- Byrd, Dani, and Elliot Saltzman. 2003. The elastic phrase: Modeling the dynamics of boundary-adjacent lengthening. *Journal of Phonetics* 31:149–180.
- Cambier-Langeveld, Gerda Martina. 1997. The domain of final lengthening in the production of Dutch. *Linguistics in the Netherlands* 14:13–24.
- Cambier-Langeveld, Gerda Martina. 2000. *Temporal marking of accents and boundaries*. Den Haag: Holland Academic Graphics.
- Campbell, W. Nick. 1990. Timing invariance in read speech. In *Speaker Characterization in Speech Technology*, 78–83.
- Campbell, W. Nick. 1992. Syllable-based segmental duration. *Talking machines: Theories, models, and designs* 211–224.
- Campbell, W. Nick, and Stephen D. Isard. 1991. Segment durations in a syllable frame. *Journal of Phonetics* 19:37–47.
- Catford, John Cunnison. 1977. *Fundamental problems in phonetics*. Midland Books.
- Chen, Matthew. 1970. Vowel length variation as a function of the voicing of the consonant environment. *Phonetica* 22:129–159.
- Cho, Taehong. 2016. Prosodic boundary strengthening in the phonetics–prosody interface. *Language and Linguistics Compass* 10:120–141.
- Cho, Taehong, and Sun-Ah Jun. 2000. Domain-initial strengthening as enhancement of laryngeal features: Aerodynamic evidence from Korean. *UCLA Working Papers in Phonetics* 57–70.
- Cho, Taehong, and Patricia Keating. 2009. Effects of initial position versus prominence in English. *Journal of Phonetics* 37:466–485.

- Choi, Jiyoung, Sahayng Kim, and Taehong Cho. 2016. Phonetic encoding of coda voicing contrast under different focus conditions in L1 vs. L2 English. *Frontiers in psychology* 7.
- Clements, G. N., and S. J. Keyser. 1983. *CV phonology: A generative theory of the syllable*. MIT Press.
- Coretta, Stefano. 2019. An exploratory study of voicing-related differences in vowel duration as compensatory temporal adjustment in Italian and Polish. *Glossa: a journal of general linguistics* 4.
- Crowther, Court S., and Virginia Mann. 1992. Native language factors affecting use of vocalic cues to final consonant voicing in English. *The Journal of the Acoustical Society of America* 92:711–722.
- Crystal, Thomas H, and Arthur S House. 1982. Segmental durations in connected speech signals: Preliminary results. *The Journal of the Acoustical Society of America* 72:705–716.
- Crystal, Thomas H, and Arthur S House. 1988. Segmental durations in connected-speech signals: Current results. *The Journal of the Acoustical Society of America* 83:1553–1573.
- Cuartero Torres, Néstor. 2002. Voicing assimilation in Catalan and English. Doctoral Dissertation, Universitat Autònoma de Barcelona.
- Davis, Stuart, and W. Van Summers. 1989. Vowel length and closure duration in word-medial VC sequences. *Journal of Phonetics* 17:339–353.
- De Jong, Kenneth J. 1991. An articulatory study of consonant-induced vowel duration changes in English. *Phonetica* 48:1–17.
- De Jong, Kenneth J. 1995. The supraglottal articulation of prominence in English: Linguistic stress as localized hyperarticulation. *The Journal of the Acoustical Society of America* 97:491–504.

- De Jong, Kenneth J. 2004. Stress, lexical focus, and segmental focus in English: Patterns of variation in vowel duration. *Journal of Phonetics* 32:493–516.
- De Jong, Kenneth J., and Bushra Zawaydeh. 2002. Comparing stress, lexical focus, and segmental focus: Patterns of variation in Arabic vowel duration. *Journal of Phonetics* 30:53–75.
- Delattre, Pierre. 1962. Some factors of vowel duration and their cross-linguistic validity. *The Journal of the Acoustical Society of America* 34:1141–1143.
- Delattre, Pierre. 1966. A comparison of syllable length conditioning among languages. *IRAL-International Review of Applied Linguistics in Language Teaching* 4:183–198.
- Denes, Peter. 1955. Effect of duration on the perception of voicing. *The Journal of the Acoustical Society of America* 27:761–764.
- Derr, Marcia A, and Dominic W. Massaro. 1980. The contribution of vowel duration, F0 contour, and frication duration as cues to the /juz/-/jus/ distinction. *Perception & Psychophysics* 27:51–59.
- Durvasula, Karthik, and Qian Luo. 2012. Voicing, aspiration, and vowel duration in Hindi. In *Proceedings of Meetings on Acoustics 164ASA*, volume 18, 060009.
- Eddington, David, Rebecca Treiman, and Dirk Elzinga. 2013. Syllabification of American English: Evidence from a large-scale experiment. Part I. *Journal of Quantitative Linguistics* 20:45–67.
- Elert, Claes Christian. 1965. *Phonologic studies of quantity in Swedish: Based on material from Stockholm speakers*. Almqvist och Wiksell.
- Fant, Gunnar. 1973. Stops in CV-syllables. Technical report, Dept. of Speech, Music and Hearing: Quarterly Progress and Status Report.
- Fidelholtz, James L. 1975. Word frequency and vowel reduction in English. In *Papers from the Eleventh Regional Meeting of the Chicago Linguistic Society*, volume 11, 200–213.

- Fischer, Rebecca M., and Ralph N. Ohde. 1990. Spectral and duration properties of front vowels as cues to final stop-consonant voicing. *The Journal of the Acoustical Society of America* 88:1250–1259.
- Fitch, Hollis Leslie. 1981. Distinguishing temporal information for speaking rate from temporal information for intervocalic stop consonant voicing. Technical report, Haskins Laboratory.
- Flege, James. 1979. Phonetic interference in second language acquisition. Doctoral Dissertation, Indiana University.
- Fosler-Lussier, Eric, and Nelson Morgan. 1999. Effects of speaking rate and word frequency on pronunciations in conversational speech. *Speech Communication* 29:137–158.
- Fougeron, C., and P. A. Keating. 1997. Articulatory strengthening at edges of prosodic domains. *The Journal of the Acoustical Society of America* 101:3728–3740.
- Fowler, Carol, Kevin Munhall, Elliot Saltzman, and Sarah Hawkins. 1986. Acoustic and articulatory evidence for consonant-vowel interactions. *The Journal of the Acoustical Society of America* 80:S96–S96.
- Fowler, Carol A. 1981. A relationship between coarticulation and compensatory shortening. *Phonetica* 38:35–50.
- Fowler, Carol A. 1992. Vowel duration and closure duration in voiced and unvoiced stops: There are no contrast effects here. *Journal of Phonetics* 20:143–165.
- Fox, Robert A., and Dale Terbeek. 1977. Dental flaps, vowel duration and rule ordering in American English. *Journal of Phonetics* 5:27–34.
- Gahl, Susanne. 2009. Homophone duration in spontaneous speech: A mixed-effects model. Technical Report 5.
- Gahl, Susanne, and Harald Baayen. 2022. Time and thyme again: Connecting spoken word duration to models of the mental lexicon .



Gahl, Susanne, Yao Yao, and Keith Johnson. 2012. Why reduce? Phonological neighborhood density and phonetic reduction in spontaneous speech. *Journal of Memory and Language* 66:789–806.

Gimson, Alfred Charles. 1970. *An introduction to the pronunciation of English*. London: Hodder Arnold.

Halle, Morris, and Kenneth Stevens. 1967. Mechanism of glottal vibration for vowels and consonants. *The Journal of the Acoustical Society of America* 41:1613–1613.

Hardcastle, William J. 1985. Some phonetic and syntactic constraints on lingual coarticulation in stop consonant sequences. *Speech Communication* 4:247–263.

Harris, MS, and Noriko Umeda. 1974. Effect of speaking mode on temporal factors in speech: Vowel duration. *The Journal of the Acoustical Society of America* 56:1016–1018.

Hillenbrand, James M., Michael J. Clark, and Robert A. Houde. 2000. Some effects of duration on vowel recognition. *The Journal of the Acoustical Society of America* 108:3013–3022.

Hillenbrand, James M., Dennis R. Ingrisano, Bruce L. Smith, and James E. Flege. 1984. Perception of the voiced–voiceless contrast in syllable-final stops. *The Journal of the Acoustical Society of America* 76:18–26.

Hofhuis, E., Carlos Gussenhoven, and Toni Rietveld. 1995. Final lengthening at prosodic boundaries in Dutch. volume 1, 154–157. Stockholm: Stockholm University.

Hogan, John T., and Anton J. Rozsypal. 1980. Evaluation of vowel duration as a cue for the voicing distinction in the following word-final consonant. *The Journal of the Acoustical Society of America* 67:1764–1771.

Hooper, Joan B. 1976. Word frequency in lexical diffusion and the source of morphophonological change. In *Current progress in historical linguistics*, ed. William Christie, 96–105. Amsterdam: North Holland.

- House, Arthur S. 1961. On vowel duration in English. *The Journal of the Acoustical Society of America* 33:1174–1178.
- House, Arthur S, and Grant Fairbanks. 1953. The influence of consonant environment upon the secondary acoustical characteristics of vowels. *The Journal of the Acoustical Society of America* 25:105–113.
- Hyman, Larry. 2019. *A theory of phonological weight*. De Gruyter Mouton.
- Iverson, Gregory K., and Joseph C. Salmons. 1995. Aspiration and laryngeal representation in Germanic. *Phonology* 12:369–396.
- Jacewicz, Ewa, Robert Al. Fox, and Samantha Lyle. 2009. Variation in stop consonant voicing in two regional varieties of American English. *Journal of the International Phonetic Association* 39:313.
- Javkin, Hector. 1977. Phonetic universals and phonological change. Doctoral Dissertation, U.C. Berkeley.
- Jessen, Michael. 2001. *Distinctive feature theory*, volume 2, chapter Phonetic implementation of the distinctive auditory features [voice] and [tense] in stop consonants, 237–294. Mouton de Gruyter Berlin.
- Johnson, Keith. 2004. Massive reduction in conversational American English. In *Proceedings of the Workshop on Spontaneous Speech: Data and Analysis*, 29–54.
- Jurafsky, D., A. Bell, M. Gregory, and W. D. Raymond. 2001. Probabilistic relations between words: Evidence from reduction in lexical production. In *Frequency and the emergence of linguistic structure*, ed. Joan L. Bybee and Paul J. Hopper, number 45 in Typological studies in language, 229–254. Amsterdam: John Benjamins.
- Jurafsky, Daniel, Alan Bell, Eric Fosler-Lussier, Cynthia Girand, and William Raymond. 1998.

Reduction of English function words in Switchboard. In *Fifth International Conference on Spoken Language Processing*.

Katsika, Argyro. 2016. The role of prominence in determining the scope of boundary-related lengthening in Greek. *Journal of phonetics* 55:149–181.

Kavitskaya, D. 2002. *Compensatory lengthening: Phonetics, phonology, diachrony*. London: Routledge.

Keating, Patricia A. 1979. A phonetic study of a voicing contrast in Polish. Doctoral Dissertation, Brown University.

Keating, Patricia A. 1985. Universal phonetics and the organization of grammars. In *Phonetic linguistics: Essays in honor of Peter Ladefoged*, ed. Victoria A. Fromkin, 115–132. Orlando: Academic Press.

Kessinger, Rachel H., and Sheila E. Blumstein. 1997. Effects of speaking rate on voice-onset time in Thai, French, and English. *Journal of Phonetics* 25:143–168.

Kim, Chin-Wu. 1970. A theory of aspiration. *Phonetica* 21:107–116.

Kim, Heejin, and Jennifer Cole. 2005. The stress foot as a unit of planned timing: evidence from shortening in the prosodic phrase. In *Ninth European Conference on Speech Communication and Technology*.

Kim, Sahyang, and Taehong Cho. 2012. Prosodic strengthening in the articulation of English /æ/. *Studies in Phonetics, Phonology and Morphology* 18:321–337.

Kingston, J., and R. L. Diehl. 1994. Phonetic Knowledge. *Language* 70:419–454.

Klatt, Dennis H. 1973. Interaction between two factors that influence vowel duration. *The Journal of the Acoustical Society of America* 54:1102–1104.

- Klatt, Dennis H. 1975. Vowel lengthening is syntactically determined in a connected discourse. *Journal of Phonetics* 3:129–140.
- Klatt, Dennis H. 1976. Linguistic uses of segmental duration in English: Acoustic and perceptual evidence. *The Journal of the Acoustical Society of America* 59:1208–1221.
- Kluender, Keith R., Randy L. Diehl, and Beverly A. Wright. 1988. Vowel-length differences before voiced and voiceless consonants: An auditory explanation. *Journal of Phonetics* 16:153–169.
- Ko, Eon-Suk. 2018. Asymmetric effects of speaking rate on the vowel/consonant ratio conditioned by coda voicing in English. *Phonetics and Speech Sciences* 10:45–50.
- Kohler, K. J. 1979. Dimensions in the perception of fortis and lenis consonants. *Phonetica* 36:332–343.
- Kohler, Klaus J. 1984. Phonetic explanation in phonology: the feature fortis/lenis. *Phonetica* 41:150–174.
- Kozhevnikov, Valeriĭ Aleksandrovich, and Liudmila Andreevna Chistovich. 1965. *Speech: Articulation and perception*. Nauka.
- Krause, Sue Ellen. 1982. Vowel duration as a perceptual cue to postvocalic consonant voicing in young children and adults. *The Journal of the Acoustical Society of America* 71:990–995.
- Kristoffersen, Gjert. 2000. *The phonology of Norwegian*. Oxford: Oxford University Press on Demand.
- Krivokapić, Jelena. 2020. *Prosodic theory and practice*, chapter Prosody in Articulatory Phonology. MIT Press. In press Cambridge, MA.
- Kulikov, Vladimir. 2012. Voicing and voice assimilation in Russian stops. Doctoral Dissertation, University of Iowa.

- Laeufer, Christiane. 1992. Patterns of voicing-conditioned vowel duration in French and English. *Journal of Phonetics* 20:411–440.
- Lehiste, Ilse. 1972. The timing of utterances and linguistic boundaries. *The Journal of the Acoustical Society of America* 51:2018–2024.
- Levy, R., and T. F. Jaeger. 2007. Speakers optimize information density through syntactic reduction. In *Advances in neural information processing systems*, ed. B. Scholkopf, J. Platt, and T. Hoffman, 849–856. MIT Press.
- Lindblom, B, and Karin Rapp. 1971. Reexamining the compensatory adjustment of vowel duration in Swedish words. *Stockholm, KTH, Speech Transmission Laboratory Quarterly Progress and Status Report* 4:19–25.
- Lisker, L. 1957a. Minimal cues for separating /w,r,l,y/ in intervocalic position. *Word* 13:256–267.
- Lisker, Leigh. 1957b. Closure duration and the intervocalic voiced-voiceless distinction in English. *Language* 33:42–49.
- Lisker, Leigh. 1974. On "explaining" vowel duration variation. *Glossa* 8:233–246.
- Lisker, Leigh. 1978. Rapid vs. Rabid: A catalogue of acoustic features that may cue the distinction. *Haskins Laboratories Status Report on Speech Research* 54:127–132.
- Lisker, Leigh. 1986. "Voicing" in English: A catalogue of acoustic features signaling /b/ versus /p/ in trochees. *Language and Speech* 29:3–11.
- Luce, Paul A. 1986. Neighborhoods of Words in the Mental Lexicon. Research on Speech Perception. Technical Report No. 6. Technical report, Indiana University.
- Luce, Paul A, and Jan Charles-Luce. 1985. Contextual effects on vowel duration, closure duration, and the consonant/vowel ratio in speech production. *The Journal of the Acoustical Society of America* 78:1949–1957.

Mack, Molly. 1982. Voicing-dependent vowel duration in English and French: Monolingual and bilingual production. *The Journal of the Acoustical Society of America* 71:173–178.

Maddieson, Ian. 1985. Phonetic cues to syllabification. *Phonetic linguistics: Essays in honor of Peter Ladefoged* 203–221.

Malécot, André. 1968. The force of articulation of American stops and fricatives as a function of position. *Phonetica* 18:95–102.

Massaro, Dominic W., and Michael M. Cohen. 1983. Consonant/vowel ratio: An improbable cue in speech. *Attention, Perception, & Psychophysics* 33:501–505.

Miller, Joanne L. 1981. *Perspectives on the study of speech*, chapter Effects of speaking rate on segmental distinctions, 39–74. Routledge.

Miller, Joanne L., and Thomas Baer. 1983. Some effects of speaking rate on the production of /b/ and /w/. *The Journal of the Acoustical Society of America* 73:1751–1755.

Miller, Joanne L., Kerry P. Green, and Adam Reeves. 1986. Speaking rate and segments: A look at the relation between speech production and speech perception for the voicing contrast. *Phonetica* 43:106–115.

Miller, Joanne L., and Lydia E Volaitis. 1989. Effect of speaking rate on the perceptual structure of a phonetic category. *Perception & Psychophysics* 46:505–512.

Moreton, Elliott. 2004. Realization of the English postvocalic [voice] contrast in F1 and F2. *Journal of Phonetics* 32:1–33.

Munhall, Kevin, Carol Fowler, Sarah Hawkins, and Elliot Saltzman. 1992. "Compensatory shortening" in monosyllables of spoken English. *Journal of Phonetics* 20:225–239.

Nagao, Kyoko, and Kenneth J. de Jong. 2007. Perceptual rate normalization in naturally produced rate-varied speech. *The Journal of the Acoustical Society of America* 121:2882–2898.

- Nam, Hosung, and Elliot Saltzman. 2003. A competitive, coupled oscillator model of syllable structure. In *Proceedings of the 15th international congress of phonetic sciences*, volume 1, 2253–2256.
- Nittrouer, Susan. 2004. The role of temporal and dynamic signal components in the perception of syllable-final stop voicing by children and adults. *The Journal of the Acoustical Society of America* 115:1777–1790.
- O'Dell, Michael, and Tommi Nieminen. 1999. Coupled oscillator model of speech rhythm. In *Proceedings of the XIVth international congress of phonetic sciences*, volume 2, 1075–1078.
- Ohala, J. J. 2011. Accommodation to the Aerodynamic Voicing Constraint and its Phonological Relevance. In *Proceedings of the 15th International Conference of Phonetic Sciences*, 64–67.
- Ohala, John J. 1981. The Listener as a source of sound change. In *Parasession on language and behavior*, ed. C. S. Masek, R. A. Hendrick, and M. F. Miller, 178–203. Chicago: Chicago Linguistics Society.
- Ohala, John J. 1983. The origin of sound patterns in vocal tract constraints. In *The production of speech*, ed. P. F. MacNeilage, 189–216. New York: Springer.
- O'Kane, Donal. 1978. Manner of vowel termination as a perceptual cue to the voicing status of postvocalic stop consonants. *Journal of Phonetics* 6:311–18.
- Oller, D. Kimbrough. 1973. The effect of position in utterance on speech segment duration in English. *The Journal of the Acoustical Society of America* 54:1235–1247.
- Peterson, Gordon E., and Ilse Lehiste. 1960. Duration of syllable nuclei in English. *The Journal of the Acoustical Society of America* 32:693–703.
- Pike, Kenneth L. 1945. *The intonation of american english*. University of Michigan Academic Press.

- Pind, Jörgen. 1995. Speaking rate, voice-onset time, and quantity: The search for higher-order invariants for two Icelandic speech cues. *Perception & Psychophysics* 57:291–304.
- Pitt, Mark A, Laura Dilley, Keith Johnson, Scott Kiesling, William Raymond, E Hume, and E Fosler-Lussier. 1997. Buckeye Corpus of Conversational Speech (2nd release). Department of Psychology, Ohio State University (Distributor), Columbus, OH, USA.
- Pitt, Mark A, Keith Johnson, Elizabeth Hume, Scott Kiesling, and William Raymond. 2005. The Buckeye corpus of conversational speech: Labeling conventions and a test of transcriber reliability. *Speech Communication* 45:89–95.
- Pluymaekers, Mark, Mirjam Ernestus, and R. Harald Baayen. 2005. Lexical frequency and acoustic reduction in spoken Dutch. *The Journal of the Acoustical Society of America* 118:2561–2569.
- Port, Robert F. 1976. The influence of speaking tempo on the duration of stressed vowel and medial stop in English trochee words. Doctoral Dissertation.
- Port, Robert F. 1979. The influence of tempo on stop closure duration as a cue for voicing and place. *Journal of Phonetics* 7(1):45–56.
- Port, Robert F. 1981. Linguistic timing factors in combination. *The Journal of the Acoustical Society of America* 69:262–274.
- Port, Robert F., and Fred Cummins. 1992. The English voicing contrast as velocity perturbation. In *Proceedings of the Second International Conference on Spoken Language Processing*, 1311–1314. Banff, Alberta, Canada.
- Port, Robert F., and Jonathan Dalby. 1982. Consonant/vowel ratio as a cue for voicing in English. *Perception & Psychophysics* 32:141–152.
- Port, Robert F, Jonathan Dalby, and Michael O’Dell. 1987. Evidence for mora timing in Japanese. *The Journal of the Acoustical Society of America* 81:1574–1585.



- Priva, Uriel Cohen. 2010. Constructing typing-time corpora: A new way to answer old questions. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 32, 43–48.
- Priva, Uriel Cohen. 2017. Not so fast: Fast speech correlates with lower lexical and structural information. *Cognition* 160:27–34.
- Raphael, Lawrence J. 1972. Preceding vowel duration as a cue to the perception of the voicing characteristic of word-final consonants in American English. *The Journal of the Acoustical Society of America* 51:1296–1303.
- Raphael, Lawrence J. 1975. The physiological control of durational differences between vowels preceding voiced and voiceless consonants in English. *Journal of Phonetics* 3:25–33.
- Raphael, Lawrence J. 1981. Durations and contexts as cues to word-final cognate opposition in English. *Phonetica* 38:126–147.
- Raphael, Lawrence J., Michael F. Dorman, Frances Freeman, and Charles Tobin. 1975. Vowel and nasal duration as cues to voicing in word-final stop consonants: Spectrographic and perceptual studies. *Journal of Speech and Hearing Research* 18:389–400.
- Raymond, William D, Mark Pitt, Keith Johnson, Elizabeth Hume, Matthew Makashay, Robin Dautricourt, and Craig Hilts. 2002. An analysis of transcription consistency in spontaneous speech from the Buckeye corpus. In *Seventh International Conference on Spoken Language Processing*.
- Repp, Bruno H., and David R. Williams. 1985. Influence of following context on perception of the voiced–voiceless distinction in syllable-final stop consonants. *The Journal of the Acoustical Society of America* 78:445–457.
- Revoile, S., J.M. Pickett, Lisa D. Holden, and David Talkin. 1982. Acoustic cues to final stop voicing for impaired- and normal- hearing listeners. *The Journal of the Acoustical Society of America* 72:1145–1154.

- Saltzman, Elliot, Hosung Nam, Jelena Krivokapic, and Louis Goldstein. 2008. A task-dynamic toolkit for modeling the effects of prosodic structure on articulation. In *Proceedings of the 4th International Conference on Speech Prosody (Speech Prosody 2008)*, Campinas, Brazil, 175–184.
- Sanker, Chelsea. 2019. Influence of coda stop features on perceived vowel duration. *Journal of Phonetics* 75:43–56.
- Sanker, Chelsea. 2020. A perceptual pathway for voicing-conditioned vowel duration. *Laboratory Phonology* 11.
- Schwartz, Geoffrey. 2010. Phonology in the speech signal-Unifying cue and prosodic licensing. *Poznań Studies in Contemporary Linguistics* 46:499–518.
- Selkirk, Elizabeth. 1982. *The structure of phonological representations part 2*, chapter The syllable, 337–383. Dordrecht: Foris.
- Sharf, Donald J. 1962. Duration of post-stress intervocalic stops and preceding vowels. *Language and speech* 5:26–30.
- Sharf, Donald J. 1964. Vowel duration in whispered and in normal speech. *Language and Speech* 7:89–97.
- Smith, Bruce L. 2002. Effects of speaking rate on temporal patterns of English. *Phonetica* 59:232–244.
- Solé, Maria-Josep. 2007. *Experimental approaches to phonology*, chapter Controlled and mechanical properties in speech, 302–321. Oxford University Press.
- Stetson, Raymond H. 1928. *Motor phonetics: A study of speech movements in action*. Dordrecht: Springer.
- Stevens, K. N., and A. S. House. 1963. Perturbation of vowel articulations by consonantal context: An acoustical study. *Journal of Speech and Hearing Research* 6:111–128.

- Summerfield, Q. 1981. Articulatory rate and perceptual constancy in phonetic perception. *Journal of Experimental Psychology: Human Perception and Performance* 7:1074–1095.
- Sweet, Henry. 1880. *A handbook of phonetics*. Clarendon Press Series. London: MacMillan and Co.
- Tanner, James, Morgan Sonderegger, Jane Stuart-Smith, and SPADE Data Consortium. 2019. Vowel duration and the voicing effect across English dialects. *Toronto Working Papers in Linguistics* 41.
- Toscano, Joseph C., and Bob McMurray. 2015. The time-course of speaking rate compensation: Effects of sentential rate and vowel length on voicing judgments. *Language, Cognition and Neuroscience* 30:529–543.
- Treiman, Rebecca, Kathleen Straub, and Patrick Laver. 1994. Syllabification of bisyllabic nonwords: Evidence from short-term memory errors. *Language and Speech* 37:45–59.
- Turk, Alice E., and Stefanie Shattuck-Hufnagel. 2007. Multiple targets of phrase-final lengthening in American English words. *Journal of Phonetics* 35:445–472.
- Umeda, Noriko. 1975. Vowel duration in American English. *The Journal of the Acoustical Society of America* 58:434–445.
- Van Heuven, Walter J.B., Pawel Mandera, Emmanuel Keuleers, and Marc Brysbaert. 2014. SUBTLEX-UK: A new and improved word frequency database for British English. *The Quarterly Journal of Experimental Psychology* 67:1176–1190.
- Van Summers, W. 1987. Effects of stress and final-consonant voicing on vowel production: Articulatory and acoustic analyses. *The Journal of the Acoustical Society of America* 82:847–863.
- Vatikiotis-Bateson, Eric. 1984. The temporal effects of homorganic medial nasal clusters. *Research in Phonetics* 4:197–233.

- Viswanathan, Navin, Annie J. Olmstead, and M. Pilar Aivar. 2019. The use of vowel length in making voicing judgments by native listeners of English and Spanish: Implications for rate normalization. *Language and Speech* 63:436–452.
- Volaitis, L. E., and J. L. Miller. 1992. Phonetic prototypes: Influence of place of articulation and speaking rate on the internal structure of voicing categories. *Journal of the Acoustical Society of America* 92:723–735.
- Walsh, Thomas, and Frank Parker. 1981. Vowel length and voicing in a following consonant. *Journal of Phonetics* 9:305–308.
- Wardrip-Fruin, Carolyn. 1982. On the status of temporal cues to phonetic categories: Preceding vowel duration as a cue to voicing in final stop consonants. *The Journal of the Acoustical Society of America* 71:187–195.
- Weismer, Gary. 1979. Sensitivity of voice-onset time (VOT) measures to certain segmental features in speech production. *Journal of Phonetics* 7:197–204.
- Wells, John C. 1982. *Accents of English*, volume 1. Cambridge: Cambridge University Press.
- Wightman, C. W., S. Shattuck-Hufnagel, M. Ostendorf, and P. J. Price. 1992. Segmental durations in the vicinity of prosodic phrase boundaries. *Journal of the Acoustical Society of America* 91:1707–1717.

	Estimate	Std. Error	df	t value	Pr(> t )
(intercept)	4.80	0.04	29.8	108	< 2e-16
Vowel Height	0.17	0.03	10.4	4.95	5.1e-4
Vowel Class	0.11	0.03	10.7	3.56	0.005
Voicing	-0.008	0.02	120	-0.33	0.74
Speaking Rate	0.67	0.02	63.8	31.6	< 2e-16
Obstruent Type	0.02	0.02	143	1.30	0.20
Frequency	-0.11	0.03	124	-3.60	4.63e-4
Vowel Height: Voicing	-0.02	0.02	208	-1.40	0.16
Voicing: Rate	0.06	0.02	63.5	2.79	0.007
Voicing: Obs. Type	-0.01	0.01	233	-0.85	0.39
Voicing: Frequency	-0.06	0.03	234	-2.18	0.03
Vowel Class: Voicing	0.03	0.02	197	1.75	0.08

**Table 1**

*Utterance Final: Largest Converging Model.*

	Estimate	Std. Error	df	t value	Pr(> t )
(intercept)	4.81	0.03	110	140	< 2e-16
Vowel Height	0.16	0.02	223	9.11	< 2e-16
Vowel Type	0.11	0.02	214	6.59	3.3e-10
Speaking Rate	0.67	0.02	70.8	29.6	< 2e-16
Voicing	0.04	0.02	210	2.30	0.02
Obstruent Type	0.02	0.02	160	1.33	0.18
Frequency	0.09	0.03	129	-3.01	0.003

**Table 2**

*Utterance-Final: Simplified Model: No interactions*

	Estimate	Std. Error	df	t value	Pr(> t )
(intercept)	4.79	0.04	41.4	109	< 2e-16
Vowel Height	0.18	0.03	10.29	6.26	8.3e-5
Vowel Type	0.13	0.03	10.5	5.06	0.0004
Voicing	-0.06	0.02	146	2.63	0.01
Speaking Rate	1.07	0.02	3.1e3	56.0	< 2e-16
Consonant Duration	-0.55	0.03	80.0	-21.4	< 2e-16
Obstruent Type	-0.07	0.02	151	-4.43	1.8e-05
Frequency	-0.07	0.03	133	-2.50	0.01
Voicing: Rate	0.14	0.02	2.8e3	7.23	6.4e-13
Voicing: Obs. Type	-0.03	0.01	278	-1.88	0.06
Voicing: Frequency	-0.04	0.02	244	-1.54	0.13
Voicing:Consonant Duration	-0.11	0.02	2.7e3	-5.67	1.7e-8

**Table 3**

*Non-Utterance Final: Largest Converging Model, with consonant duration added.*

	Estimate	Std. Error	df	t value	Pr(> t )
(intercept)	4.80	0.04	82.7	126	< 2e-16
Vowel Height	0.16	0.02	201	10.4	< 2e-16
Vowel Type	0.12	0.01	180	8.62	3.5e-15
Speaking Rate	1.03	0.02	146	42.7	< 2e-16
Consonant Duration	-0.52	0.02	62.9	-21.3	< 2e-16
Voicing	-0.007	0.01	181	-0.45	0.67
Obstruent Type	-0.05	0.02	137	-3.42	0.0008
Frequency	-0.06	0.03	122	-2.41	0.02

**Table 4**

*Utterance-Final: Simplified Model with consonant duration.*



	Estimate	Std. Error	df	t value	Pr(> t )
(intercept)	4.37	0.03	109	137	< 2e-16
Vowel Height	0.17	-0.01	310	12.35	< 2e-16
Vowel Type	0.15	0.01	323	11.60	< 2e-16
Speaking Rate	0.74	0.007	1.6e4	100	< 2e-16
Consonant Duration	-0.22	0.01	61.6	-18.9	< 2e-16
Obstruent Type	-0.06	0.01	212	-5.57	5.5e-08
Frequency	-0.07	0.02	262	-3.46	6.4e-4

**Table 5**

*Non-Utterance-Final only: Simplified Model with consonant duration. Voicing excluded.*

	Estimate	Std. Error	df	t value	Pr(> t )
(intercept)	-0.10	0.03	139	-2.93	0.004
Voicing	-0.06	0.03	103	-2.35	0.02
Speaking Rate	0.66	0.02	113	33.9	< 2e-16
Obstruent Type	-0.11	0.01	279	-9.21	< 2e-16
Frequency	-0.02	0.02	254	-1.02	0.309
Voicing: Rate	-0.03	0.01	119	-2.02	0.05
Voicing: Frequency	0.05	0.02	264	2.46	0.01
Voicing: Obs. Type	-0.02	0.01	1098	-1.7	0.09

**Table 6**

*Non-Utterance-Final only: Consonant Duration model.*

**Table 7**

*Mixed-Effects Linear Regression Model of vowel duration as a function of speaking rate and consonant duration and their interaction.*

	Estimate	Std. Error	estimated df	t-value	p-value
(Intercept)	355.0	25.05	27.31	14.17	4.13e-14
rate.L	324.5	15.04	969.9	21.58	< 2e-16
rate.Q	42.49	15.16	968.7	2.802	0.005
Consonant Duration	-0.169	0.059	977.7	-2.879	0.004
rate.L:C Duration	-0.637	0.113	969.3	-6.657	2.03e-8

**Table 8**

*Mixed-Effects Linear Regression Model of vowel duration as a function of speaking rate, voicing, and consonant duration with full interactions.*

	Estimate	Std. Error	estimated df	t-value	p-value
(Intercept)	349.6	27.57	24.77	12.68	2.51e-12
rate.L	385.23	24.88	960.8	15.48	< 2e-16
rate.Q	116.6	24.98	959.5	4.688	3.48e-6
Voicing	5.086	30.48	10.79	0.167	0.871
Consonant Duration	-0.281	0.068	967.8	-4.109	4.3e-5
rate.L:Voicing	120.2	33.99	959.9	-3.536	.0004
rate.Q:Voicing	-107.8	34.57	959.5	-3.119	0.0019
rate.L:Consonant Duration	-0.884	0.148	960.5	-5.990	2.96e-9
rate.Q:Consonant Duration	-0.392	0.143	959.7	-2.733	0.006
voicing: consonant duration	0.332	0.123	963.4	2.700	0.007
rate.L:voicing:Consonant Duration	0.694	0.276	959.7	2.513	0.012

**Table 9**

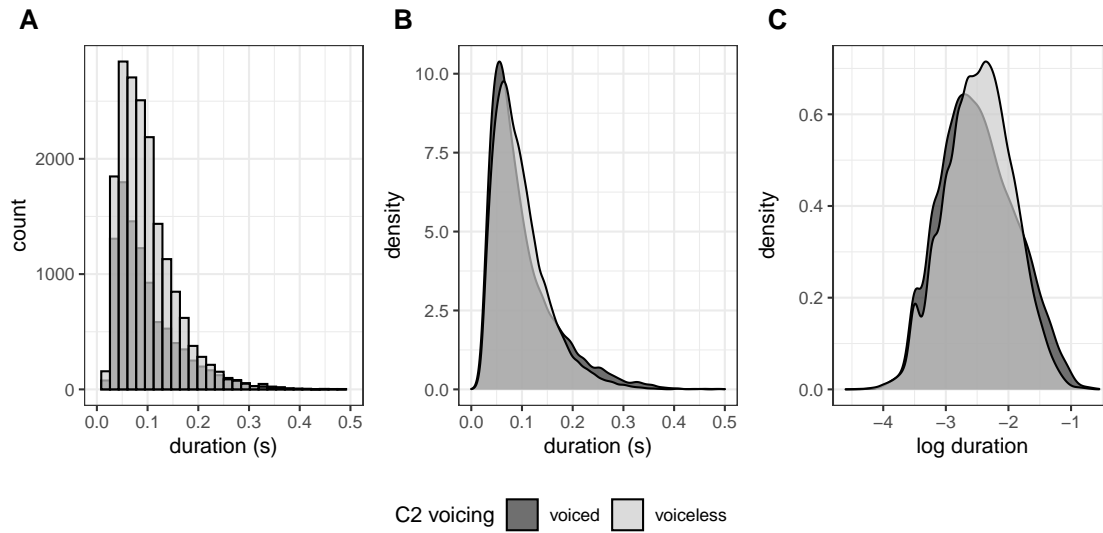
*Mixed-Effects Linear Regression Model of consonant duration as a function of speaking rate, voicing, and manner, with full interactions.*

	Estimate	Std. Error	estimated df	t-value	p-value
(Intercept)	160.1	7.253	7.175	22.07	7.37e-8
rate.L	87.84	6.835	960.2	12.85	< 2e-16
voicing	-54.96	8.804	3.987	-6.243	0.003
manner	28.83	8.812	4.001	3.272	0.031
rate.L:voicing	-46.76	9.702	960.2	-4.819	1.67e-6
rate.L:manner	19.59	9.715	960.2	2.016	0.044

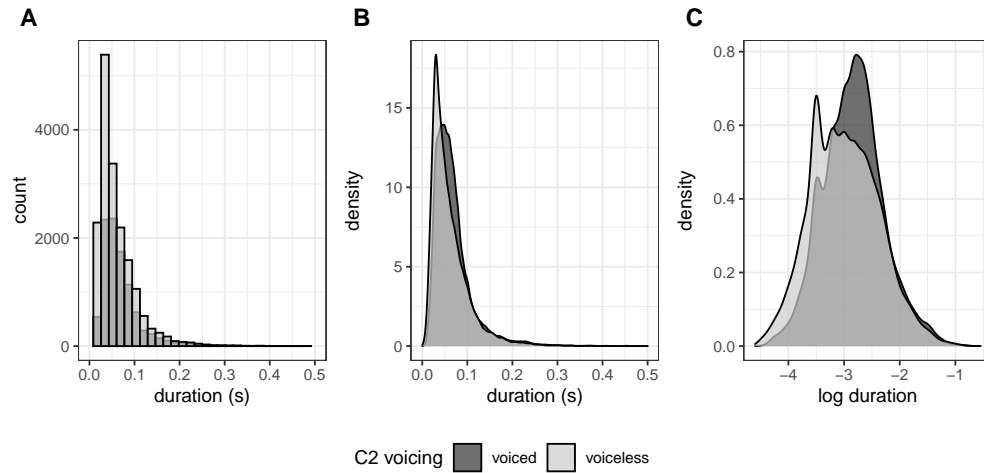
**Table 10**

*Mixed-Effects Linear Regression Model of vowel duration difference as a function of consonant duration difference, manner and rate, with full interactions.*

	Estimate	Std. Error	estimated <i>df</i>	<i>t</i> -value	<i>p</i> -value
(Intercept)	-57.37	12.34	6.891	-4.649	0.002
$\Delta C$	-0.217	0.064	475.8	-3.390	0.001
rate.Q	30.11	10.94	451.8	2.751	0.006
rate.L: $\Delta C$	-0.376	0.139	460.4	-2.701	0.007
rate.Q: $\Delta C$	-0.277	0.133	455.8	-2.081	0.038

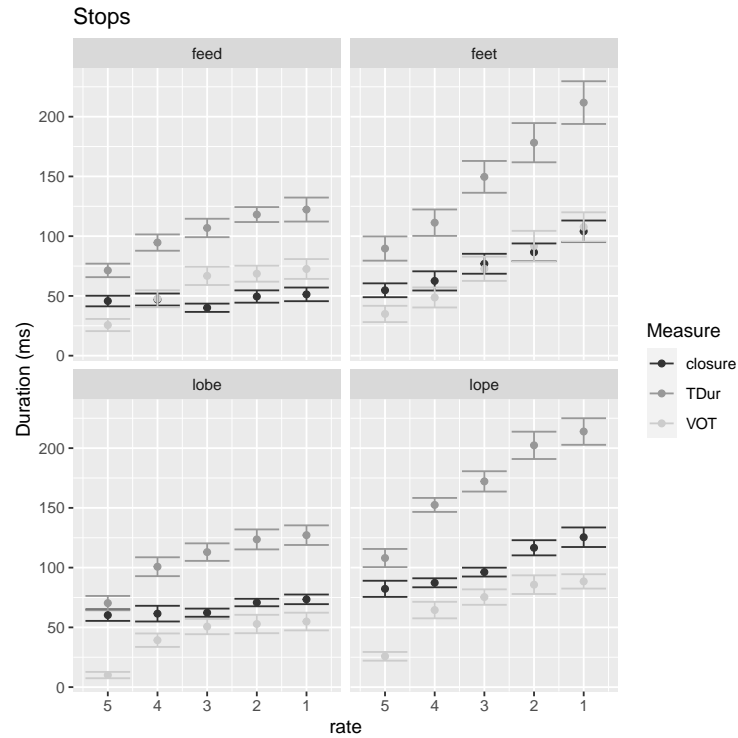


**Figure 1**  
*All CVC vowel durations*

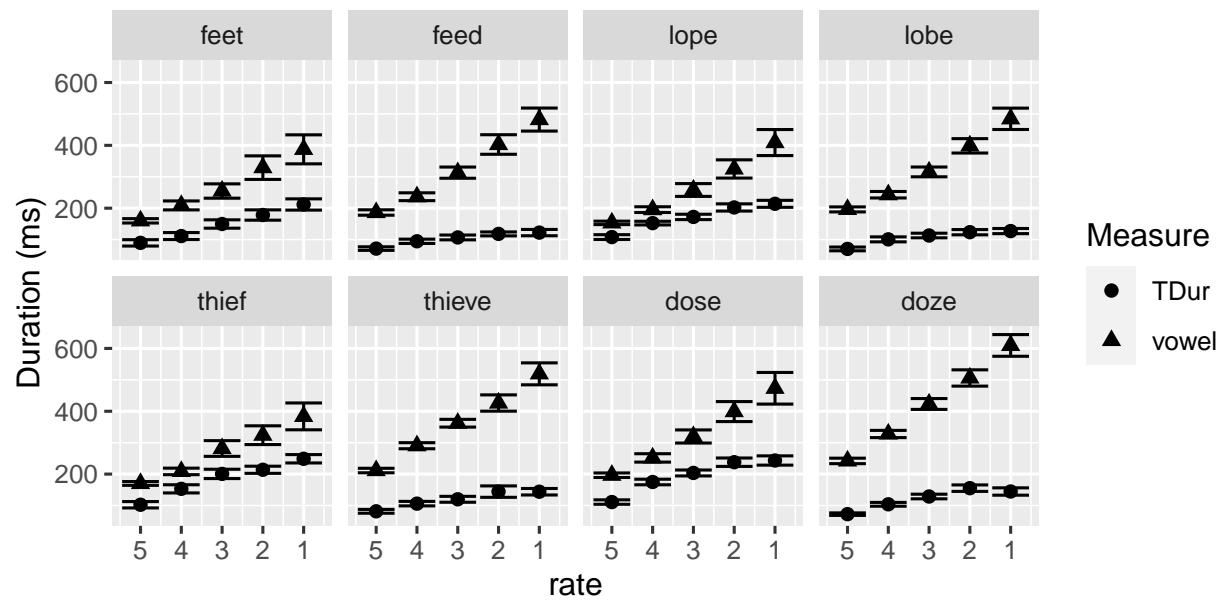


**Figure 2**  
*Obstruent durations by fricative and stop (flaps excluded).*

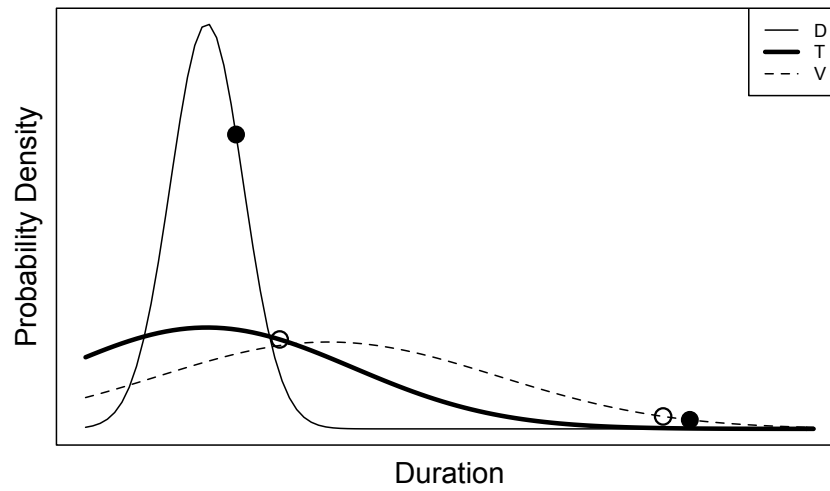


**Figure 3**

*Closure duration, VOT and Total duration (TDur) for final stops as a function of repetition rate (decreasing from left to right). Means and standard error bars.*

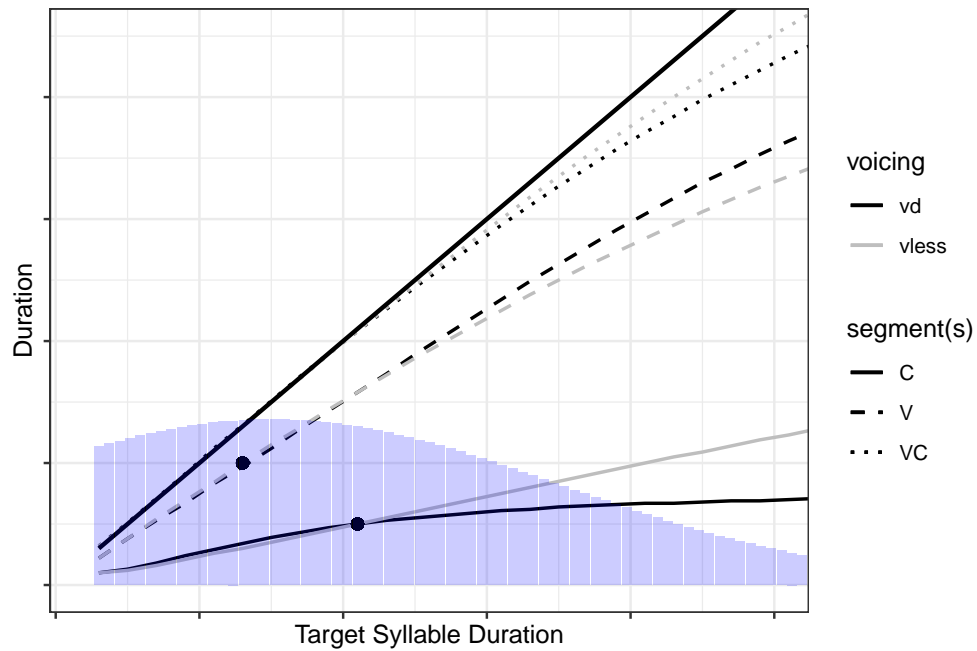
**Figure 4**

*Total consonant duration, preceding vowel duration, and rhyme duration (V+C), as a function of repetition rate. Means and standard error bars.*

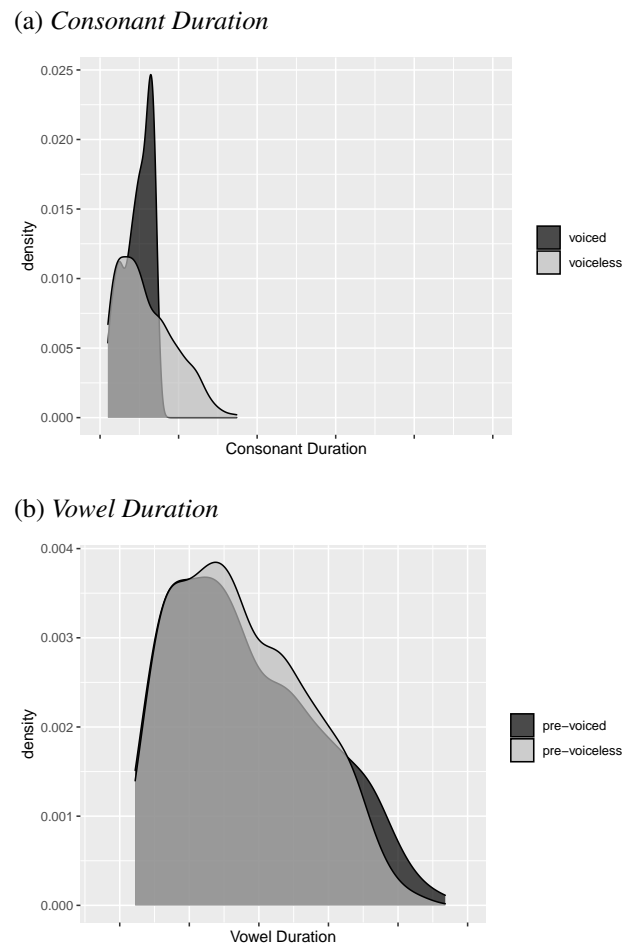


**Figure 5**

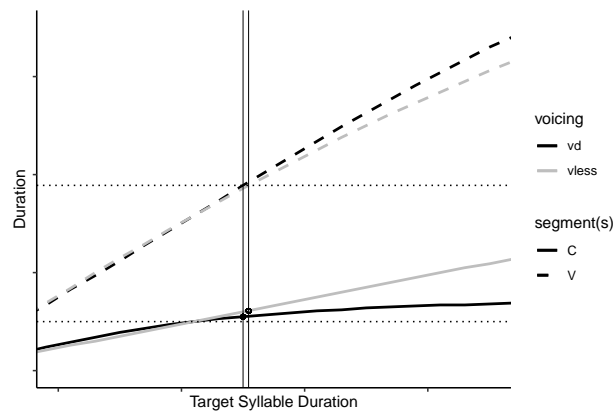
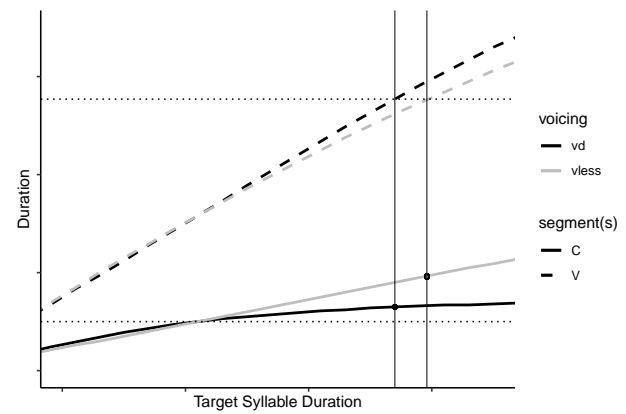
*Probability densities for: voiced obstruent (solid); voiceless obstruent (thick solid); vowel (dashed). Open circles: VT syllable with target syllable duration of 330 ms. Filled circles: VD syllable with target syllable duration of 330 ms.*

**Figure 6**

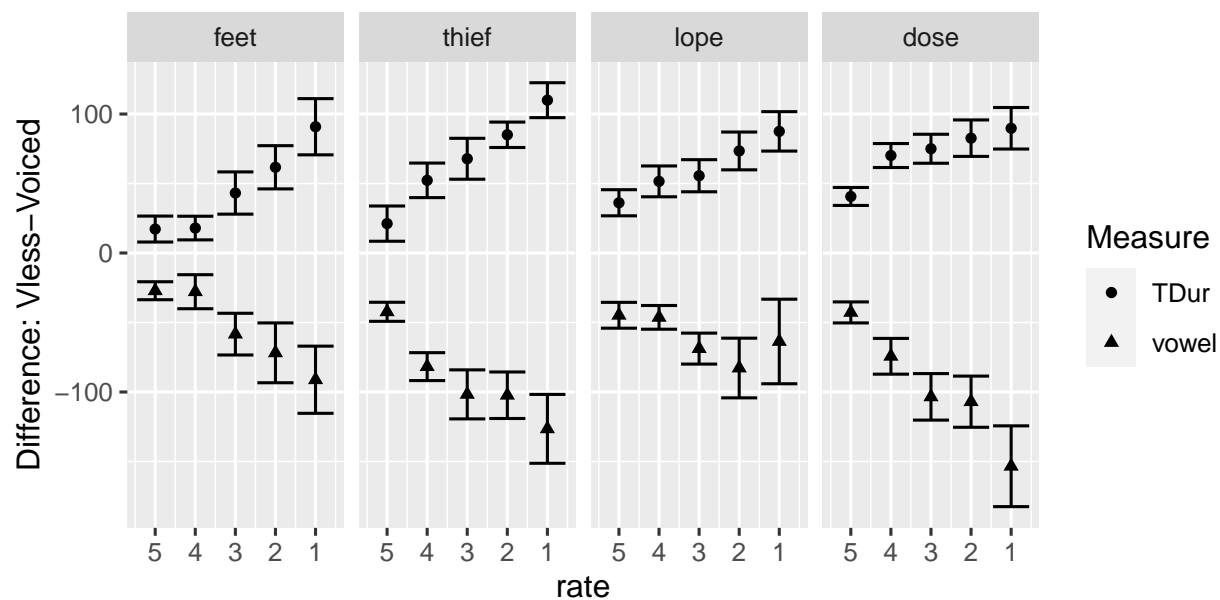
*Behavior of the Competing Constraints Model of segment duration as a function of target syllable duration. Actual and target syllable duration are equal along the upper solid black line. Vertical blue bars show the distribution used to represent the durations/rates in the Buckeye Corpus (used for the corpus simulation).*



**Figure 7**  
*Corpus Simulation. Target syllable durations randomly sampled from the Normal distribution shown in Fig. 6.*

(a) *Short vowel token*(b) *Long vowel token***Figure 8**

*Competition Simulation: Observed vowel duration is marked by the upper dotted line in both figures. Vertical solid lines intersect expected target syllable duration (speaking rate), and expected stop duration. Left line: voiced stop coda; Right line: voiceless stop coda. The lower dotted line indicates the actual duration of the following stop.*



**Figure 9**  
*Production differences by voicing for each minimal pair.*

## Appendix A

### Word Lists

#### CV Stop.

**Voiced (3071 tokens; 83 unique words) , with individual counts:** bad(111), bag(3), bed(14), big(154), bob(2), cab(3), cad(1), cod(1), code(2), could(206), cub(1), dab(2), dad(60), dead(11), did(232), died(8), dig(1), dog(18), dude(2), fed(4), feed(3), fog(4), food(20), gig(1), god(35), good(245), guide(1), had(384), head(16), hid(1), hide(4), hood(1), hub(1), hug(2), hyde(2), jed(1), jedd(2), job(104), kid(92), knob(1), lab(1), lag(1), laid(5), lead(8), league(8), led(2), leg(3), lied(1), load(2), loud(10), mad(5), made(63), med(4), meg(1), mid(4), mud(1), need(150), paid(31), pig(2), read(43), red(7), rid(1), ride(9), road(22), rob(2), rub(1), sad(5), said(311), shed(2), should(140), showed(5), side(32), sued(2), tag(2), ted(1), tied(1), todd(1), tub(1), tube(1), web(6), weed(1), wide(1), would(416)

**Voiceless (8663 tokens; 173 unique words), with individual counts:** back(341), beat(9), beep(1), bet(9), bike(1), bit(27), bite(1), boat(3), book(23), bought(9), buck(3), but(880), butt(2), cake(4), cap(2), cape(1), cat(2), caught(4), chalk(1), cheap(11), check(15), chick(1), chip(1), coke(1), cook(11), cop(9), cope(1), cup(1), cut(12), date(5), deck(4), deep(9), dip(1), dot(6), doubt(1), duck(1), duke(1), fake(5), fat(3), feet(4), fight(5), fit(10), folk(3), foot(1), fuck(1), gap(1), gate(2), get(315), got(156), gut(1), hate(14), heat(1), heck(14), height(1), hick(1), hip(2), hit(11), hook(4), hop(4), hope(29), hot(6), hype(2), jack(1), jeep(2), jet(1), jock(2), joke(6), keep(88), kick(3), knit(2), lack(12), lake(7), lap(1), late(6), let(28), light(5), like(2537), lock(10), look(149), lot(75), luck(1), luke(1), mac(1), make(225), map(2), meet(9), met(17), might(44), mike(3), mock(2), nap(1), neat(8), neck(4), net(1), night(14), nope(6), nose(5), not(322), note(1), nut(1), pack(5), peek(1), pet(1), pete(1), pick(31), pipe(4), poke(1), pop(1), pope(2), pot(1), psych(2), puck(1), put(2), rat(3), rate(2), rec(2), right(197), rock(6), rope(1), route(3), sake(3), sat(6), seat(1), set(13), shake(1), shape(7), sheet(2), ship(4), shit(4), shock(4), shoot(17), shop(6), shot(4), shut(1), sick(7), sit(39), site(3), soap(7), soup(5), suit(3), take(255), talk(125), tap(1), tape(11), taught(8), tech(5), that(1413), thick(1), this(49), thought(91), tight(1),



tip(4), took(88), top(25), type(48), vote(20), wait(11), wake(2), week(78), weight(1), wet(2), whack(3), what(346), whip(2), white(6), wick(1), woke(2), wreck(1), wright(1), write(13), wrote(6), yet(22), zip(4)

### **7.0.1 CV Fricative**

**Voiced (4912 tokens; 69 unique words), with individual counts:** b's(4), boys(32), c's(4), cahs(1), cause(53), cave(1), cheese(4), choose(15), chose(3), cows(1), d's(8), days(50), dies(2), does(116), dos(1), faze(1), five(182), gave(24), gays(10), give(100), goes(116), guys(71), has(194), have(980), hayes(4), haze(1), his(154), jazz(3), joe's(1), keys(2), knees(2), knows(31), laws(21), leave(37), live(137), lose(10), love(85), move(61), news(38), noise(2), p's(1), pays(2), phase(1), raise(19), rave(1), rise(1), rose(1), save(9), says(78), seas(1), sees(9), shave(1), shoes(13), shows(12), size(8), t's(1), taj(1), these(216), those(212), ties(1), toes(1), toys(2), twos(3), use(63), was(1634), wave(1), ways(41), whose(7), wise(10)

**Voiceless (1823 tokens; 69 unique words), with individual counts:** base(9), bash(1), bass(1), beef(1), biff(1), boss(3), bus(17), bush(8), calf(1), case(28), cash(3), chess(1), chief(7), choice(24), cuff(5), cuss(1), dose(1), face(23), fish(3), fuss(2), gas(6), geese(1), goose(1), gosh(28), guess(140), half(68), hash(1), house(134), joyce(2), juice(1), kiss(1), knife(1), las(1), laugh(2), lease(8), less(38), life(137), loose(1), los(1), mass(4), mess(6), mice(1), miss(13), moss(1), nice(86), niche(4), niece(3), pace(1), peace(7), piece(10), piss(1), push(12), race(6), rash(1), reese(1), rice(1), rough(12), rush(1), safe(5), this(707), tiff(1), toss(1), tough(21), vice(4), voice(7), wash(4), wife(47), wish(21), yes(122)

## Appendix B

### Campbell's Elasticity Model

The model in Campbell (1992) is based on the hypothesis that a single expansion coefficient ( $\epsilon_k$ ) can be applied to all segments ( $S_i$ ) within a given syllable ( $\sigma_k$ ). Duration is determined by baseline duration and segment elasticity. Mean segment duration is used as the baseline, and standard deviation is used as proxy for elasticity. Thus each segment is expanded by the same number of standard deviations ( $\kappa_i$ ) from its own mean. See Equation (6).

$$S_i = \bar{S}_i + \kappa_i \epsilon_k \quad (6)$$

The appropriate expansion coefficient is found by taking the difference between the baseline syllable duration ( $\bar{\sigma}$ ) and the target syllable duration ( $\sigma_T$ ), divided by the sum of the elasticities of all segments within the syllable. See Equation (7).

$$\epsilon_k(\sigma_T) = \frac{\sigma_T - \bar{\sigma}_k}{\sum_{i \in k} \kappa_i} \quad (7)$$

For a voiced/voiceless minimal VC syllable pair (VD, VT), with the same target duration,

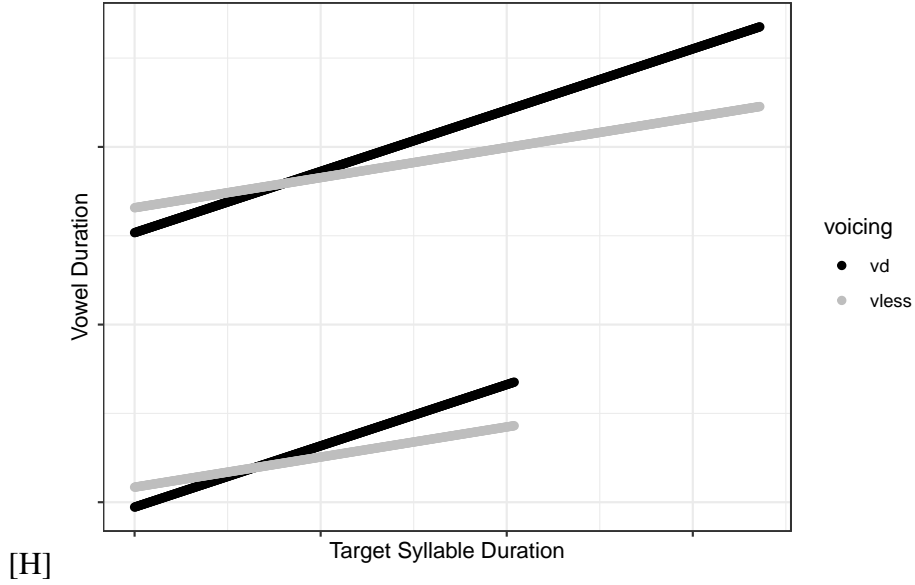
$\epsilon_{vd} = \frac{\sigma_T - (\bar{V} + \bar{D})}{\kappa_V + \kappa_D}$ , and  $\epsilon_{vless} = \frac{\sigma_T - (\bar{V} + \bar{T})}{\kappa_V + \kappa_T}$ . Therefore,  $\epsilon_{vless} = \frac{\epsilon_{vd}(\kappa_V + \kappa_D) + (\bar{D} - \bar{T})}{\kappa_V + \kappa_T}$ . Because the elasticity of the voiced obstruent is less than the elasticity of the voiceless obstruent,  $\frac{\kappa_V + \kappa_D}{\kappa_V + \kappa_T}$  is less than 1. Therefore, as target syllable duration increases,  $\epsilon_{vless}$  increases more slowly than  $\epsilon_{vd}$ .

However, when target syllable duration is short enough to lead to compression (negative values for  $\epsilon_{vd}$ ), the opposite relation holds. A higher elasticity means a segment can both lengthen more, and compress more. The difference between the more elastic and the less elastic segment will continue to increase in both directions away from the mean.

The voicing effect,  $V_{vd} - V_{vless}$ , is given by  $\kappa_V \epsilon_{vd} - \kappa_V \epsilon_{vless}$ , which can be rewritten as  $\epsilon_{vd} \kappa_V (1 - \frac{\kappa_V + \kappa_D}{\kappa_V + \kappa_T}) + \frac{\kappa_V (\bar{D} - \bar{T})}{\kappa_V + \kappa_T}$ .  $\epsilon_{vd}$ , in turn, is a linear function of  $\sigma_T$ . Thus the magnitude of the voicing effect, for any given target syllable duration, can be determined via:

$\frac{\sigma_T - (\bar{V} + \bar{D})}{\kappa_V + \kappa_D} \kappa_V (1 - \frac{\kappa_V + \kappa_D}{\kappa_V + \kappa_T}) + \frac{\kappa_V (\bar{D} - \bar{T})}{\kappa_V + \kappa_T}$ . As long as  $\epsilon_{vd}$  is positive, the difference in duration between

pre-voiced and pre-voiceless vowels will increase linearly with target syllable duration, at a rate given by  $\kappa_V(\frac{1}{\kappa_V + \kappa_D} - \frac{1}{\kappa_V + \kappa_T})$ . When  $\epsilon_{vd}$  is negative, the voicing effect reverses at the same rate. A somewhat simplified version of Campbell's model,<sup>37</sup> using only the above equations, was used to generate Fig. B1. The upper and lower sets of lines correspond to inherently longer and shorter vowels, respectively.



**Figure B1**  
Campbell's Simplified Model: VC syllables.

For a VC syllable,  $\sigma = V + C$ . If the change in duration for a given segment, S, is denoted by  $\Delta S$ , then  $\sigma_T = \bar{V} + \Delta V + \bar{C} + \Delta C$ . For a given target duration (larger than the mean), a larger expansion coefficient is required for the voiced syllable, which has the effect of lengthening the pre-voiced vowel more than the pre-voiceless, and is the source of the voicing effect. More specifically, the vowel in the voiced syllable must be lengthened by precisely the amount necessary to compensate both for the discrepancy between the expansion of the voiceless and

<sup>37</sup> Target syllable duration was treated as a random variable, ranging over multiples of the baseline syllable duration, rather than being fit to the phonological properties of the syllable.

voiced obstruent, and for the difference between their mean durations:

$$\Delta V_{vd} - \Delta V_{vless} = (\bar{T} - \bar{D}) + \Delta T - \Delta D.$$

The difference between inherently long and inherently short vowels is modeled at the syllable level by assigning different mean durations to the two kinds of syllables. The result is that “short” vowels undergo less lengthening, on average, than “long” vowels.<sup>38</sup> This also means that a difference in the magnitude of the voicing effect for shorter versus longer vowels appears only in the aggregate data. Pair-wise comparisons (at the same target duration) between long and short vowel syllables will show no difference in the magnitude of the voicing effect. Note that all results rely on assigning a smaller mean and variance to the voiced obstruent than to the voiceless, even though differences between the two are small to non-existent in the Buckeye Corpus.

---

<sup>38</sup> It is worth noting that the difference in vowel type cannot be captured at the level of the segment in Campbell’s model. The segment-level modeling implicitly requires the ability to lengthen to any degree. A larger expansion coefficient is simply applied to less elastic segments in order to achieve the same length. Thus, not only would short vowels get as long as long vowels, a larger voicing effect would occur in short versus long-vowel syllables because segments would be subjected to a larger  $\varepsilon$  on average, the opposite of what is observed.

## Appendix C

### Competing Constraints Model

**VC syllables.** Constraints in this model are realized as Normally distributed probability densities. Probability decreases in either direction away from a maximum at the segment's preferred duration ( $\mu$ ); the rate of decrease is determined by the variance of the distribution, which is determined by the elasticity of the segment. Probability densities function as gradient constraints under optimization of the joint probability. When preferred segment durations conflict with one another, the highest joint probability is achieved by violating lower-ranked constraints: i.e., shifting segments with higher elasticity further away from their preferred durations so that lower elasticity segments can remain closer to theirs.

The full set of constraints for the competing constraints model is given in (8), along with the parameter values used for the simulations. The mean values for the D, T and V distributions are roughly in line with observed values. The same is true of the relative variances: D has the smallest, then T, and V with the largest. The actual values for the variances, however, were chosen to produce differences large enough to exhibit the desired behavior. This is not problematic because actual duration variance is not equivalent to elasticity.

$$(8) \quad P\left(\frac{C}{V}\right) \sim \mathcal{N}(\mu = .3, \sigma = .1)$$

$$P(D) \sim \mathcal{N}(\mu = 50, \sigma = 15)$$

$$P(T) \sim \mathcal{N}(\mu = 50, \sigma = 60)$$

$$P(V) \sim \mathcal{N}(\mu = 100, \sigma = 70)$$

$$P\left(\frac{\sigma_T - \sigma}{\sigma_T}\right) \sim \mathcal{N}(\mu = 0, \sigma = .07)$$

$$\frac{V}{\sigma} : V \text{ cannot be shorter than half the total syllable duration}$$

The optimization function for this model, as a function of  $\sigma_T$ , and for any consonant, vowel pair  $(y, z)$ , and under the assumption of independence, is given as

$$p(y, z, \sigma_T) = p\left(\frac{C}{V} = \frac{y}{z}\right) \cdot p(C = y) \cdot p(V = z) \cdot p\left(\frac{\sigma_T - (y + z)}{\sigma_T}\right) \quad (9)$$

To reduce run time, the constraint ( $\frac{V}{\sigma}$ ) is implemented by simply restricting the search space.<sup>39</sup> The `dnorm()` functions in R (v 1.4.1106) are used for the probability functions, with means and variances specified above.

The target syllable durations used for the corpus simulation were sampled from a Normal distribution with  $\mu$  equal to 150 ms, and a  $\sigma$  of 200. These parameters were chosen to reflect the fact that almost all corpus vowel durations fell below the experimental cross-over point between voiceless and voiced percepts. Thus, sampled syllable durations are chosen to cluster in a range where there was little difference between the duration of the two obstruents (durations are not allowed to fall below 30 ms.). The same 1000 point sample of targets was used for both the voiced and voiceless distributions. The results are given in Section 4.2. See Figures 6 and 7.

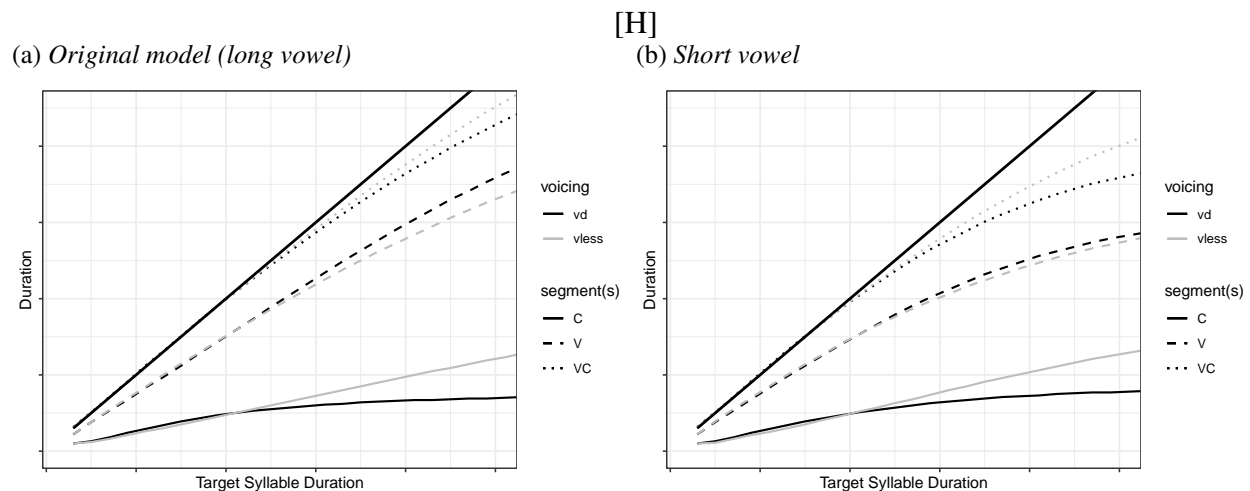
**Short Vowels:** We assume that target syllable duration is determined by a combination of speaking rate and other prosodic factors, such as phrase-final lengthening (see, e.g., Byrd and Saltzman 2003). In Campbell (1992), shorter vowels were essentially given a smaller range of possible target durations relative to longer vowels. This approach is not without justification, given that the nucleus type is often treated as equivalent to the syllable type, and different syllable types may have their own associated duration ranges. A similar approach could be implemented in the competing constraints model. However, the architecture of our model offers an alternative way to differentiate long and short vowels. A short vowel, like a short consonant, can be specified with a lower elasticity. Because this is not a simple temporal compensation model, lower vowel elasticity does not automatically lead to significantly longer consonant durations.

Figure C1 shows the result of reducing the variance of the vowel probability distribution ( $\sigma = 40$ : intermediate between that of the voiceless obstruent, and that of the voiced obstruent). The original (long vowel) model results are included for comparison. All other parameters remained the same, including the mean of the vowel distribution. The result is a smaller duration difference between the paired voiced/voiceless vowels, and shorter syllables over-all. Both

---

<sup>39</sup> Restricting the range effectively removes all values below a certain probability from consideration. Because this occurs before the joint probability is calculated, the restriction cannot be altered, making this constraint inviolable.

obstruent types become slightly longer under these conditions, but the largest change is in how closely the target syllable duration is approximated. In this model, greater target undershoot results in a higher joint probability than increased lengthening of either consonant. Qualitatively, this behavior is consistent with the finding that the voicing effect is significantly reduced in preceding vowels that are inherently short (Umeda 1975; Crystal and House 1982; De Jong 2004). Note that the difference in duration between the obstruents themselves can, in principle, still grow quite large. Because very few studies on the voicing effect report final obstruent durations, it remains to be seen whether this prediction is borne out.



**Figure C1**  
*Competing Constraints Model*

**CV syllables.** Given that differences in duration between voiced and voiceless obstruents in initial position have been shown to have some effect on the duration of following, tautosyllabic vowels, it should be possible to develop a broadly similar competing constraints model that accounts for these differences. Pre-vocalic and post-vocalic consonantal gestures are phased differently in English; singleton consonants in onset are activated at the same time as the vowel, whereas coda consonants are activated at the offset of the vowel (e.g., Browman and

Goldstein 1988). See Section 4.2. The onset-vowel phasing relationship is also less variable (e.g., Selkirk 1982), resisting adjustments that would shift the two segments apart. As a result, part of the vowel is consistently masked by the consonant, and thus acoustically shorter than a bare vowel. Inherently longer consonants will lead to greater masking than inherently shorter consonants. Thus, it is predicted that the vowel following a voiceless obstruent will be acoustically shorter than a vowel following a voiced one, but only in the range where voiceless obstruents are longer than their voiced counterparts.

By making two changes to the VC model, the predicted behavior of the onset “voicing” effect can be reasonably well-captured. The target-matching constraint is altered to apply only to the vowel, and not to the onset of the CV syllable. This assumption is necessary to produce a different outcome from the VC case. But it is also based on the fact that onsets do not typically take part in prosodic phenomena, being irrelevant to the calculation of syllable weight, for example (e.g., Hyman 2019). Changing the relevant unit from syllable to rhyme in Eq. 9 will cover both the VC and CV cases. The rhyme in the VC case is calculated by adding the durations of the consonant and vowel (assuming no overlap).<sup>40</sup> The rhyme in the CV case is calculated from the articulatory duration of the vowel; the acoustic duration of the vowel is given by subtracting consonant duration from articulatory duration (assuming full overlap).

The second change is a reduction in the variance of the C/V constraint (by 60%, measured with respect to the articulatory duration of the vowel). By making the variance smaller, the outcome becomes more strongly biased towards the preferred C/V value than it is towards perfect rhyme duration matching.<sup>41</sup> This is also consistent with the lower variability in phasing between onset and nucleus, versus nucleus and coda. All other model parameters remain the same.

These changes alter the model behavior in the desired ways. See Figure C2. The VC model (long vowel) is included for comparison. Even for the small set of constraints used here,

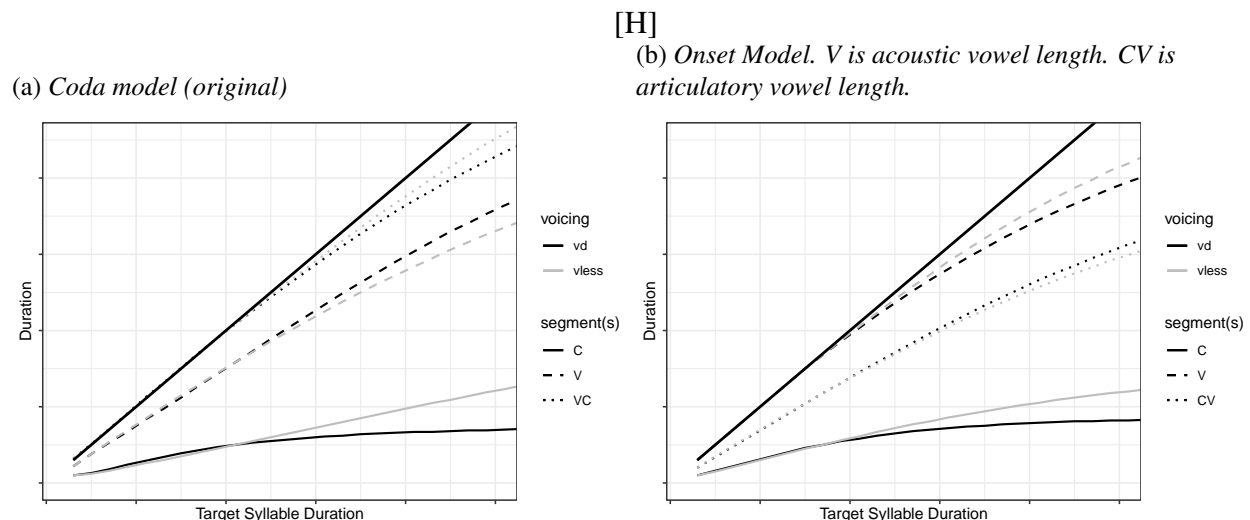
---

<sup>40</sup> If the operative unit is the rhyme, and duration is specified at that level, then this model can also account quite simply for the finding that vowels in open syllables are typically longer than vowels in closed syllables. For the same target rhyme duration, and a highly weighted target matching constraint, the same duration is distributed over two segments in the (C)VC case, and only one in the (C)V case.

<sup>41</sup> A smaller variance will increase resistance both to C/V values that are too large, as well as those that are too small.



the interactions are complex. However, we can broadly outline the effects of changing the parameters in the way described. In the coda model, duration differences between the obstruents arise because of the smaller variance of the voiced obstruent duration constraint. The interaction of this constraint with the targeted rhyme duration constraint gives rise to the complementary difference in preceding vowel duration. In the onset model, the obstruent duration constraints remain the same, causing lengthening of the voiced obstruent to be more costly (reduce probability more), than lengthening of the voiceless obstruent. However, onset duration is not relevant to the target rhyme constraint in the CV case, so there is no interaction. The only pressure to lengthen the consonants comes from the C/V constraint. Therefore, consonants only lengthen in order to achieve the best possible C/V ratio. The consonant durations, however, do not differ much from the previous models. It is the vowel duration that is most strongly affected by the re-weighting of the C/V constraint. The articulatory vowel faces pressure to lengthen, but undershoots the target more than in the original VC model, due to the greater influence of the C/V constraint.



**Figure C2**  
*Competing Constraints Model*

The general behavior of the onset model looks very similar to the coda model with the inherently short vowel. However, the mechanism is quite different; the relevant variable is a ratio ( $C/V$ ), rather than a sum ( $Rhyme = V + C$ ). Therefore, the relationship between vowel and consonant duration is not negatively correlated, but positively correlated: the articulatory vowel is *shorter* because the tautosyllabic consonant is shorter. Vowels following voiced consonants are therefore shorter than vowels following voiceless consonants (CV in Fig. C2). However, the in-phase timing between onset and nucleus means that the articulatory vowel will be masked to a greater degree by the longer (voiceless) consonant. In this case the effect of masking is slightly larger than the  $C/V$  effect. Therefore, the net result is a slightly longer post-voiced than post-voiceless acoustic vowel. The voicing effect is smaller, but it shows the same dependence on total duration as the other models. These outcomes are consistent with the literature summarized in Section 1.6.

## Appendix D

### Corpus Models

1. **Tanner et al. Simplified:** Vowel Duration ~ Vowel Height + Vowel Type + Voicing + Speaking Rate + Obstruent Type + Frequency + Voicing:(V Height + Speaking Rate Deviation + Obstruent Type + Frequency + V Type) + (Voicing + Obstruent Type + Frequency | Speaker) + (Speaking Rate Deviation | Word) + (1 | Vowel Quality)
2. **No-Interaction Model:** Vowel Duration ~ Vowel Height + Vowel Type + Speaking Rate + Voicing + Obstruent Type + Frequency + (Obstruent Type + Frequency | Speaker) + (Speaking Rate | Word)
3. **Model 1 with consonant duration added:** Vowel Duration ~ Vowel Height + Vowel Type + Voicing + Speaking Rate + Consonant Duration + Obstruent Type + Frequency + Voicing:(Speaking Rate + Obstruent Type + Frequency + Consonant Duration) + (Voicing + Obstruent Type + Frequency + Consonant Duration | Speaker) + (1 | Word) + (1 | Vowel Quality)
4. **Model 2 with consonant duration added:** Vowel Duration ~ Vowel Height + Vowel Type + Speaking Rate + Consonant Duration + Voicing + Obstruent Type + Frequency + (Obstruent Type + Frequency + Consonant Duration | Speaker) + (Speaking Rate | Word)
5. **Medial tokens; consonant duration included; voicing excluded:** Vowel Duration ~ Vowel Height + Vowel Type + Speaking Rate + Consonant Duration + Voicing + Obstruent Type + Frequency + (Obstruent Type + Frequency + Consonant Duration | Speaker) + (1 | Word)
6. **Consonant duration model:** Consonant Duration ~ Voicing + Speaking Rate + Obstruent Type + Frequency + Voicing:(Speaking Rate + Frequency + Obstruent Type) + (Speaking Rate + Frequency + Obstruent Type | Speaker) + (Speaking Rate + Frequency | Word)