



**Predictability of stop consonant phonetics across talkers:  
 Between-category and within-category dependencies  
 among cues for place and voice**

Journal:	<i>Linguistics Vanguard</i>
Manuscript ID	LingVan.2017.0047
Manuscript Type:	research-article
Classifications:	Phonetics & Phonology
Keywords:	stop consonants, talker variability, phonetic covariation, factor analysis, predictability
Abstract:	<p>The present study investigated patterns of covariation among acoustic properties of stop consonants in a large multi-talker corpus of American English connected speech. Relations among talker means for different stops on the same dimension (between-category covariation) were considerably stronger than those for different dimensions of the same stop (within-category covariation). The existence of between-category covariation supports a uniformity principle that restricts the mapping from phonological features to phonetic targets in the sound system of each speaker. This principle was formalized with factor analysis, in which observed covariation derives from a lower-dimensional space of talker variation. Knowledge of between-category phonetic covariation could facilitate perceptual adaptation to novel talkers by providing a rational basis for generalizing idiosyncratic properties to several sounds on the basis of limited exposure.</p>

SCHOLARONE™  
 Manuscripts

## Predictability of stop consonant phonetics across talkers:

Between-category and within-category dependencies among cues for place and voice

**Abstract**

The present study investigated patterns of covariation among acoustic properties of stop consonants in a large multi-talker corpus of American English connected speech. Relations among talker means for different stops on the same dimension (between-category covariation) were considerably stronger than those for different dimensions of the same stop (within-category covariation). The existence of between-category covariation supports a uniformity principle that restricts the mapping from phonological features to phonetic targets in the sound system of each speaker. This principle was formalized with factor analysis, in which observed covariation derives from a lower-dimensional space of talker variation. Knowledge of between-category phonetic covariation could facilitate perceptual adaptation to novel talkers by providing a rational basis for generalizing idiosyncratic properties to several sounds on the basis of limited exposure.

**Keywords:** stop consonants; talker variability; phonetic covariation; factor analysis; predictability

## 1 Introduction

The phonetic realization of an individual sound category can vary substantially according to contextual, lexical, dialectal, and talker-specific influences. This variation is highly structured: previous research has documented strong dependencies among phonetic properties, as well as between phonetic properties and many sociolinguistic factors (e.g., Labov, 1966; Foulkes et al., 2001; Foulkes & Docherty, 2006; Guy & Hinskens, 2016; Fruehwald, 2017; Sonderegger et al., 2017). The present study focused on two prominent types of linear dependency in phonetic variation. The first type of dependency holds among multiple categories along a single phonetic dimension ('between-category' covariation); the second holds among multiple phonetic dimensions within individual categories ('within-category' covariation).<sup>1</sup>

Instances of between-category phonetic dependencies have been observed among several speech sounds. Talker-specific vowel systems differ extensively in the log F1 x F2 formant plane, but the systems are highly parallel, suggesting covariation along these dimensions (e.g., Joos, 1948; Nearey, 1978; Nearey & Assmann, 2007). Furthermore, largely constant spectral and temporal ratios are preserved among vowel categories across speaking rates and styles (Smiljanic & Bradlow, 2008; DiCanio et al., 2015). Relations among vowels can also be preserved during diachronic sound change, as when multiple vowels undergo parallel shifts in their phonetic realization (e.g., Fruehwald, 2013, 2017). Among fricatives, the spectral center of gravity of [s] and [ʃ] vary substantially across talkers, yet within a talker, the mean COG of [s] is systematically higher than the corresponding mean COG of [ʃ] (Newman et al., 2001). Strong covariation of mean voice onset time (VOT) has also been observed among stop consonants

---

<sup>1</sup> There may be other forms of statistical dependency, beyond linear relations, among phonetic variables. It is also conceivable that the values of one category on a given dimension could covary with those of another category on a different dimension. We considered this alternative notion of 'between-category' covariation but found limited evidence for it in the present data.

1  
2  
3 across speakers of the same language (e.g., Zlatin, 1974; Koenig, 2000; Newman, 2003; Solé,  
4  
5 2007; Theodore et al., 2009; Chodroff & Wilson, 2017). Theodore et al. (2009) observed a  
6  
7 similar difference between the mean VOT values of [p<sup>h</sup>] and [k<sup>h</sup>] across talkers. Chodroff &  
8  
9 Wilson (2017) extended the study of between-category VOT covariation to all word-initial stop  
10  
11 categories of American English (AE), in both isolated and connected speech, while controlling  
12  
13 for many other sources of VOT variation (e.g., utterance position, following vowel, lexical  
14  
15 properties). Correlations of talker VOT means were particularly strong among the voiceless  
16  
17 aspirated stops, and moderate among the voiced stops and homorganic voiced pairs.  
18  
19

20  
21  
22 There is substantially less evidence for within-category phonetic dependencies across  
23  
24 tokens and across talker means. For example, while vowel height (as indexed by F1) and vowel  
25  
26 duration are known to covary across vowels in many languages (e.g., Lindblom, 1967;  
27  
28 Maddieson, 1997), this relation does not appear to hold across individual tokens of the same  
29  
30 vowel category (Toivonen et al., 2015).<sup>2</sup> Several studies have also examined the possibility that  
31  
32 different cues to the voicing contrast, such as VOT and fundamental frequency (f0) at following  
33  
34 vowel onset, covary within stop categories (e.g., Shultz et al., 2012; Dmitrieva et al., 2015; Kirby  
35  
36 & Ladd, 2015, 2016; Clayards, 2018). Positive correlations of the relevant cues would indicate  
37  
38 enhancement of the contrast, whereas negative correlations would suggest cue-trading or  
39  
40 compensation relations. These relationships could hold across tokens or across talker means, and  
41  
42 would indicate systematic relations in the *use* of phonetic dimensions. The observed within-  
43  
44 category dependencies tend to be weak and vary considerably by language and sample. Kirby &  
45  
46 Ladd (2015) found a significant negative correlation between VOT and f0 across tokens of word-  
47  
48  
49  
50  
51  
52

---

53  
54 <sup>2</sup> Strong pairwise correlations of talker mean log f0, F1, F2, and F3 have been observed when  
55  
56 aggregated over all vowels (Nearey, 1989; Assmann et al., 2008; see also Rose, 2010 for F2 and  
57  
58 F3, and Whalen & Levitt, 1995 for f0 and F1). It remains unclear, however, which particular  
59  
60 vowel categories exhibit these dependencies most strongly.

1  
2  
3 initial [p] in Italian, but this correlation did not reach significance in French. A weak negative  
4  
5 correlation between VOT and  $f_0$  was observed across tokens of AE [p<sup>h</sup>] by Dmitrieva et al.  
6  
7 (2015), but another study of the same language yielded no significant linear relation between  
8  
9 those dimensions for [p<sup>h</sup>] or [b] across tokens or talker means (Clayards, 2018).  
10  
11

12  
13 Some factors involved in phonetic realization may induce both between- and within-  
14  
15 category covariation. For example, Koenig (2000) found that median VOT and following vowel  
16  
17 duration were highly correlated across talkers for both [p<sup>h</sup>] ( $r = 0.72$ ) and [t<sup>h</sup>] ( $r = 0.77$ ). This  
18  
19 likely reflects talker specificity in global speaking rate, leading to the expectation that  
20  
21 correlations would also be found between the stop categories and for other duration-sensitive  
22  
23 cues. Covariation among cues that occurs both between and within categories may be reducible  
24  
25 to global factors such as speaking rate or airflow rate; this point will be considered further in the  
26  
27 discussion.  
28  
29  
30

31  
32 We examined covariation within and among the six AE word-initial stop consonants,  
33  
34 focusing on three well-known cues to the place and voice contrasts: spectral center of gravity  
35  
36 (COG; e.g., Forrest et al., 1988), positive VOT (e.g., Lisker & Abramson, 1964), and  $f_0$  at vowel  
37  
38 onset (e.g., Haggard et al., 1970; Ohde, 1984). Correlation analyses (section 2) revealed  
39  
40 considerably stronger between-category covariation than within-category covariation. This  
41  
42 accords with previous findings but is comprehensively demonstrated for the first time here, with  
43  
44 all measurements performed on the same multi-talker data set. The between-category  
45  
46 correlations can be accounted for with a principle of *uniformity* that constrains the mapping from  
47  
48 phonological feature values to talker-specific phonetic targets (Chodroff & Wilson, 2017;  
49  
50 Chodroff, 2017). This principle was quantitatively formalized and evaluated against the observed  
51  
52 correlations within the dimensionality-reduction framework of factor analysis (section 3). As we  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 discuss in section 4, covariation in phonetic realization across speakers implies predictability for  
4  
5 listeners: listeners could use phonetic dependencies among stop categories to generalize from  
6  
7 limited experience with a novel talker.  
8  
9

## 10 11 **2 Correlation analysis**

### 12 13 **2.1 Methods**

14  
15 The data was extracted from an audited subset of the Mixer 6 corpus (Brandschain et al.,  
16  
17 2010, 2013; Chodroff et al., 2016) containing approximately 45 minutes of read speech from 180  
18  
19 native AE speakers (102 female). Transcripts were aligned to the corresponding WAV files with  
20  
21 the Penn Forced Aligner (Yuan & Liberman, 2008), and all word-initial prevocalic stop  
22  
23 consonants were further processed with AutoVOT (Keshet et al., 2014). AutoVOT automatically  
24  
25 identifies the stop release and following vowel onset within a user-specified window of analysis.  
26  
27 Further details about the talkers, read sentences, and boundary alignments can be found in  
28  
29 Chodroff & Wilson (2017).<sup>3</sup>  
30  
31  
32  
33  
34

35  
36 COG, positive VOT, and onset  $f_0$  in the following vowel were measured for each stop.  
37  
38 COG was calculated from a smoothed spectrum over the initial portion of the release burst. Each  
39  
40 spectrum was computed by averaging FFTs from seven consecutive 3 ms windows, with the first  
41  
42 window centered on the burst transient and a window shift of 1 ms (Hanson & Stevens, 2003;  
43  
44 Flemming, 2007; Chodroff & Wilson, 2014). Positive VOT was defined as the duration from  
45  
46 stop release to the onset of periodicity in the vowel; this was automatically extracted from the  
47  
48 AutoVOT boundaries or from manually-corrected boundaries when available. The  $f_0$  value was  
49  
50  
51

---

52  
53 <sup>3</sup> The data set here was somewhat larger than that analyzed in Chodroff & Wilson (2017), which  
54  
55 included only stops at the beginning of *stressed* word-initial syllables. Because all talkers  
56  
57 recorded the same set of sentences, any effects of stress (or other contextual factors) on VOT and  
58  
59 other acoustic properties should be approximately consistent. We aimed to include as many  
60  
61 tokens as possible to maximize the power available to identify phonetic dependencies.

1  
2  
3 the first one measured by Praat (Boersma & Weenink, 2016) within 50 ms after the following  
4  
5 vowel onset.  
6  
7

8 For each stop category and cue separately, values 2.5 standard deviations above or below  
9  
10 the talker-specific mean were excluded. Consequently, the total number of valid tokens varied  
11  
12 somewhat by stop category and measurement (COG: 87,968 tokens; VOT: 96,357; onset f0:  
13  
14 76,144). Table 1 summarizes the data available per talker for each stop and cue combination.  
15  
16

17  
18 ===== TABLE 1 =====  
19

## 20 21 **2.2 Results**

22 For each stop and cue combination, talker means were calculated from all available  
23  
24 tokens (e.g., a talker's mean COG for [p<sup>h</sup>] was calculated from all of his or her productions of  
25  
26 that stop with non-outlier COG values). Pearson correlations were performed on the talker means  
27  
28 between stop categories along a single dimension and, separately, between dimensions within  
29  
30 each stop category. The bias-corrected and accelerated percentile (BCa) method was used to  
31  
32 form 95% bootstrapped confidence intervals for the correlations (1000 replicates; Efron, 1987).  
33  
34 As shown in Table 2, the between-category correlations among stops were positive and generally  
35  
36 high.<sup>4</sup> For the COG dimension, moderate correlations were observed among the voiceless stops,  
37  
38 and strong correlations were observed among the voiced stops and between homorganic stop  
39  
40 pairs (e.g., [k<sup>h</sup>]-[g]). For VOT, the pattern of correlations replicated that found in Chodroff &  
41  
42 Wilson (2017): relations were very strong among the voiceless stops, and moderate to weak  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55

---

56  
57 <sup>4</sup> The significance level was adjusted for multiple comparisons to the conservative value of  $\alpha =$   
58 0.001.  
59  
60

among the voiced stops and homorganic pairs. Finally, talker mean  $f_0$  was almost perfectly correlated for all stops.<sup>5</sup>

===== TABLE 2 =====

Within-category correlations of talker means are shown in Table 3. Note that for  $f_0$ , the correlations were calculated separately for male and female talkers. COG and VOT means showed significant positive correlations within each of the stop categories, but the strength of the relation depended on the category. These two cues were weakly correlated within the voiceless stops and [b], but strongly correlated within [d] and moderately within [g]. Within-category correlations between COG and  $f_0$ , and VOT and  $f_0$ , were numerically quite weak and none reached significance. Additional within-category correlations calculated over tokens (instead of over talker means) are reported in the Appendix.

===== TABLE 3 =====

### 2.3 Discussion

We found strong between-category covariation on each dimension for the segments and phonetic cues investigated here. In part, this surely reflects anatomical differences among talkers: differences that affect the articulation and resulting acoustics of many sounds (e.g., the strong dependencies in onset  $f_0$  are partly due to cross-talker variation in vocal fold length and tissue density; Titze, 2011). But anatomy does not wholly determine phonetic realization. Each talker could in principle have shown greater between-category differences in the phonetic targets that are indexed by COG (e.g., by specifying tongue tip position and contact width for coronal [t] differently than for [d]), VOT (e.g., by planning the duration or timing of glottal spreading for

---

<sup>5</sup> Correlations are described using modifiers based on recommendations in Evans (1996): a ‘strong’ correlation is one above 0.59, a ‘moderate’ correlation is between 0.40 and 0.59, and a ‘weak’ correlation is below 0.40.



1  
2  
3 [p<sup>h</sup>] differently than for [k<sup>h</sup>]), and even f<sub>0</sub> (e.g., by having a different pitch target for vowels  
4 following [d] than for those after [g]). Indeed, research on language- and dialect- specific  
5 phonetics (e.g., Lisker & Abramson, 1964), as well as the dual phonetic systems of bilinguals  
6 (e.g., Flege, 1991; Grosjean & Miller, 1994; MacLeod & Stoel-Gammon, 2005; Chang et al.,  
7 2011), has demonstrated that the phonetic targets associated with any given category, such as  
8 [p<sup>h</sup>] or [b], can differ in ways that anatomy alone could never explain. Some additional principle  
9 must restrict the variation in phonetic targets for a given individual when speaking a given  
10 language.

11  
12 One version of the principle would require speakers to have the same (or highly similar)  
13 *patterns* of phonetic targets for sounds that bear a given phonological feature value. For example,  
14 the principle could require each talker's laryngeal target for [p<sup>h</sup>] to be systematically related to  
15 the same talker's laryngeal target for [k<sup>h</sup>], on the grounds that both sounds are specified [-voice]  
16 (or [+spread glottis]). A more stringent version of the principle requires the talker-specific  
17 phonetic target corresponding to a given phonological feature value to be identical or *uniform*  
18 across all sounds bearing that specification.<sup>6</sup> In this version, a talker cannot independently  
19 specify the properties of the laryngeal spreading gesture and associated timing relations for [p<sup>h</sup>],  
20 [t<sup>h</sup>], and [k<sup>h</sup>]: the mapping from phonological feature to phonetic target must be the same, in  
21 these respects, for all three voiceless aspirated stop categories. (Observed variation in cue values  
22 for identically specified sounds must then arise 'automatically' from independent differences in  
23 other targets, as discussed in section 4 below.)

24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53 <sup>6</sup> Either version of the principle must apply separately to each language / dialect / register that is  
54 controlled by a speaker. More generally, the principle should allow wide latitude in the  
55 contextual factors that can affect phonetic realization (e.g., prosody and speaking rate alongside  
56 language and register), requiring only that such factors have uniform effects on all sounds that  
57 are identically specified with respect to the relevant phonological feature.  
58  
59  
60

1  
2  
3 Within-category covariation was generally quite weak and we do not propose a separate  
4 principle of phonetic implementation to explain the dependencies that were found. The positive  
5 COG-VOT correlation observed for several stops may be attributable to aerodynamics: a higher  
6 airflow rate may give rise to an upwards shift in the energy distribution across the spectrum  
7 (Zue, 1976; Koenig et al., 2013; Chodroff & Wilson, 2014) and to longer aspiration durations.  
8 Notably, we found that this dependency held not only within stop categories but also when  
9 comparing the same two cues across them (e.g., for mean COG of [d] and mean VOT of [p<sup>h</sup>],  $r =$   
10 0.46,  $p < 0.001$ ). Perhaps each talker has a relatively higher or lower airflow rate across all stops  
11 (plausibly due to uniform realization of the shared feature [–continuant]), and this affects all  
12 COG and VOT values accordingly. A remaining question is why the correlations are particularly  
13 high within [d] and [g]. We speculate that the presence of voicing during stop closure for some  
14 talkers may lower COG and VOT. The weaker within-category correlation for [b] is likely due to  
15 the fact that the (positive) VOT of this stop does not vary as much across talkers (Chodroff &  
16 Wilson, 2017).  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

### 37 **3 Factor analysis**

38  
39 The strong between-category correlations documented above indicate that population-  
40 level variation in the phonetic realization of AE stops can be accurately modeled with a  
41 relatively small number of *latent values* for each talker. This idea could be formalized with many  
42 dimensionality reduction methods, including traditional principle component analysis (PCA; e.g.,  
43 Pols et al., 1973; van Nierop et al., 1973) and more recently proposed eigenvoice decompositions  
44 (Kuhn et al., 1998). The model developed in this section is an instance of *factor analysis* (FA;  
45 e.g., Harshman et al., 1977; Clopper & Paolillo, 2006; Leinonen, 2008), a formally simple  
46 method that can express easily interpretable hypotheses about the content of the latent values.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

In FA generally, each *observed* vector ( $\mathbf{x}_i$ ) is modeled as drawn from a multivariate normal distribution with mean  $\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Psi}$ . The *factor loading matrix*  $\mathbf{W}$  represents a linear map from a latent vector ( $\mathbf{z}_i$ ) into the observation space, where the dimensionality of  $\mathbf{z}_i$  is smaller than that of  $\mathbf{x}_i$ . The offset vector  $\boldsymbol{\mu}$  represents aspects of the mean that are, according to the model, the same across all individuals. Two additional restrictions are enforced: (a)  $\boldsymbol{\Psi}$  is required to be diagonal, so that the components of an observed vector  $\mathbf{x}_i$  are independent conditional on the latent vector  $\mathbf{z}_i$  and (b) the distribution over latent vectors is a multivariate normal with zero mean and unit covariance, so that the components of  $\mathbf{z}_i$  are standardized and independent from one another. In summary, for each individual  $i = 1, \dots, n$

$$p(\mathbf{x}_i) = N(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}), \text{ where } \boldsymbol{\Psi} \text{ is diagonal}$$

$$p(\mathbf{z}_i) = N(\mathbf{0}, \mathbf{I})$$

It follows that two components of the observed vector are predicted to covary only if they ‘load on’ (have non-zero influence from) one or more common latent variables. In this way, factor analysis represents correlated variables in a higher-dimensional space with uncorrelated variables in a lower-dimensional space through a simple (i.e., linear) transformation.

FA provides a method for formalizing the hypothesis that, within the set of AE stops, the talker-specific contribution to phonetic targets is uniform for each phonological feature specification (as suggested in section 2). To evaluate this hypothesis, we take each observed  $\mathbf{x}_i$  to be the vector of stop cue means for a particular talker (i.e.,  $\mathbf{x}_i = [p_i^{\text{COG}}, b_i^{\text{COG}}, t_i^{\text{COG}}, \dots, p_i^{\text{VOT}}, \dots, p_i^{\text{Onset-f0}}, \dots, g_i^{\text{Onset-f0}}]^T$ , where  $x_i^c$  is the sample mean for cue  $c$  of stop  $x$  computed from the productions of the  $i$ th talker). The latent vector  $\mathbf{z}_i$  represents the talker-specific contributions to the phonetic targets for all six stops, as reflected in the acoustic cues. Under the idealization that each cue reflects the talker-specific target contribution for one phonological feature, the

uniformity principle implies a factor loading matrix  $\mathbf{W}$  with the following highly sparse structure (where zero entries are left blank and the matrix is transposed for display purposes):

$$\mathbf{W}^T = \begin{matrix} \begin{matrix} lab \\ cor \\ dor \\ vcl \\ vcd \\ pitch \end{matrix} & \begin{pmatrix} p^{COG} & b^{COG} & t^{COG} & d^{COG} & k^{COG} & g^{COG} & p^{VOT} & b^{VOT} & t^{VOT} & d^{VOT} & k^{VOT} & g^{VOT} & p^{f0} & b^{f0} & t^{f0} & d^{f0} & k^{f0} & g^{f0} \\ w_{lab} & w_{lab} & & & & & & & & & & & & & & & & & \\ & & w_{cor} & w_{cor} & & & & & & & & & & & & & & & \\ & & & & w_{dor} & w_{dor} & & & & & & & & & & & & & \\ & & & & & & w_{vcl} & & w_{vcl} & & w_{vcl} & & w_{vcl} & & & & & & \\ & & & & & & & w_{vcd} & & w_{vcd} & & w_{vcd} & & w_{vcd} & & & & & \\ & & & & & & & & & & & & & & & & & & & w_p & w_p & w_p & w_p & w_p & w_p \end{pmatrix} \end{matrix}$$

For example, according to  $\mathbf{W}$  the latent factor identified with the feature *vcl* (i.e., [-voice] or [+spread glottis]) has the same influence ( $w_{vcl}$ ) on the VOT of the three voiceless stops, and zero effect on all other stop-cue combinations. Similar logic applies to the other feature-cue combinations (e.g., given the idealization that COG reflects only place features); note that ‘pitch’ is an ad-hoc feature that predicts covariation of talker-specific  $f_0$  means across all vowels. In our implementation of the model, we in fact multiplied each row of  $\mathbf{W}$  by the empirical standard deviation of its label (e.g., the first row was multiplied by the standard deviation of the talker means for COG of [p<sup>h</sup>]). This induces a marginal distribution over  $\mathbf{x}_i$  in which each  $w_k^2$  is interpretable as a positive and pooled correlation coefficient.

The free parameters of this version of the FA model (i.e., the coefficients of  $\mathbf{W}$ , the offset  $\boldsymbol{\mu}$ , and the diagonal elements of  $\boldsymbol{\Psi}$ ) were fit to the talker mean vectors measured earlier. As expected, the factor-loading coefficients indicated strong correlations among the stop-cue combinations that reflect a common feature value ( $w_{lab}^2 = .60$ ,  $w_{cor}^2 = .67$ ,  $w_{dor}^2 = .69$ ,  $w_{vcl}^2 = .72$ ,  $w_{vcd}^2 = .20$ ,  $w_{pitch}^2 = .91$ ; these values should be compared to the correlations in Table 2). The values in the offset  $\boldsymbol{\mu}$  account for talker-general differences in the values of stop-cue combinations that are otherwise unexpected given  $\mathbf{W}$ ; for example, the offset for the VOT of [b] ( $\mu = 8.42$ ) is lower than that for the VOT of [g] ( $\mu = 16.86$ ). One interpretation of such offset differences is that they reflect ‘automatic’ articulatory and acoustic effects—influences on the

1  
2  
3 measurements that would be present even if the underlying phonetic targets studied here were  
4  
5 exactly uniform within a talker. For example, effects of place on stop closure duration could  
6  
7 contribute to differences in VOT values even if laryngeal targets and their timing with respect to  
8  
9 supralaryngeal gestures are uniform (e.g., Weismer, 1980; Maddieson, 1997).  
10  
11

12  
13 We compared the FA model above (the *target uniformity* model) with several alternatives  
14  
15 that differed in the factor loading matrix: a *null covariation* model, in which  $\mathbf{W}$  was the diagonal  
16  
17 matrix; a sample of 500 *row permutation* models derived by randomly exchanging the rows of  
18  
19 the target uniformity model; and an *exploratory* model in which the factor loading matrix had six  
20  
21 columns with all cell values were fit to the data. Table 4 reports the marginal negative log-  
22  
23 likelihood values of the talker mean vectors for each model. The target uniformity model  
24  
25 provided a significant improvement over the null model and all of the row-permutation variants,  
26  
27 while the exploratory model was superior to target uniformity (similar results were obtained with  
28  
29 other model comparison measures and with cross-validation). These results suggest that target  
30  
31 uniformity is an important (but unlikely the only) principle of phonetic implementation that  
32  
33 constrains the covariation of stop consonants within talkers.<sup>7</sup>  
34  
35  
36  
37  
38

39 ===== TABLE 4 =====  
40  
41  
42  
43  
44

45  
46 <sup>7</sup> Two columns of the loading matrix in the exploratory model were quite similar to columns in  
47  
48 the theoretically-determined matrix  $\mathbf{W}$ . The first had large values only for the VOT means of the  
49  
50 voiceless stops, closely emulating uniform realization of [-voice] or [+spread glottis]. The  
51  
52 second had values near unity for all of the f0 means, and much smaller values elsewhere, in line  
53  
54 with uniformity with respect to the ‘pitch’ feature. A third column seemed to combine the place  
55  
56 effects of  $\mathbf{W}$ , with particularly large values for the COG mean of [k<sup>h</sup>] and [g], intermediate  
57  
58 values for the other COG means, and smaller values for all other means. Two of the remaining  
59  
60 columns appear to express correlations among COG and VOT for [d] and [g] separately; these  
within-category relations were found in our statistical analysis (see section 2) but do not follow  
from the uniformity principle. The final column was generally difficult to interpret but assigned a  
particularly large value to the VOT mean of [b], the stop most likely to have closure voicing.

#### 4 Covariation and predictability in perceptual adaptation

Patterns of phonetic covariation such as the one observed above have important implications not only for the mapping between phonology and phonetics, but also for their potential role in perceptual adaptation. On the basis of strong between-category covariation, listeners could reasonably predict a talker-specific target for one sound category having heard productions only of one or more covarying categories. The findings of many studies of perceptual generalization are consistent with this idea. For instance, listeners generalize talker-specific spectral characteristics from exposure vowels to previously unheard vowels (e.g., Ladefoged & Broadbent, 1957; Maye et al., 2008; Chládková et al., 2017) and have been shown to extrapolate a talker's characteristic VOT from [p<sup>h</sup>] to [k<sup>h</sup>] based on direct evidence about [p<sup>h</sup>] only (e.g., Theodore & Miller, 2010; Nielsen 2011). A listener may have low prior expectations for within-category covariation, but could infer talker-specific relations among cues through distributional learning (e.g., Clayards et al., 2008).

The FA model presented above encodes covariation with a lower-dimensional set of latent variables. If listeners attempted to infer the latent factor values of a novel talker, they would generalize across segments with shared phonological feature specifications in a way that is consistent with the population-level correlations. In this sense, the FA model could be interpreted as a cognitive model of adaptation. Between-category covariation has been incorporated to varying degrees in previous models of adaptation. Models that employ variants of mean subtraction or *z*-scoring within each phonetic dimension implicitly enforce covariation, provided the calculation incorporates values from several speech sounds (e.g., Sliding Template Model of Vowel Normalization: Nearey & Assmann, 2007; c-CuRE: McMurray & Jongman, 2011; VOT generalization: Nielsen & Wilson, 2008). These models, however, have assumed that

1  
2  
3 each talker has a single ‘offset’ value per cue and thus assumed *perfect* covariation among *all*  
4  
5 speech sounds represented on a given acoustic-phonetic dimension. To the extent that the  
6  
7 empirical covariation is weak for some of the relevant sounds (e.g., as observed for VOT  
8  
9 between stops contrasting in voice), this strong assumption could lead to suboptimal  
10  
11 performance. Alternative models of talker adaptation, such as exemplar models or the ideal  
12  
13 adaptor model, do not currently encode covariation of phonetic properties, and may therefore fail  
14  
15 to model aspects of perceptual generalization across speech sounds (e.g., Johnson, 1997;  
16  
17 Kleinschmidt & Jaeger, 2015). In comparison, FA can model selective, and ideally theoretically-  
18  
19 grounded, patterns of between-category covariation, as opposed to assuming perfect covariation  
20  
21 or omitting covariation altogether. Future research should be directed towards understanding the  
22  
23 relation between measured phonetic covariation and patterns of perceptual generalization by  
24  
25 human listeners.  
26  
27  
28  
29  
30  
31

## 32 **5 Conclusion**

33  
34  
35 The analyses of talker means for COG, VOT, and onset  $f_0$  within and among stop  
36  
37 categories revealed greater between-category covariation in comparison to within-category  
38  
39 covariation. As examined in section 3, the observed covariation among phonetic categories may  
40  
41 arise from a constraint of uniformity on the mapping from phonological features to phonetic  
42  
43 targets underlying acoustic-phonetic properties. Further research is required to evaluate the  
44  
45 predictions of uniformity as it applies to other segments and languages (Chodroff, 2017). In  
46  
47 addition, perceptual knowledge of covariation could facilitate prediction in perceptual  
48  
49 processing, and more generally, the measured covariation serves as a testable hypothesis of  
50  
51 perceptual knowledge in generalized adaptation in speech perception.  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Appendix

The analyses in section 2 involved correlations of talker means; however, many previous studies have also examined correlations across individual tokens (e.g., Dmitrieva et al., 2015; Kirby & Ladd, 2015, 2016; Clayards, 2018). For comparison with these studies, token-by-token correlations between phonetic cues were calculated for each stop category within and across talkers. Only stop consonants with non-outlier values for both cues were retained for these correlations. There were 71,852 stops for the COG-VOT analysis, 57,737 stops for the COG-f<sub>0</sub> analysis, and 74,916 stops for the VOT-f<sub>0</sub> analysis. The first correlation analysis, reported in Table A1, was conducted across all tokens (see also Dmitrieva et al., 2015; Clayards, 2018). These correlations largely resembled the correlations of talker means in magnitude (especially between COG and VOT for [b], [d], and [g]); while many of these correlations reached significance, they were nevertheless quite weak. In the second analysis, correlations were limited to talkers with more than 20 tokens per stop category. The median number of talkers excluded from each analysis was four and the maximum was 55 talkers (between COG and f<sub>0</sub> for [t<sup>h</sup>]). Table A2 presents the median token-by-token correlation for each of the cue pairs and stop consonants, as well as the range across talkers. Consistent with findings in Kirby & Ladd (2016) for French and Italian intervocalic stops, the magnitude and direction of the by-speaker correlations varied substantially across talkers. Together, these findings indicate that while there may exist weak relationships across talker means, the token-by-token relationships within talker-specific productions are highly variable.

===== TABLE A1 =====  
===== TABLE A2 =====



Table 1. For each stop category and measurement separately, the median number of tokens per talker (left column) and the range of tokens per talker (right column).

	COG		VOT		f0	
	Median	Range	Median	Range	Median	Range
p <sup>h</sup>	75	45 – 98	77	44 – 100	59.5	6 – 96
b	86	57 – 127	98	64 – 138	80	9 – 127
t <sup>h</sup>	45	16 – 75	46	17 – 77	37	4 – 67
d	115	53 – 170	140	64 – 192	113	12 – 173
k <sup>h</sup>	91	48 – 112	93	50 – 114	77	4 – 110
g	82	52 – 116	91	54 – 122	78	5 – 111

Table 2. Pearson correlation coefficients and 95% BCa bootstrap confidence intervals of stop-specific talker means for COG, VOT, and f0. All *p*-values were less than 0.001 unless otherwise indicated.

	COG		VOT		f0	
p <sup>h</sup> -b	0.60	[0.49, 0.69]	0.17, <i>n.s.</i>	[0.00, 0.32]	0.98	[0.97, 0.98]
t <sup>h</sup> -d	0.66	[0.57, 0.73]	0.53	[0.43, 0.63]	0.97	[0.94, 0.98]
k <sup>h</sup> -g	0.69	[0.58, 0.76]	0.43	[0.32, 0.52]	0.97	[0.94, 0.98]
p <sup>h</sup> -t <sup>h</sup>	0.40	[0.26, 0.51]	0.83	[0.77, 0.88]	0.98	[0.97, 0.99]
t <sup>h</sup> -k <sup>h</sup>	0.47	[0.34, 0.58]	0.79	[0.73, 0.83]	0.98	[0.95, 0.99]
k <sup>h</sup> -p <sup>h</sup>	0.57	[0.45, 0.65]	0.83	[0.78, 0.87]	0.98	[0.98, 0.99]
b-d	0.55	[0.40, 0.66]	0.11, <i>n.s.</i>	[-0.03, 0.24]	0.98	[0.98, 0.99]
d-g	0.67	[0.58, 0.75]	0.37	[0.24, 0.50]	0.98	[0.94, 0.99]
g-b	0.63	[0.51, 0.72]	0.48	[0.37, 0.58]	0.98	[0.95, 0.98]

Table 3. Pearson correlation coefficients and 95% BCa bootstrap confidence intervals of stop-specific talker means between COG, VOT, and f0. For f0, the correlations are reported separately for female and male talkers. Correlations with *p*-values less than 0.001 are identified with an asterisk.

	COG-VOT		COG-f0 (female)		COG-f0 (male)		VOT-f0 (female)		VOT-f0 (male)	
p <sup>h</sup>	0.32*	[0.17, 0.44]	0.11	[-0.08, 0.29]	0.12	[-0.09, 0.37]	0.09	[-0.09, 0.28]	-0.03	[-0.22, 0.20]
b	0.37*	[0.25, 0.50]	0.03	[-0.17, 0.22]	-0.16	[-0.34, 0.05]	-0.11	[-0.33, 0.12]	-0.09	[-0.30, 0.13]
t <sup>h</sup>	0.22	[0.07, 0.38]	0.10	[-0.08, 0.30]	0.11	[-0.13, 0.34]	0.03	[-0.19, 0.23]	-0.02	[-0.25, 0.21]
d	0.74*	[0.66, 0.80]	0.03	[-0.17, 0.24]	-0.19	[-0.39, 0.02]	0.00	[-0.20, 0.21]	0.00	[-0.24, 0.23]
k <sup>h</sup>	0.23	[0.09, 0.35]	0.07	[-0.12, 0.25]	0.05	[-0.18, 0.28]	0.17	[-0.02, 0.36]	-0.08	[-0.30, 0.15]
g	0.54*	[0.43, 0.63]	0.12	[-0.09, 0.27]	-0.14	[-0.34, 0.09]	0.08	[-0.13, 0.26]	-0.06	[-0.26, 0.17]

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Table 4. Marginal negative log-likelihood for observed talker mean vectors under different versions of the factor analysis model.

Model	Negative log-likelihood
target uniformity	14836.99
null (diagonal)	16694.72
row permutation	16438.98 (range: 15696.90-16583.83)
exploratory	14419.77

For Review Only

Table A1. Token-by-token correlations for each cue pair and stop category aggregated over all talkers. An asterisk reflects  $p < 0.001$ .

	COG-VOT	COG-f0 (female)	COG-f0 (male)	VOT-f0 (female)	VOT-f0 (male)
p <sup>h</sup>	0.18*	-0.01	0.09*	-0.05*	-0.02
b	0.34*	0.00	-0.05*	-0.03	-0.01
t <sup>h</sup>	0.09*	0.07*	0.12*	-0.11*	-0.06*
d	0.57*	-0.11*	-0.06*	-0.06*	-0.01
k <sup>h</sup>	0.17*	0.10*	0.09*	-0.15*	-0.13*
g	0.52*	0.03	-0.01	-0.03	0.01

Table A2. For each stop category and cue pair separately, the median talker-specific token-by-token correlation (left column) and range of talker-specific token-by-token correlations (right column).

	COG-VOT		COG-f0		VOT-f0	
	Median	Range	Median	Range	Median	Range
p <sup>h</sup>	0.17	-0.41 – 0.62	0.09	-0.41 – 0.44	-0.04	-0.47 – 0.54
b	0.33	-0.14 – 0.77	-0.01	-0.53 – 0.64	0.00	-0.37 – 0.41
t <sup>h</sup>	0.06	-0.46 – 0.59	0.10	-0.34 – 0.51	-0.16	-0.59 – 0.41
d	0.54	-0.18 – 0.75	-0.06	-0.50 – 0.45	-0.03	-0.35 – 0.42
k <sup>h</sup>	0.16	-0.31 – 0.57	0.10	-0.34 – 0.52	-0.20	-0.61 – 0.39
g	0.54	-0.04 – 0.79	0.01	-0.51 – 0.39	0.00	-0.38 – 0.33