Dynamic acoustic vowel distances within and across dialects

Cynthia G. Clopper^{1, a}

¹Department of Linguistics, Ohio State University, Columbus, Ohio 43210, USA

Vowels vary in their acoustic similarity across regional dialects of American English, such that some vowels are more similar to one another in some dialects than others. Acoustic vowel distance measures typically evaluate vowel similarity at a discrete time point, resulting in distance estimates that may not fully capture vowel similarity in formant trajectory dynamics. In the current study, language and accent distance measures, which evaluate acoustic distances between talkers over time, were applied to the evaluation of vowel category similarity within talkers. These vowel category distances were then compared across dialects and their utility in capturing predicted patterns of regional dialect variation in American English was examined. Dynamic time warping of melfrequency cepstral coefficients was used to assess acoustic distance across the frequency spectrum and captured predicted Southern American English vowel similarity. Root-mean-square distance and generalized additive mixed models were used to assess acoustic distance for selected formant trajectories and captured predicted Southern, New England, and Northern American English vowel similarity. Generalized additive mixed models captured the most predicted variation, but, unlike the other measures, do not return a single acoustic distance value. All three measures are potentially useful for understanding variation in vowel category similarity across dialects.

^a Email: clopper.1@osu.edu

I. INTRODUCTION

1

22

23

24

2 Regional dialects of American English vary in the acoustic-phonetic realization of vowel 3 categories (Labov et al., 2006). This variation leads to vowel categories that are more similar in 4 acoustic-phonetic space in some dialects than others. For example, the Northern Cities vowel shift 5 leads to greater acoustic-phonetic similarity of $/\epsilon \propto /$ in the Northern dialect than in the Midland 6 dialect (Clopper and Tamati, 2014), the Southern vowel shift leads to greater acoustic-phonetic 7 similarity of /eɪ ɛ/ in the Southern dialect than in the Western dialect (Farrington et al., 2018), and 8 the low-back merger leads to greater acoustic-phonetic similarity of /a o/ in California English than 9 in New York City English (Nycz and Hall-Lew, 2013). Acoustic vowel category distance measures¹ are typically defined at a discrete time point within 10 11 the vowels of interest (e.g., Hay et al., 2006; Kendall and Fridland, 2012; Wassink, 2006). The goal of 12 the current study was to assess dynamic acoustic vowel category distance measures that capture 13 formant trajectories over time within the vowels of interest (e.g., Fox and Jacewicz, 2009; Renwick 14 and Stanley, 2020). These dynamic measures of vowel category distance within regional dialects of 15 American English were evaluated in comparison to qualitative descriptions of regional variation 16 (Labov et al., 2006). The results demonstrated that dynamic vowel category distance measures, 17 including dynamic time warping (DTW) of mel-frequency cepstral coefficients (MFCCs; Bartelds et 18 al., 2020; Lind-Combs et al., 2023; Mielke, 2012), root-mean-square distance (RMSD; Cole et al., 19 2023; Kaland, 2023), and generalized additive mixed modeling (GAMM; Kirkham et al., 2019; 20 Renwick and Stanley, 2020), can all capture variation in vowel category distances within and across 21 dialects of American English, complementing prior qualitative descriptions.

A. Vowel category distance measures

Regional variation in acoustic-phonetic similarity of vowel categories has been quantified in previous work using Euclidean distances and Pillai scores of formant frequency estimates (see Kelley

and Tucker, 2020; Nycz and Hall-Lew, 2013, for reviews). Euclidean distances between vowel categories are typically defined in the two-dimensional first formant (F1) x second formant (F2) space, using formant estimates from a discrete time point (e.g., vowel midpoint). Euclidean distances between vowels have been used to predict sound change (Wieling et al., 2012), as a proxy for individual participation in ongoing vowel shifts (Farrington et al., 2018; Kendall and Fridland, 2012), to predict accentedness ratings (Gunter et al., 2020), and to assess the effects of lexical competition on dialect variation (Clopper and Tamati, 2014). Pillai scores are also typically calculated from estimates of F1 and F2 from a discrete time point. Unlike Euclidean distances, Pillai scores quantify the overlap of the vowel category distributions in the F1 x F2 space, instead of simply the distance between the two category means. Pillai scores are used most commonly as a measure of vowel category merger within individuals (Gunter et al., 2020; Hay et al., 2006; Nycz and Hall-Lew, 2013). A number of related measures of multidimensional category overlap have also been proposed. These approaches typically involve estimates of F1 and F2 of the target vowels from a discrete time point, but may also include duration (Wassink, 2006), discrete cosine transformations of the formant frequency trajectories (Elvin et al., 2016), or formant frequencies sampled at multiple discrete time points (Morrison, 2008). These approaches have been used to identify the dimensions of variation that are necessary to distinguish vowel categories within a language or dialect (Elvin et al., 2016; Haynes and Taylor, 2014; Morrison, 2008; Wassink, 2006). In principle, these vowel category distance measures are appropriate for assessing variation in vowel distances within and across varieties of American English. For example, Euclidean distances between /eI ɛ/ have been used to distinguish Southern from Western varieties of American English (Farrington et al., 2018) and Pillai scores for α / have been used to distinguish degrees of the lowback vowel merger in California and New York City varieties (Nycz and Hall-Lew, 2013). However, at least as they are typically implemented, these measures are based on a single sample of each

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

formant estimated at a discrete time point for each vowel token. This implementation therefore does not capture vowel category distances at other time points in the formant trajectories or over the entire formant trajectory, despite clear evidence that vowels vary within and across dialects of American English in their formant trajectory dynamics (Farrington et al., 2018; Fox and Jacewicz, 2009). In contrast to this work on vowel category distances, dynamic acoustic distance measures have been used to assess consonant category distances (Mielke, 2012) and lexical distances (Kelley, 2023; Kelley and Tucker, 2022). Mielke (2012) used dynamic time warping of mel-frequency cepstral coefficients, estimated over VCV utterances, to assess the relationship between acoustic and phonological consonant similarity. Similarly, both Kelley (2023) and Kelley and Tucker (2022) used dynamic time warping of mel-frequency cepstral coefficients, estimated over words, to predict lexical competition in speech processing. In the current study, this DTW approach was applied to vowel category distances to assess variation within and across dialects of American English. Acoustic distances between languages, dialects, and accents have also been quantified in previous work using dynamic acoustic distance measures (e.g., Bartelds et al., 2020; Chernyak et al., 2024; Heeringa et al., 2009). These acoustic language and accent distance measures capture distances between different talkers' productions of the same linguistic content (e.g., words, sentences), whereas the acoustic vowel distance measures, such as Euclidean distances and Pillai scores, capture distances between different vowel categories produced by the same talkers. In the current study, the measures of between-talker language and accent distances were applied to within-talker vowel category distances to assess dynamic vowel category similarity within dialects of American English and to compare those vowel similarities across dialects. These comparisons across dialects were then evaluated in the context of previous qualitative descriptions of regional variation (Labov et al., 2006).

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

B. Language and accent distance measures

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

Three general approaches to quantifying acoustic language and accent distance have been proposed in the literature: dynamic time warping (Bartelds et al., 2020), mean distances over time (Heeringa et al., 2009), and generalized additive mixed modeling (Kirkham et al., 2019). Bartelds et al. (2020) and Lind-Combs et al. (2023) both used dynamic time warping of MFCCs, estimated over words and sentences, respectively, to assess native and non-native accent distances. This approach differs from the Euclidean distance and Pillai score measures of vowel category distance, and from the mean distance and GAMM approaches to language and accent distance, in that it includes a summary representation of the spectrum over the frequency analysis range (i.e., up to the Nyquist frequency), instead of focusing on estimated formant frequencies and their trajectories over time. DTW distances of MFCCs therefore capture dynamic formant trajectory distance, but may also capture distances between consonants, f0 contours, voice quality, and even recording devices and background noise (Bartelds et al., 2020; Lind-Combs et al., 2023). MFCCs have also been used without DTW to assess dialect similarity using multidimensional scaling and clustering techniques (Ferragne and Pellegrino, 2010). DTW of MFCCs was adopted in the current study as a measure of overall vowel distance, including dynamic formant, f0, and voice quality information, within and across dialects. The MFCCs were estimated over the target vowels only, so that consonant effects were minimized, and all distances were calculated within-talker, so that effects of recording device and background noise were also minimized. Heeringa et al. (2009) used Manhattan (or city-block) distances of formant estimates over time to quantify Norwegian dialect distances. This approach is conceptually similar to Euclidean distance measures of vowel category distance, except that (1) distances were calculated every 10 ms over entire words, including both consonants and vowels, and (2) Manhattan distances were used instead of Euclidean distances. Another conceptually similar measure is root-mean-square distance, in which

distance is defined as the mean Euclidean distance over time. This measure has been used to quantify the similarity of intonation contours (Cole et al., 2023; Kaland, 2023) and was adopted in the current study as a measure of dynamic formant trajectory distance within and across dialects.

Generalized additive mixed models of formant trajectories involve fitting curves to formant estimates over time and then identifying when in time the trajectories overlap and when they diverge. This approach differs from the other distance measures in that a discrete, positive distance measure is not returned. Rather, GAMMs provide an indication of (1) when in time the two formant trajectories differ, and (2) the direction of this difference (i.e., positive or negative). GAMMs have been used to compare vowel-lateral sequences in varieties of UK English (Kirkham et al., 2019), features of Southern American English among white and Black talkers in the American South (Renwick and Stanley, 2020), and Australian vowel change over time (Cox et al., 2024). GAMMs were adopted in the current study to assess the temporal properties of formant trajectory differences within and across regional dialects of American English.

C. The current study

Dynamic time warping of MFCCs, root-mean-square distance, and generalized additive mixed modeling were applied to three small datasets to assess their utility in quantifying dynamic acoustic vowel distances within and across dialects of American English. Each dataset comprised the stimulus materials from a previous perception task designed to examine cross-dialect lexical processing (Clark et al., 2022; Clopper and Walker, 2017; Ross and Clopper, 2023). Each dataset included minimal pair tokens for two vowel contrasts that were expected to differ in their acoustic similarity across two talker dialects, based on previous descriptions of regional variation in American English (Labov et al., 2006). Across datasets, the target vowel contrasts included /æ ɛ/ and /I e/ to capture features of the Northern Cities vowel shift, /aɪ a/ and /I e/ to capture features of the Southern vowel shift, and /aɪ a/ to capture non-rhoticity in New England. The three dynamic

distance measures were evaluated against the predicted patterns of vowel category similarity for these vowel contrasts. The vowel contrasts of interest were not predicted to be merged in any of the target varieties and the central question was therefore how to capture varying non-zero distances between vowel categories across regional dialects.

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

Given that the DTW approach successfully predicts perceptual native and non-native accent distances (Bartelds et al., 2020; Lind-Combs et al., 2023), DTW distances were expected to produce the predicted patterns of differential vowel distance across dialects. However, the cross-dialect vowel distance predictions are based on patterns of overall vowel shifts in one or two dimensions of the F1 x F2 vowel space (e.g., raising and fronting of /æ/ in the Northern Cities shift) and/or patterns of vowel formant trajectory dynamics in one or two dimensions of the F1 x F2 x F3 vowel space (e.g., lowering of F3 in /al/ sequences in rhotic varieties). Given that DTW distances were calculated over MFCCs, capturing more spectral information than just one or two formant trajectories, the predicted patterns of differential vowel distance across dialects may be masked by other similarities and differences in the spectra across tokens. For example, predicted differences in F1 might be masked by similarities in F2 and higher formants. The RMSD and GAMM analyses were therefore based on target formant frequencies that were selected to highlight the expected cross-dialect variation. This intentional focus on the target variation was expected to produce strong evidence of the predicted patterns of differential vowel distance across dialects for these measures, as in previous related work (Cox et al., 2024; Kirkham et al., 2019; Renwick and Stanley, 2020). Finally, given that DTW distances and RMSDs were estimated over the entire vowel trajectory, the predicted patterns of differential vowel distance across dialects may not be observed if they are limited to a short temporal span of the trajectory. For example, predicted differences in the offglide of a diphthong might be masked by similarities in the nucleus of the diphthong. The GAMM analysis allowed for a

consideration of this temporal detail and was expected to produce the strongest evidence of the predicted patterns of differential vowel distance across dialects.

II. METHODS

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

A. Datasets

The first dataset (Midland/Northern) comprises the stimulus materials from Clopper and Walker's (2017) cross-modal lexical decision task. The stimulus talkers were three Midland and three Northern female young adults. The target vowel contrasts were $/ \alpha \epsilon /$ and $/ 1 \epsilon I /$. These vowel contrasts were selected to capture variation across the two dialects due to the Northern Cities vowel shift. In particular, in the Northern Cities shift, $/\alpha$ is raised and fronted and $/\epsilon$ I/ are lowered and backed (Labov et al., 2006). The raising and fronting of $/\alpha$ / and the lowering and backing of $/\epsilon$ / is predicted to lead to a smaller distance between these vowels for Northern talkers than Midland talkers. In both varieties, /I/ is lower and backer than /eI/, so the lowering and backing of /I/ in the Northern Cities shift is predicted to lead to a larger distance between these vowels for Northern talkers than Midland talkers. The stimulus materials were 48 minimal pairs (24 per vowel contrast) for each of the six talkers. There were two missing tokens and the minimal pair token for each of these missing tokens was also excluded, leaving 286 minimal pairs (47-48 per talker) for analysis. The second dataset (Southern/Northern Virginian) comprises the stimulus materials from Clark et al.'s (2022) cross-modal lexical decision task. The stimulus talkers were four Southern and four Northern Virginian female young adults. The target vowel contrasts were $/\alpha I \alpha / \alpha I \alpha / \alpha I \alpha I$. These vowel contrasts were selected to capture variation across the two dialects due to the Southern vowel shift. In particular, in the Southern vowel shift, /ai/ is monophthongized and /I ɛ/ are raised and fronted, with greater raising and fronting of ϵ than /I/ (Labov et al., 2006). The monophthongization of $/\alpha I/$ is predicted to lead to a smaller distance between $/\alpha I$ $\alpha I/$ for Southern

talkers than Northern Virginian talkers. The more advanced raising and fronting of $/\epsilon$ / than /I/ is likewise predicted to lead to a smaller distance between these vowels for Southern talkers than Northern Virginian talkers. The stimulus materials were 69 word pairs (33 for $/\alpha I$ α / and 36 for /I ϵ /) for each of the eight talkers. There were eight missing tokens and the minimal pair token for each of these missing tokens was also excluded, leaving 544 minimal pairs (67-69 per talker) for analysis.

B. Acoustic distance analysis

For each dataset, the word tokens were stored in separate digital sound files that were segmented to the word onset and offset. These sound files were down-sampled to 16 kHz with 16-bit quantization for analysis. Each token was forced-aligned using the Penn Phonetics Lab Forced Aligner (Yuan and Liberman, 2008) to obtain a preliminary segmentation of the target vowel. These vowel alignments were then hand-corrected, following the segmentation guidance provided by

Peterson and Lehiste (1960). For the New England/Northern dataset, the target vowel /a/ was merged with the following /1/ for all /al/ tokens so that /al a/ distances could be estimated.

1. Dynamic time warping

The first distance measure was dynamic time warping of mel-frequency cepstral coefficients, following Bartelds et al. (2020) and Lind-Combs et al. (2023). The target vowel was extracted from each token with a Hamming window using the hand-corrected boundaries. MFCCs were then extracted using the *mfcc* function in the *python_speech_features* package in Python (Lyons et al., 2020) from each 25 ms window in each vowel token in 10 ms steps. For each window, 39 coefficients were extracted, including the overall energy and the first 12 cepstral coefficients, along with the first- and second-order derivatives of these measures. The spectrum used to calculate the MFCCs was extracted using a 1024-point FFT with 0.97 pre-emphasis. The MFCCs were z-scored over time separately for each coefficient for each token. Dynamic time warping was then used to estimate the distance between each minimal pair for each talker. The DTW distances were time-normalized to account for differences in overall duration of the paired tokens (Bartelds et al., 2020).

The DTW distances were analyzed using separate linear mixed-effects models for each dataset with talker dialect, vowel contrast, and their interaction as fixed effects. The maximal data-driven by-talker and by-minimal-pair random effects were used for each model. The fixed effects were sum-contrast coded. Statistical significance was assessed using the Satterthwaite approximation of degrees of freedom, as implemented in the *lmerTest* package in R (Kuznetsova et al., 2017). Post-hoc pairwise comparisons of significant effects were conducted using Tukey's HSD tests in the *emmeans* package in R (Lenth, 2024).

2. Root-mean-square distance

The second distance measure was root-mean-square distance of selected formant trajectories over time, following Cole et al.'s (2023) and Kaland's (2023) analyses of intonation contours. For

each token, the first three formant frequencies (F1, F2, and F3) were estimated at 10% intervals over the vowel duration from 0-100% (11 estimates per vowel) using a 12th order LPC analysis (Burg method) in the frequency range of 0-5500 Hz in Praat (Boersma and Weenink, 2023). All formant estimates were converted to Bark to facilitate comparisons across talkers and formants (Traunmüller, 1990). Given a large number of missing formant estimates at vowel onset (0%) and offset (100%), RMSDs were calculated for each minimal pair for each talker over the middle 80% of the vowel duration (10-90%). RMSDs were calculated separately for each formant and one formant was selected for each vowel contrast for analysis. For the Midland/Northern dataset, F1 was selected for both $/ \alpha \epsilon /$ and $/ I \epsilon I /$ distances to capture effects of raising of $/ \alpha /$ and lowering of $/ \epsilon$ I/ in the Northern Cities shift. For the Southern/Northern Virginian dataset, F2 was selected for /aɪ a/ distances to capture effects of /aɪ/ monophthongization (i.e., reduction of the fronting offglide) in the Southern vowel shift and F1 was selected for /I ɛ/ distances to capture effects of raising of /I ɛ/ in the Southern vowel shift. For the New England/Northern dataset, F3 was selected for /ax a/ distances to capture effects of non-rhoticity in New England and F1 was selected for $/\infty \epsilon$ distances to capture effects of raising of $/\infty$ and lowering of $/\epsilon$ in the Northern Cities shift, as in the Midland/Northern dataset. RMSDs were analyzed using separate linear mixed-effects models with the same specifications as in the DTW analysis. Due to missing formant estimates, the analysis of the Midland/Northern dataset was based on 283 minimal pairs, the analysis of the Southern/Northern Virginian dataset

was based on 543 minimal pairs, and the analysis of the New England/Northern dataset was based on 189 minimal pairs.

3. Generalized additive mixed modeling

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

The third distance measure involved generalized additive mixed models of selected formant trajectories over time, following Renwick and Stanley (2020). The same selected formant estimates from the RMSD analysis were used in the GAMMs, except that the entire formant trajectory (0-100%) was used because GAMMs are more robust to missing data than RMSDs. For each minimal pair for each talker, the formant trajectory difference was calculated by subtracting the formant estimate at each time point for one member of the minimal pair from the other. These formant difference trajectories were then analyzed separately for each vowel contrast using GAMMs with a parametric effect of talker dialect and talker dialect smooths over time as fixed effects and with bytalker smooths over time and minimal pair x talker dialect smooths over time as random effects. A correction term for autocorrelation at lag 1 was also included in each model.

III. RESULTS

A. Dynamic time warping

three datasets is shown in Fig. 1. The by-talker means overlap considerably for both vowel contrasts in the Midland/Northern (Fig. 1a) and New England/Northern (Fig. 1c) datasets, whereas they show clear separation by talker dialect in the Southern/Northern Virginian (Fig. 1b) dataset. As expected, the distances are smaller for both the / α I α / and /I ϵ / contrasts for the Southern talkers than the Northern Virginian talkers.

The linear mixed effects model predicting DTW distances from talker dialect and vowel contrast for the Midland/Northern dataset revealed a marginal main effect of talker dialect (β = .08, F(1, 236.5) = 3.60, p = .059). The Midland talkers had larger DTW distances overall, as expected for the / α ϵ / contrast, but contrary to the prediction for the /I ϵ I/ contrast. Neither the main effect of vowel contrast nor the interaction were significant. The model for the Southern/Northern Virginian dataset revealed a significant main effect of talker dialect (β = .25, F(1, 7.1) = 29.82, p < .001). The Northern Virginian talkers had larger DTW distances overall, as expected for both vowel contrasts.

A summary of the mean DTW distances for each talker for each vowel contrast in each of the

Neither the main effect of vowel contrast nor the interaction were significant. The model for the New England/Northern dataset revealed a significant interaction between talker dialect and vowel contrast (β = -0.10, F(1, 140.0) = 4.80, p = .030). The interaction reflects the expected cross-over pattern, in which DTW distances were larger for the Northern talkers than the New England talkers for the /ai a/ contrast (β = -.29), but larger for the New England talkers than the Northern talkers for the /æ ϵ / contrast (β = .10). However, post-hoc comparisons of estimated marginal means revealed that the talker dialect effect was not significant for either vowel contrast. The main effects of talker dialect and vowel contrast were also not significant.

Overall, the DTW results are suggestive of the predicted patterns, but are only statistically robust with these small samples in the Southern/Northern Virginian dataset. Moreover, the marginal result in the Midland/Northern dataset is contrary to predictions.

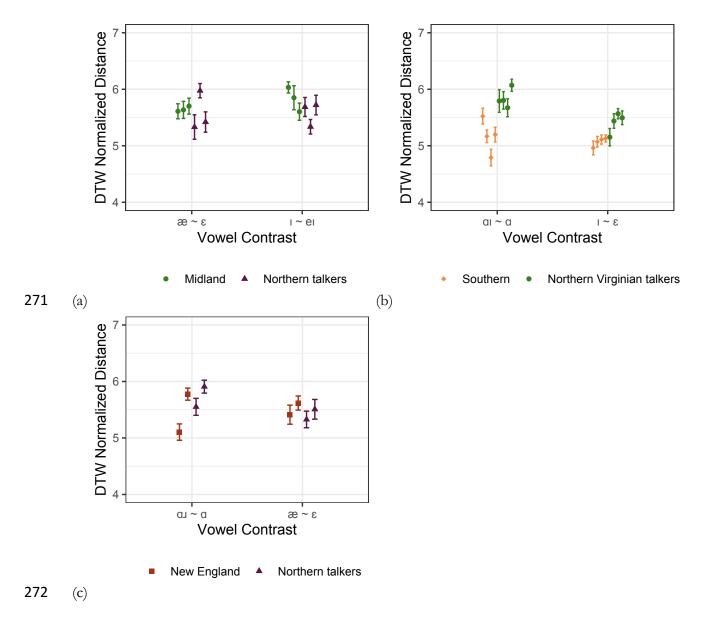


FIG 1. (Color online). Mean DTW distances for each talker for each vowel contrast for the Midland/Northern (a), Southern/Northern Virginian (b), and New England/Northern (c) datasets. Error bars are standard errors.

B. Root-mean-square distance

A summary of the mean RMSDs for each talker for each vowel contrast in each of the three datasets is shown in Fig. 2. The by-talker means show the expected talker dialect differences for the /I eI/ contrast in the Midland/Northern (Fig. 2a) dataset, for both vowel contrasts in the

Southern/Northern Virginian (Fig. 2b) dataset, and for the /αι α/ contrast in the New

England/Northern (Fig. 2c) dataset. The by-talker means overlap considerably for the /æ ε/

contrast in both the Midland/Northern and New England/Northern datasets.

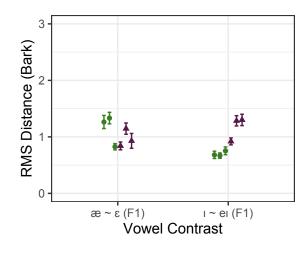
The linear mixed effects model predicting RMSDs from talker dialect and vowel contrast for the Midland/Northern dataset revealed a significant interaction between talker dialect and vowel contrast (β = -.16, F(1, 5.2) = 8.54, p = .032). Post-hoc comparisons of estimated marginal means for each vowel contrast revealed a significant effect of talker dialect for the /I eI/ contrast (β = -.46, t(5.5) = -3.34, p = .018). As expected, the Northern talkers had larger RMSDs than the Midland talkers. The talker dialect effect was not significant for the / α ϵ / contrast. The main effects of talker dialect and vowel contrast were also not significant.

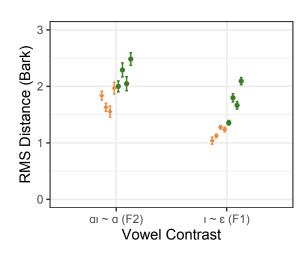
The model for the Southern/Northern Virginian dataset revealed significant main effects of talker dialect (β = .25, F(1, 6.4) = 12.75, p = .011) and vowel contrast (β = .26, F(1, 25.1) = 28.73, p < .001). As expected, the Northern Virginian talkers had larger RMSDs than the Southern talkers overall. The overall RMSDs were also larger for the / α I α / contrast than the /I ϵ / contrast. This vowel contrast effect likely reflects inherent differences between the two contrasts: whereas the / α I α / contrast involves a difference between a diphthong and a monophthong, where larger distances might be expected, the /I ϵ / contrast involves a difference between two lax vowels, where shorter distances might be expected. The interaction between talker dialect and vowel contrast was not significant.

The model for the New England/Northern dataset revealed a significant main effect of vowel contrast (β = .35, F(1, 183.0) = 107.27, p < .001) and a significant interaction between talker dialect and vowel contrast (β = -.24, F(1, 183.0) = 52.23, p < .001). The overall RMSDs were larger for the / α I α / contrast than the / α E/ contrast. As in the Southern/Northern Virginian dataset, this effect likely reflects inherent properties of the two contrasts, including a vowel-rhotic sequence in the / α I

 α contrast versus two lax vowels in the $/\alpha \epsilon$ contrast. Post-hoc comparisons of estimated marginal means for each vowel contrast revealed a marginal effect of talker dialect for the $/\alpha a$ contrast (β = -1.01, t(2.2) = -3.51, p = .062). As expected, the Northern talkers had larger RMSDs than the New England talkers. The talker dialect effect was not significant for the $/ \alpha \epsilon /$ contrast.

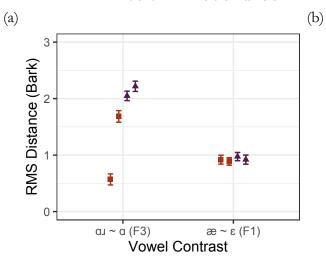






Midland Northern talkers

Northern Virginian talkers Southern



New England Northern talkers

310 (c)

304

305

306

307

308

FIG 2. (Color online). Mean RMSDs for each talker for each vowel contrast for the Midland/Northern (a), Southern/Northern Virginian (b), and New England/Northern (c) datasets.

Error bars are standard errors.

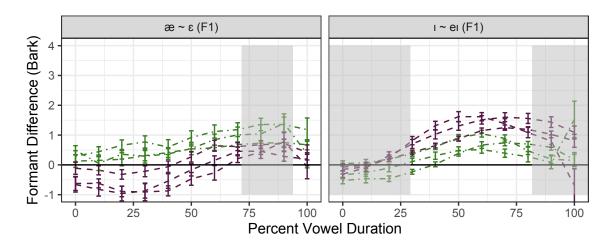
Like the DTW results, the RMSD results are suggestive of the predicted patterns, although more of the predicted patterns are statistically robust and there is no evidence of unexpected patterns. In particular, in the RMSD analysis, significant talker dialect effects in the predicted direction are observed for the /I eI/ contrast in the Midland/Northern dataset and for both vowel contrasts in the Southern/Northern Virginian dataset. The predicted effect is marginal for the /aɪ a/ contrast in the New England/Northern dataset.

C. Generalized additive mixed modeling

A summary of the mean formant trajectory differences for each talker for each vowel contrast in each of the three datasets is shown in Fig. 3. The by-talker means show the expected talker dialect differences for both contrasts in the Midland/Northern (Fig. 3a) dataset, for both vowel contrasts in the Southern/Northern Virginian (Fig. 3b) dataset, and for the / α I α / contrast in the New England/Northern (Fig. 3c) dataset. The by-talker means overlap considerably for the / α E/ contrast in the New England/Northern dataset.

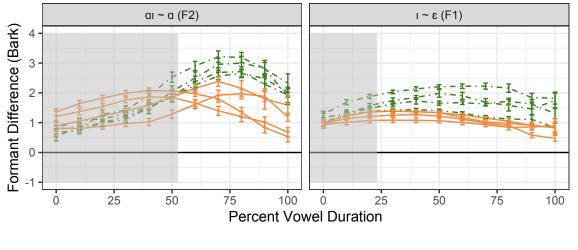
The GAMMs predicting formant trajectory differences over time for the Midland/Northern dataset revealed significant main effects of talker dialect for both the / α ϵ / contrast (β = -.77, t = -3.00, p = .003) and the /I eI/ contrast (β = .48, t = 3.21, p = .001). Significant talker dialect differences in the trajectories were observed for the / α ϵ / contrast between 0% and 72% and between 94% and 100% of the vowel duration and for the /I eI/ contrast between 29% and 82% of the vowel duration. Formant differences were larger for the Midland talkers than the Northern talkers for the / α ϵ / contrast and for the Northern talkers than the Midland talkers for the /I eI/

contrast, as expected. The GAMMs for the Southern/Northern Virginian dataset revealed significant main effects of talker dialect for both the / α I α / contrast (β = -.39, t = -2.39, p = .017) and the /I ϵ / contrast (β = -.51, t = -3.17, p = .002). Significant talker dialect differences in the trajectories were observed for the / α I α / contrast between 53% and 100% of the vowel duration and for the /I ϵ / contrast between 23% and 100% of the vowel duration. Formant trajectory differences were larger for the Northern Virginian talkers than the Southern talkers for both vowel contrasts, as expected. The GAMMs for the New England/Northern dataset revealed a significant main effect of talker dialect for the / α I α / contrast (β = 1.11, t = 2.13, p = .034). Significant talker dialect differences in the trajectories were observed for the / α I α / contrast between 42% and 100% of the vowel duration. The talker dialect effect was not significant for the / α E/ contrast. Formant trajectory differences were larger for the Northern talkers than the New England talkers for the / α I α / contrast, as expected.



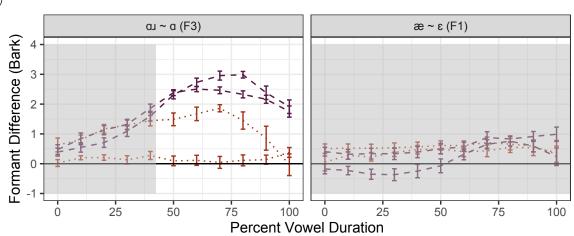
· - · Midland - - Northern talkers

346 (a)



Southern · - · Northern Virginian talkers

347 (b)



· · · New England - - Northern talkers

348 (c)

349

350

351

352

353

354

355

FIG 3. (Color online). Mean formant trajectory differences for each talker for each vowel contrast for the Midland/Northern (a), Southern/Northern Virginian (b), and New England/Northern (c) datasets. Error bars are standard errors. Temporal regions with non-significant talker dialect differences in the GAMMs are shaded in gray.

Like the DTW and RMSD results, the GAMMs are suggestive of the predicted patterns, although more of the predicted patterns are statistically robust and there is no evidence of unexpected patterns. In particular, in the GAMM analysis, significant talker dialect effects in the

predicted direction are observed for both vowel contrasts in the Midland/Northern dataset, for both vowel contrasts in the Southern/Northern Virginian dataset, and for the /aı a/ contrast in the New England/Northern dataset.

IV. DISCUSSION

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

All three of the acoustic distance measures considered in this study—DTW, RMSD, and GAMMs—led to predicted patterns of vowel category distance within and across dialects. DTWs revealed significant effects in the expected direction for both the $\langle \alpha i \alpha \alpha \rangle$ and $\langle i \alpha \alpha \rangle$ contrasts in the Southern/Northern Virginian dataset. RMSDs revealed significant effects in the expected direction for these two contrasts, as well as for the /I eI/ contrast in the Midland/Northern dataset. GAMMs revealed significant effects in the expected direction for these three contrasts, as well as for the /æ ε/ contrast in the Midland/Northern dataset and the /αι α/ contrast in the New England/Northern dataset. The only predicted effect that was not observed across any of the distance measures was for the $/\infty \epsilon$ contrast in the New England/Northern dataset. This null result likely reflects the variability between the two Northern talkers in this dataset. As shown in Fig. 3c, one of the Northern talkers produced formant trajectory differences for the $/ \alpha \epsilon / contrast$ that are similar to the Northern talkers in the Midland/Northern dataset (see Fig. 3a), whereas the other Northern talker's mean formant trajectory difference for the $/ \approx \epsilon /$ contrast was similar to the two New England talkers in the New England/Northern dataset. Although the three datasets that were analyzed in the current study are all rather small, the results suggest that dynamic acoustic distance measures that have been used to assess language and accent distances in previous work (e.g., Bartelds et al., 2020; Heeringa et al., 2009; Renwick and Stanley, 2020) can usefully be applied to assessments of vowel category distance within and across dialects. This quantification of dynamic vowel category distance in the current study complements

previous qualitative descriptions of regional vowel variation in American English (e.g., Labov et al., 2006), as well as previous measures of vowel category distance at discrete time points (e.g., Hay et al., 2006; Kendall and Fridland, 2012; Wassink, 2006).

Moreover, although the magnitude of the DTW distances are not readily interpretable, the RMSD and GAMM analyses were conducted in Bark and the magnitudes of the effects are therefore interpretable from the model coefficients. The significant talker dialect effects ranged in magnitude from .25 Bark in the RMSD analysis of the Southern/Northern Virginian dataset to 1.11 Bark in the GAMM analysis of the New England/Northern dataset. These differences are likely perceptible, given vowel discrimination thresholds of approximately .28 Bark for American English listeners (Kewley-Port and Zheng, 1999). Thus, just as these kinds of distance measures predict perception of language and accent distances (e.g., Heeringa et al., 2009; Gunter et al., 2020; Porretta et al., 2015) and performance in cross-accent speech processing tasks (e.g., Chernyak et al., 2024; Hay et al., 2006), the dynamic measures examined in this study are also likely to predict judgments about regional accents and cross-dialect speech processing performance. Exploring the relationship between these dynamic vowel category distance measures and the perception of regional dialects is a critical next step for linking acoustic distances between vowel categories to speech processing performance (see also Bent et al., 2021; Kelley, 2023; Kelley and Tucker, 2022).

The increase in the number of predicted patterns that were observed from the DTW analysis to the RMSD analysis to the GAMM analysis is consistent with the information that is used to estimate distance in each approach. DTW of MFCCs captures more information in the spectrum than the formant-based RMSD and GAMM analyses, including voice quality, f0, and higher formant structure. The MFCCs also include first- and second-order derivatives of the cepstral coefficients, which explicitly capture the slope and rate of change of these coefficients over time, in addition to the time-varying coefficients themselves. These sources of additional information may mask the

predicted differences, which are characterized by one or two time-varying formants, leading to fewer predicted results. However, listeners also have access to this broader array of information in the spectrum and so DTW distances may better predict perception than RMSD or GAMMs. Future research exploring the relationship between these measures and the perception of regional dialects is critical for understanding the roles of general acoustic distance and specific formant trajectory distance in cross-dialect speech processing.

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

The RMSD and GAMM analyses in the current study captured formant trajectory information from a single formant that was selected for each vowel contrast to highlight the predicted variation in vowel category distance within and across dialects. Despite this similarity, these two analyses differ in two primary ways. First, RMSDs return a discrete measure of distance, whereas GAMMs assess difference throughout the formant trajectory. Second, RMSDs return a positive distance value, which can mask variation in the direction of the formant trajectory difference, whereas GAMMs assess signed (positive and negative) differences. The former difference in the two approaches likely explains the significant difference that was observed for the $/\alpha I \alpha I$ contrast in the New England/Northern dataset in the GAMM analysis, but not in the RMSD analysis. The formant trajectory differences were only significantly different for the two talker dialects in the second half of the contrast, consistent with the contrast between the presence vs. absence of the following rhotic for Northern vs. New England talkers, respectively (see Fig. 3c). This difference is captured by the GAMM analysis, but lost in the RMSD analysis, which returns just a single distance value, collapsed over time. However, this single distance value makes RMSDs more suitable than GAMMs for predicting perception because RMSDs can be straightforwardly entered as predictors in statistical models.

The latter difference between the two approaches likely explains the significant difference that was observed for the $/\alpha \epsilon$ contrast in the Midland/Northern dataset in the GAMM analysis, but

not in the RMSD analysis. As shown in Fig. 3a, whereas the Midland talkers had a positive formant trajectory difference for this contrast throughout the entire trajectory, the Northern talkers had a negative formant trajectory difference in the first half of the contrast and a positive formant trajectory difference in the second half of the contrast. The overall magnitudes of the formant trajectory differences were comparable across dialects, so that the absolute, positive distance measure in the RMSD analysis could not distinguish the two trajectory differences, whereas the signed differences emerge as significant in the GAMM analysis. Future research exploring the relationship between these measures and the perception of regional dialects is critical for understanding the roles of overall formant trajectory distance and the direction and magnitude of formant trajectory difference in cross-dialect speech processing.

The three analyses that were conducted in this study are a subset of the possible applications of dynamic distance measures to vowel category distances within and across dialects. First, one formant was selected for analysis in the RMSD and GAMM analyses, whereas other formants may also exhibit interesting differences in distance for these vowel contrasts. For example, the analysis of the /aɪ a/ contrast in the Southern/Northern Virginian analysis focused on F2 to capture fronting of the offglide, but it would also be reasonable to analyze F1 to capture raising of the offglide. Second, the DTW analysis was based on MFCCs, whereas the RMSD and GAMM analyses were based on formant estimates. A DTW analysis could also be conducted on the formant trajectories to disentangle DTW as a distance measure from the information about the spectrum included in the analysis.² Finally, several recent studies have used discrete cosine transformations (DCTs) to capture formant trajectories (e.g., Cox et al., 2024; Elvin et al., 2016). Euclidean distance could be applied to DCTs to further disentangle the distance measures from the characterization of the spectrum.

In summary, vowel categories are acoustically more similar to one another in some regional dialects of American English than in others, even in the absence of vowel category mergers (Labov

et al., 2006). The application of dynamic acoustic distance measures, including DTW, RMSD, and GAMMs, can capture this variation in vowel category acoustic distance within and across dialects. The analyses that focused on specific formant trajectory distance (i.e., RMSD, GAMMs) captured more of these predicted patterns than the analysis that focused on overall acoustic distance (i.e., DTW), although overall acoustic distance may be a better predictor of perception. Within the analyses that focused on specific formant trajectory distance, the analysis that focused on formant trajectory difference (i.e., GAMMs) captured more of these predicted patterns than the analysis that focused on formant trajectory distance (i.e., RMSD), although RMSD may be more useful for predicting perception because it returns a discrete distance value for each minimal pair. These approaches therefore address varying aspects of the nature of dynamic vowel category distance and may be most useful for distinct kinds of questions. The relationship between these acoustic distance measures and human perception remains to be explored.

ACKNOWLEDGMENTS

This work was partially supported by the National Science Foundation (BCS-1843454).

AUTHOR DECLARATIONS

Conflict of Interest

The author has no conflicts to disclose.

Ethics Approval

The analysis involved existing datasets. The Midland/Northern and Southern/Northern Virginian datasets are de-identified. The analysis of the New England/Northern dataset was approved by the Institutional Review Board at Ohio State University (2017B0576) and informed consent was obtained from all participants.

473 DATA AVAILABILITY 474 The data that support the findings of this study are available on the OSF repository for this 475 project: https://osf.io/hp6xc/ 476 **ENDNOTES** 477 ¹ Although many of the approaches discussed here are not distance metrics in the mathematical 478 sense (including dynamic time warping, Pillai scores, and generalized additive mixed models), they 479 can usefully be applied to questions related to vowel category acoustic distance and so are referred to 480 here as "distance measures," as in previous work (e.g., Bartelds et al., 2020; Mielke, 2012). 481 ² I would like to thank Matthew Kelley for this suggestion. 482 REFERENCES 483 Bartelds, M., Richter, C., Liberman, M., and Wieling, M. (2020). "A new acoustic-based 484 pronunciation distance measure," Front. Artif. Intell. 3(39), 1–10. 485 Bent, T., Holt, R. F., Van Engen, K. J., Jamsek, I. A., Arzbecker, L. J., Liang, L., and Brown, E. (2021). "How pronunciation distance impacts word recognition in children and adults," J. 486 487 Acoust. Soc. Am. 150, 4103–4117. 488 Boersma, P., and Weenink, D. (2023). "Praat: Doing phonetics by computer," Computer program, 489 Version 6.4.01. 490 Chernyak, B. R., Bradlow, A. R., Keshet, J., and Goldrick, M. (2024). "A perceptual similarity space 491 for speech based on self-supervised speech representations," J. Acoust. Soc. Am. 155, 3915— 492 3929. 493 Clark, H., Bissell, M., Clopper, C. G., and Walker, A. (2022). "Effects of lexical competition, dialect 494 familiarity, and dialect exposure on lexical processing," Poster presented at Laboratory

495

Phonology 18, virtual, June 23–25.

- Clopper, C. G., and Tamati, T. N. (2014). "Effects of local lexical competition and regional dialect
- on vowel production," J. Acoust. Soc. Am. 136, 1–4.
- 498 Clopper, C. G., and Walker, A. (2017). "Effects of lexical competition and dialect exposure on
- 499 phonological priming," Lang. Speech 60, 85–109.
- 500 Cole, J., Steffman, J., Shattuck-Hufnagel, S., and Tilsen, S. (2023). "Hierarchical distinctions in the
- production and perception of nuclear tunes in American English," *Lab. Phonol.* **14**(1), 1–51.
- 502 Cox, F., Penney, J., and Palethorpe, S. (2024). "Australian English monophthong change across 50
- years: Static vs. dynamic measures," *Languages* **9**(99), 1–35.
- Elvin, J., Williams, D., and Escudero, P. (2016). "Dynamic acoustic properties of monophthongs
- and diphthongs in Western Sydney Australian English," J. Acoust. Soc. Am. 140, 576–581.
- Farrington, C., Kendall, T., and Fridland, V. (2018). "Vowel dynamics in the Southern Vowel Shift,"
- 507 *Am. Speech* **93**, 186–222.
- Ferragne, E., and Pellegrino, F. (2010). "Vowel systems and accent similarity in the British Isles:
- Exploiting multidimensional acoustic distances in phonetics," *J. Phonetics* **38**, 526–539.
- 510 Fox, R. A., and Jacewicz, E. (2009). "Cross-dialectal variation in formant dynamics of American
- English vowels," *J. Acoust. Soc. Am.* **126**, 2603–2618.
- 512 Gunter, K. M., Vaughn, C. R., and Kendall, T. S. (2020). "Perceiving Southernness: Vowel
- categories and acoustic cues in Southernness ratings," J. Acoust. Soc. Am. 147, 643–656.
- Hay, J., Warren, P., and Drager, K. (2006). "Factors influencing speech perception in the context of
- a merger-in-progress," J. Phonetics 34, 458–484.
- Haynes, E. F., and Taylor, M. (2014). "An assessment of acoustic contrast between long and short
- 517 vowels using convex hulls," *J. Acoust. Soc. Am.* **136**, 883–891.
- Heeringa, W., Johnson, K., and Gooskens, C. (2009). "Measuring Norwegian dialect distances using
- 519 acoustic features," Speech Comm. 51, 167–183.

- Kaland, C. (2023). "Intonation contour similarity: f0 representations and distance measures compared
- to human perception in two languages," J. Acoust. Soc. Am. 154, 95–107.
- Kelley, M. C. (2023). "Acoustic absement in detail: Quantifying acoustic differences across time-
- series representations of speech data," *Proc. 20th Int. Congress Phonetic Sci.* 679–683.
- Kelley, M. C., and Tucker, B. V. (2020). "A comparison of four vowel overlap measures," J. Acoust.
- 525 *Soc. Am.* 147, 137–145.
- Kelley, M. C., and Tucker, B. V. (2022). "Using acoustic distance and acoustic absenuent to quantify
- 527 lexical competition," *J. Acoust. Soc. Am.* **151**, 1367–1379.
- Kendall, T., and Fridland, V. (2012). "Variation in perception and production of mid front vowels in
- the U.S. Southern Vowel Shift," J. Phonetics 40, 289–306.
- Kewley-Port, D. K., and Zheng, Y. (1999). "Vowel formant discrimination: Towards more ordinary
- 531 listening conditions," *J. Acoust. Soc. Am.* **106**, 2945–2958.
- Kirkham, S., Nance, C., Littlewood, B., Lightfoot, K., and Groarke. E. (2019). "Dialect variation in
- formant dynamics: The acoustics of lateral and vowel sequences in Manchester and
- Liverpool English," *J. Acoust. Soc. Am.* **145**, 784–794.
- Kuznetsova, A., Brockhoff, P. B., and Christensen, R. H. B. (2017). "ImerTest package: Tests in
- linear mixed effects models," J. Stat. Software 82(13), 1–26.
- 537 Labov, W., Ash, S., and Boberg, C. (2006). Atlas of North American English (Mouton de Gruyter, New
- 538 York).
- 539 Lenth, R. (2024). "emmeans: Estimated marginal means, aka least-squares means," R package,
- Version 1.10.0.
- Lind-Combs, H. C., Bent, T., Holt, R. F., Clopper, C. G., and Brown, E. (2023). "Comparing
- Levenshtein distance and dynamic time warping in predicting listeners' judgments of accent
- 543 distance," *Speech Comm.* **155**(102987), 1–14.

- Lyons, J., Wang, D. Y.-B., Gianluca, Shteingart, H., Mavrinac, E., Gaurkar, Y., Watcharawisetkul,
- 545 W., Birch, S., Lu. Z., Hölzl, J., Lesinskis, J., Almér, H., Lord, C., and Stark, A. (2020).
- 546 "jameslyons/python_speech_features: release v0.6.1," *Zenodo*.
- 547 https://doi.org/10.5281/zenodo.3607820.
- 548 Mielke, J. (2012). "A phonetically based metric of sound similarity," *Lingua* 122, 145–163.
- Morrison, G. S. (2008). "Comment on 'A geometric representation of spectral and temporal vowel
- features: Quantification of vowel overlap in three linguistic varieties' []. Acoust. Soc. Am.
- 551 119, 2334–2350 (2006)] (L)," J. Acoust. Soc. Am. 123, 37–40.
- Nycz, J., and Hall-Lew, L. (2013). "Best practices in measuring vowel merger," *Proc. Mtgs. Acoust.*
- **20**(060008), 1–19.
- Peterson, G. E., and Lehiste, I. (1960). "Duration of syllable nuclei in English," J. Acoust. Soc. Am.
- **32**, 693–703.
- Porretta, V., Kyröläinen, A.-J., and Tucker, B. V. (2015). "Perceived foreign accentedness: Acoustic
- 557 distances and lexical properties," *Atten. Percept. Psychophys.* 77, 2438–2451.
- Renwick, M. E. L., and Stanley, J. A. (2020). "Modeling dynamic trajectories of front vowels in the
- 559 American South," *J. Acoust. Soc. Am.* **147**, 579–595.
- Ross, J., and Clopper, C. G. (2023). "Talker variability in cross-dialect lexical processing," *Proc.* 20th
- Int. Congress Phonetic Sci. 152–156.
- Traunmüller, H. (1990). "Analytical expressions for the tonotopic sensory scale," J. Acoust. Soc. Am.
- **88**, 97–100.
- Wassink, A. B. (2006). "A geometric representation of spectral and temporal vowel features:
- Quantification of vowel overlap in three linguistic varieties," J. Acoust. Soc. Am. 119, 2334—
- **566** 2350.

Wieling, M., Margaretha, E., and Nerbonne, J. (2012). "Inducing a measure of phonetic similarity
from pronunciation variation," *J. Phonetics* 40, 307–314.
Yuan, J. and M. Liberman. (2008). "Speaker identification on the SCOTUS corpus," *Proc. Mtgs.*Acoust. 2008, 5687–5690.