

A New Acoustic-based Pronunciation Distance Measure

Martijn Bartelds^{1,*}, Caitlin Richter², Mark Liberman² and Martijn Wieling¹

¹Center for Language and Cognition, Faculty of Arts, University of Groningen, Groningen, Netherlands

²Department of Linguistics, University of Pennsylvania, Philadelphia, PA, USA

Correspondence*:
Martijn Bartelds
m.bartelds@rug.nl

2 ABSTRACT

3 We present an acoustic distance measure for comparing pronunciations, and apply the measure
4 to assess foreign accent strength in American-English by comparing speech of non-native
5 American-English speakers to a collection of native American-English speakers. An acoustic-
6 only measure is valuable as it does not require the time-consuming and error-prone process
7 of phonetically transcribing speech samples which is necessary for current edit distance-based
8 approaches. We minimize speaker variability in the data set by employing speaker-based cepstral
9 mean and variance normalization, and compute word-based acoustic distances using the dynamic
10 time warping algorithm. Our results indicate a strong correlation of $r = -0.71$ ($p < 0.0001$) between
11 the acoustic distances and human judgments of native-likeness provided by more than 1,100
12 native American-English raters. Therefore, the convenient acoustic measure performs only slightly
13 lower than the state-of-the-art transcription-based performance of $r = -0.77$. We also report the
14 results of several small experiments which show that the acoustic measure is not only sensitive
15 to segmental differences, but also to intonational differences and durational differences. However,
16 it is not immune to unwanted differences caused by using a different recording device.

17 **Keywords:** Acoustic measure, Acoustic features, Foreign accent, Mel frequency cepstral coefficient, Pronunciation, Spoken language
18 processing, Validation

19 **Word count:** 4762.

20 **Number of Figures:** 3.

21 **Number of Tables:** 4.

INTRODUCTION

22 The strength of foreign accent in a second language is mainly caused by the first language background
23 of non-native speakers, and is influenced by a wide variety of variables with the most valuable predictor
24 being the age of second-language learning (Asher and García, 1969; Leather, 1983; Flege, 1988; Arslan
25 and Hansen, 1997). Understanding the factors that affect the degree of foreign accent may be essential
26 for second language teaching, and knowledge about the acoustic features of foreign-accented speech can
27 improve speech recognition models (Arslan and Hansen, 1996; Piske et al., 2001). Computational methods
28 that investigate foreign accent strength are, however, scarce.

29 Studies that investigate and compare different pronunciations often use transcribed speech (Nerbonne and
30 Heeringa, 1997; Livescu and Glass, 2000; Gooskens and Heeringa, 2004; Heeringa, 2004; Wieling et al.
31 2011; Chen et al., 2016; Jeszenszky et al., 2017). For example, Kessler (1995) presented the Levenshtein
32 distance for finding linguistic distances between language varieties. To calculate the Levenshtein distance,
33 speech samples have to be manually transcribed using a phonetic alphabet, but this process is very
34 time consuming and labor intensive (Hakkani-Tür et al., 2002; Novotney and Callison-Burch, 2010).
35 Furthermore, transcribing speech is prone to errors, and interference from transcriber variation might
36 lead to a sub-optimal distance calculation when differences in transcribers' habits cannot be distinguished
37 from differences in speakers' productions (Bucholtz, 2007). Another limitation of this approach is that
38 the set of discrete symbols used in phonetic transcriptions is unable to capture all the acoustic details
39 that are relevant for studying accented pronunciations (Cucchiari, 1996). As Mermelstein (1976) notes,
40 transcribing speech results in a loss of information whereby perceptually distinct differences between
41 sounds diminish or largely disappear. For example, problems may arise when fine-grained pronunciation
42 differences cannot be represented by the set of transcription symbols (Duckworth et al., 1990), or when an
43 important dimension of difference between accents is their use of tone, but no tone or pitch information is
44 transcribed (Heeringa et al., 2009). It is therefore potentially useful to develop an acoustic-only method
45 to study pronunciation differences, such as foreign accent strength in the speech of non-native speakers.
46 Fine-grained characteristics of human speech are preserved in the speech representations, while at the same
47 time a time consuming and costly process may be omitted.

48 To evaluate computational methods of determining accent differences, validation against reliable data
49 regarding these differences is necessary, which usually consists of comparing the automatically obtained
50 ratings to human judgments of accent strength. Derwing and Munro (2009) stress the importance of
51 including human judgments, since these provide the most appropriate method to evaluate these measurement
52 techniques. Studies that compare human perceptual judgments to a computational difference measure which
53 is not based on the alignment of phonetic transcriptions are uncommon, despite the potential advantages of
54 this approach. This may be due to the challenges of directly comparing speech samples, as there exists a
55 considerable amount of variability in the signal. A substantial amount of variability in the structure of a
56 speech signal is also dependent on non-linguistic characteristics of the speakers, which may mask relevant
57 phonetic information in acoustic measurements (Goslin et al., 2012). For example, Heeringa et al. (2009)
58 calculated speaker-dependent pronunciation distances for a set of fifteen speakers from different Norwegian
59 varieties and for a subset of eleven female speakers. The Manhattan distance was computed between the
60 frequency values of the first three formants per vowel in each word. Correlations between their procedure
61 and human judgments of native-likeness only ranged from $r = 0.36$ to $r = 0.60$ ($p < 0.001$). Given that they
62 only obtained a moderate correlation with the human judgments, their acoustic-based measure could not
63 serve as a reliable alternative to transcription-based methods for assessing accent differences.

64 The primary goal of this study is therefore to develop an improved acoustic pronunciation distance
65 measure that computes pronunciation distances without requiring phonetic transcriptions. To assess whether
66 the acoustic distance measure is a valid measurement technique to measure accent strength (compared to
67 native speakers), we compare the acoustic distances to a collection of human native-likeness judgments
68 that were collected by Wieling et al. (2014) to evaluate a phonetic transcription-based method. The core
69 of the acoustic distance measure is to use dynamic time warping (DTW) to compare non-native accented
70 American-English to native-accented American-English speech samples represented as Mel-frequency

71 cepstral coefficients (MFCCs)¹ In short, our approach consists of obtaining word-level acoustic differences,
72 which are averaged to obtain speaker-based acoustic differences. To make the comparison less dependent
73 on individual speaker characteristics, we use speaker-based cepstral mean and variance normalization
74 before calculating the word-level acoustic differences. We evaluate the method by comparing the acoustic
75 distances to both transcription-based pronunciation distances and human perception. To better understand
76 what (desired and less desired) differences are captured by our acoustic difference measure, we conduct
77 several small-scale experiments.

MATERIALS AND METHODS

78 Speech Accent Archive

79 We use data from the Speech Accent Archive, which contains over 2000 speech samples from both native
80 and non-native American-English speakers (Weinberger, 2015). For each participant an acoustic voice
81 recording of the same standard 69-word-paragraph is present. The paragraph is primarily composed of
82 common English words, and contains a wide variety of consonants and vowels that can be found in the
83 English language. The paragraph is shown in (1).

84 (1) *Please call Stella. Ask her to bring these things with her from the store: Six spoons of fresh snow*
85 *peas, five thick slabs of blue cheese, and maybe a snack for her brother Bob. We also need a small*
86 *plastic snake and a big toy frog for the kids. She can scoop these things into three red bags, and we*
87 *will go meet her Wednesday at the train station.*

88 The availability of data from both native and non-native speakers of American-English enables us to
89 compare the accents of a broad range of different speakers of English (Weinberger and Kunath, 2011).
90 Speech samples from 280 non-native American-English speakers make up our target non-native speaker data
91 set, and 115 speech samples from U.S.-born L1 speakers of English serve as our reference native speaker
92 data set. For each non-native speaker the goal is to determine how different that speaker's pronunciation
93 is on average from the native American-English speakers in the reference native speaker data set. We do
94 not rely on choosing a single native American-English reference speaker, as there is considerable regional
95 variability in the data set. (The native American-English speakers who rated the non-native speech samples
96 also had different regional backgrounds.)

97 The data we include in this study is similar to the data used for evaluating a transcription-based
98 measurement in the study of Wieling et al. (2014). As in some cases a word was produced twice by
99 a speaker, or two words were merged into one word, we removed duplicate words from the speech samples
100 by deleting one of the repeated words, and merged words were split such that each speech sample consisted
101 of 69 separate words.

102 Our data set contains slightly more male speakers (206) than female speakers (189). The average age of
103 all speakers in our data set is 32.6 years with a standard deviation of 13.5 years. In the target non-native
104 speaker data set, the average age of starting to learn English is 10.5 years with a standard deviation of
105 6.6 years. The 280 non-native English speakers have a total of 99 different native languages. The most
106 frequent native languages in the target data set of non-native English speakers are Spanish ($N = 17$),
107 French ($N = 13$), and Arabic ($N = 12$). A total of 46 languages is only spoken by a single speaker.

¹ This method is an extension of the (unpublished) approach reported by Richter (2017). In her approach, she also calculated pronunciation distances using MFCCs, but compared the whole utterances directly without segmenting them into words first.

108 Human judgments of native-likeness

109 Perceptual data have been widely used to assess the degree of foreign-accentedness (Koster and Koet,
110 1993; Munro, 1995; Magen, 1998; Munro and Derwing, 2001). We therefore use human judgments of
111 native-likeness that were collected in the study of Wieling et al. (2014). They created an online questionnaire
112 in which native speakers of American-English were asked to rate the accent strength of 50 speech samples
113 extracted from the Speech Accent Archive. The degree of native-likeness of the speech samples was judged
114 on a 7-point Likert scale. A score of 1 was assigned to a speaker that was perceived as very foreign-
115 sounding, and a score of 7 was assigned to a speaker that was perceived as having native American-English
116 speaking abilities. The speech samples presented to the participants were not duplicated, so each participant
117 rated each sample at most once. The set of samples available for different participants to judge was changed
118 several times during the period the questionnaire was online. To increase the reliability of the ratings, not
119 all speech samples from the Speech Accent Archive were included in the questionnaire, so that each speech
120 sample could be judged by multiple participants. It was also not compulsory to rate all 50 samples, because
121 the participants could decide to rate a subset of the speech samples.

122 The questionnaire of Wieling et al. (2014) was distributed by asking colleagues and friends to forward
123 it to native speakers of American-English. The questionnaire was also mentioned in a blog post of Mark
124 Liberman² which led to a considerable amount of responses. In total, 1,143 participants provided native-
125 likeness ratings (57.6% men and 42.4% woman). On average, they rated 41 samples with a standard
126 deviation of 14 samples. The participants had a mean age of 36.2 years with a standard deviation of 13.9
127 years, and people most frequently came from California (13.2%), New York (10.1%), and Massachusetts
128 (5.9%).

129 Experimental setup

130 Segmentation

131 We obtain acoustic distances comparing speakers from the target data set to the speakers in the reference
132 data set. The data sets we use contain recordings of the entire 69 word paragraph (henceforth referred
133 to as the complete speech sample). These complete speech samples do not only contain the 69 word
134 pronunciations, but also speech disfluencies. Examples of these disfluencies include, but are not limited to,
135 (filled) pauses, false starts, word order changes, or mispronunciations.

136 To only compare corresponding segments of speech (the approach of Richter (2017) comparing the
137 complete speech samples was not very successful), we segment each complete speech sample into words.
138 While this segmentation procedure may be performed manually, this is very time consuming (Goldman,
139 2011). We therefore employ the Penn Phonetics Lab Forced Aligner (P2FA) to time-align the speech
140 samples with a word-level orthographic transcription (Yuan and Liberman, 2008). The P2FA is an automatic
141 phonetic alignment toolkit that is based on the Hidden Markov Toolkit (HTK). Prior to creating the forced
142 alignments, we resample each of the speech samples to 11,025 Hz (Yuan and Liberman, 2008). The forced
143 alignment approach identifies the word boundaries in the speech samples, and by using this information
144 we automatically divide the complete speech samples of the target and reference data set into separate
145 words. Each word corresponds to a word from the elicitation paragraph presented in (1). In this way, we
146 also remove non-speech elements that exist between these word boundaries, preventing them from entering
147 the acoustic distance calculation. After the forced alignment procedure, we have a target data set that for
148 each of the 280 speakers contains 69 segmented speech samples, as well as a reference data set of 115
149 speakers with for each speaker 69 corresponding segmented speech samples. A detailed explanation of

² <https://languagelog.ldc.upenn.edu/nll/?p=3967>, May 19, 2012, "Rating American English Accents."

150 the theoretical framework behind the forced alignment procedure is provided in the studies of Young and
151 Young (1993) and Bailey (2016).

152 Feature representation

153 For each segmented speech sample in both data sets, we calculate a numerical feature representation
154 based on Mel-frequency cepstral coefficients (MFCCs). MFCCs have shown their robustness, as these
155 speech features are widely used as representations of phonetic content in automatic speech recognition
156 systems (Davis and Mermelstein, 1980).

157 We visualize the computation of each MFCC feature representation in **Figure 1**. The first, commonly
158 executed, step in calculating this numerical feature representation is to compensate for the negative spectral
159 slope of each speech sample (Sluijter and Van Heuven, 1996). The nature of the glottal pulses causes
160 voiced segments in the audio signal to contain more energy at the lower frequencies compared to the higher
161 frequencies (Vergin and O'Shaughnessy, 1995). We remove some of these glottal effects from the spectrum
162 of the vocal tract by applying a filter to the audio signal (see equation 1). This filter emphasizes the higher
163 frequencies, and as a result a more balanced spectrum of the speech sample is obtained. This is usually
164 referred to as the pre-emphasis step (Muda et al., 2010).

$$H(z) = 1 - 0.97 * z^{-1} \quad (1)$$

165 We then divide each speech sample into short frames of time using a windowing function. These frames
166 of analysis are important since the characteristics of an audio signal are fairly stable when a short frame
167 of time is taken into account (Zhu and Alwan, 2000). We create overlapping frames of a 25 millisecond
168 time interval using a 10 millisecond step size. A set of cepstral coefficients is computed for each of these
169 windowed frames per speech sample. The Hamming windowing function is used to extract each frame
170 from the audio signal (Deller Jr et al., 1993).

171 The Discrete Fourier Transform (DFT) is then taken from each of these windowed frames to transform
172 the audio signal from the time domain to the frequency domain (Zheng et al., 2001). Taking the DFT of the
173 windowed frames is related to the way sound is perceived by human beings. The oscillation of the human
174 cochlea depends on the frequency of incoming sounds, and these oscillations inform the human brain that
175 certain frequencies are present in the audio signal. With the application of DFT, the process that occurs
176 within the human auditory system is simulated (Dave, 2013).

177 After the DFT is taken from the windowed frames, the Mel spectrum is computed. The DFT-transformed
178 audio signal is modified by passing it through a collection of triangular band-pass filters. These filters are
179 also known as the Mel filter bank, and each processes frequencies that occur within a certain range while
180 discarding frequencies that are outside that range (Muda et al., 2010). The Mel filter bank then provides
181 information about the amount of energy that is present near certain frequency regions (Rao and Manjunath,
182 2017). The width of the filter banks is determined via Mel-scaling. Units on the Mel scale are based on
183 the way frequencies are perceived by the human auditory system. These Mel units do not correspond to
184 tone frequencies in a linear way, as the human auditory system does not perceive frequencies linearly.
185 Instead, the Mel scale is composed such that the frequencies below 1,000 Hertz are approximately linearly
186 spaced, and the frequencies above 1,000 Hertz are distributed according to a logarithmic scale (Stevens
187 et al., 1937).

188 The first filters of the Mel-filter bank are most strict, since the low frequencies are the most informative
189 in speech perception (Raut and Shah, 2015). The energy of voiced speech is mostly concentrated at

190 the lower frequencies (Seltzer et al., 2004). After the DFT-transformed audio signal goes through the
 191 triangular-shaped band-pass filters, the logarithm is taken of the energies that are returned by the Mel-filter
 192 bank. This procedure is also in accordance with the human auditory system, since humans do not perceive
 193 the loudness of an incoming audio signal linearly. The final result of this procedure is a signal that is
 194 represented in the cepstral domain (Oppenheim and Schaffer, 2004).

195 The logarithmically transformed filter bank energy representations do, however, overlap. To provide
 196 a solution to the overlapping filter banks, the discrete cosine transform (DCT) is computed from the
 197 logarithmically transformed filter bank output. The result of the DCT is a set of cepstral coefficients.
 198 Following an established standard, we chose to solely include the first 12 cepstral coefficients and energy
 199 in each frame, which characterise the most relevant information of the speech signal (Picone, 1993). In
 200 addition, we calculate the first-order and second-order derivatives from each of the cepstral coefficients and
 201 energy features (Furui, 1981). We therefore have 12 first-order and 12 second-order derivatives that are
 202 associated with the 12 cepstral coefficients, and one first-order and second-order derivative related to the
 203 energy feature. These first-order and second-order derivatives, or (double) delta coefficients, model the
 204 changes between the frames over time (Muda et al., 2010). A total of 39 coefficients is computed at each
 205 10 millisecond step per speech sample, to represent the most important phonetic information embedded
 206 within each 25 millisecond windowed frame. The MFCC feature representation per segmented speech
 207 sample is obtained by concatenating its corresponding vectors of 39 coefficients computed for each of the
 208 windowed frames.

209 Normalization

210 Ganapathy et al. (2011) and Shafik et al. (2009) showed that the quality of the MFCC feature
 211 representation is highly influenced by the presence of noise in the speech samples. To reduce the effect of
 212 noise, cepstral mean and variance normalization is applied to the feature representations (Auckenthaler
 213 et al., 2000). In addition to the robustness in the presence of noisy input, cepstral mean and variance
 214 normalization reduces the word error rate in automatic speech recognition implementations, and improves
 215 the generalization across speakers (Haeb-Umbach, 1999; Molau et al., 2003; Tsakalidis and Byrne, 2005).
 216 Adank et al. (2004) showed that cepstral mean and variance normalization can be used to highlight the
 217 linguistic content of the feature representations.

218 We implement cepstral mean and variance normalization by applying a linear transformation to the
 219 coefficients of the MFCC feature representations (Lu et al., 2009). The MFCC feature representations are
 220 standardized per speaker by removing the speaker's mean, and scaling to unit variance. The equation that
 221 we use to calculate the cepstral mean and variance normalized feature representations is shown in equation
 222 2.

$$\hat{c}(i, t) = \frac{c(i, t) - \bar{c}(i, t)}{\sigma(i)} \quad (2)$$

223 In this equation, the i -th cepstral coefficient at time index t is represented by $c(i, t)$. The mean value of each
 224 feature representation, and the corresponding standard deviation are given by $\bar{c}(i, t)$ and $\sigma(i)$, respectively.
 225 In equations 3 and 4 we show how the mean value and standard deviation are obtained. In these equations,
 226 N corresponds to the number of windows used in processing the speech sample.

$$\bar{c}(i, t) = \frac{1}{N} * \sum_{t=1}^N c(i, t) \quad (3)$$

$$\sigma(i) = \sqrt{\frac{1}{N} * \sum_{t=1}^N (c(i, t) - \bar{c}(i, t))^2} \quad (4)$$

227 Dynamic time warping

228 The acoustic word distances are computed using the dynamic time warping (DTW) algorithm. This
 229 algorithm compares two MFCC feature representations, and returns their degree of similarity as a distance
 230 score (Galbally and Galbally, 2015). DTW has already been widely used in the domain of speech
 231 recognition, and is also used for sequence comparison in many other research domains, such as computer
 232 vision and protein structure matching (Sakoe et al., 1990; Bahlmann and Burkhardt, 2004; Efrat et al.,
 233 2007).

234 To compare a target pronunciation with a reference pronunciation, the DTW algorithm uses the
 235 corresponding target and reference MFCC feature representations. These are shown in equations 5 and 6

$$\text{target} = (x_1, x_2, \dots, x_n) \quad (5)$$

236

$$\text{reference} = (y_1, y_2, \dots, y_m) \quad (6)$$

237 An $m * n$ cost matrix is created to align the target MFCC feature representation with the reference MFCC
 238 feature representation (Muda et al., 2010). This cost matrix is filled with the Euclidean distances between
 239 every pair of points (frames) in both the target and reference MFCC feature representations (Danielsson,
 240 1980). For example, element (i, j) of the cost matrix contains the distance d that is given by equation 7

$$d(\text{target}_i, \text{reference}_j) = (\text{target}_i - \text{reference}_j)^2 \quad (7)$$

241 The optimal alignment between the MFCC feature representations corresponds to the shortest path
 242 through the cost matrix, and is therefore to some extent comparable to the edit distance. The DTW
 243 algorithm computes the shortest path using an iterative method that calculates the minimum cumulative
 244 distance $\gamma(i, j)$ (Keogh and Pazzani, 2001). The cumulative distance is composed of the distance in the
 245 current cell $d(\text{target}_i, \text{reference}_j)$ and the minimum of the cumulative distance found in the adjacent cells
 246 (shown in equation 8).

$$\gamma(i, j) = d(\text{target}_i, \text{reference}_j) + \min(\gamma(i - 1, j - 1), \gamma(i - 1, j), \gamma(i, j - 1)) \quad (8)$$

247 After the cumulative distance is computed, it is divided by the length of the target feature representation
 248 and the reference feature representation $(n + m)$. It is important to normalize the computed distances,
 249 since the speech samples we work with do not necessarily have the same length. Without normalization
 250 applied to DTW, longer alignment paths (from longer recordings) would have higher distances than shorter
 251 alignments, because they have more frames to accumulate cost (Giorgino et al., 2009).

252 The final speaker pronunciation distances are obtained by first calculating the acoustic distance for each
 253 of the 69 words pronounced by a non-native speaker of American-English and a single native speaker
 254 of American-English in the reference data set. We subsequently average these word-based distances to
 255 measure the between-speaker acoustic distance. The difference between the pronunciation of a non-native
 256 speaker and native American-English in general, is determined by calculating the between-speaker acoustic

257 distances compared to all 115 native American-English speakers, and subsequently averaging these. We
258 compute these acoustic distances for all foreign-accented speech samples by applying this same procedure
259 to each of the 280 non-native speakers of American-English in the target data set. To evaluate our measure,
260 the correlation between the native-likeness ratings and the acoustic distances is computed. We evaluate the
261 impact of the (size of the) set of reference speakers, by calculating the correlation for successively smaller
262 subsets of reference speakers.

263 Understanding the acoustic distance measure

264 In addition to the main experiment, we perform a variety of other analyses to obtain a more complete
265 understanding of the acoustic details captured by the acoustic distance measure.

266 First we use a multiple linear regression model to predict the human native-likeness ratings on the basis
267 of our acoustic distance measure, but also using the transcription-based distances reported by [Wieling](#)
268 [et al. \(2014\)](#), and the (manually counted) number of mispronunciations a speaker made, as these might be
269 important for native-likeness ratings [\(Flege, 1981\)](#), but are not included in either of the two other measures.

270 Second, to assess whether our acoustic distance measure adequately captures fine-grained segmental
271 differences, we compute acoustic differences between ten repetitions of hVd words (e.g., [hid]) pronounced
272 by a single speaker. We subsequently correlate these differences with differences based on the first and
273 second formant measured at the mid-point of the vowel of the recordings. We follow [Wieling et al. \(2012\)](#)
274 in Bark-scaling the formant-based distances. We use a total of 12 Dutch monophthongs in the vowel context
275 (a, ɑ, ɛ, e, ø, i, i, ɔ, u, o, ʏ, y). We visualize the differences (both the formant-based distances, and the
276 acoustic-based distances) using multidimensional scaling [\(Torgerson, 1952\)](#).

277 Third and finally, to assess whether non-segmental variability is also captured by our acoustic method,
278 we compute acoustic distances between four series of recordings (ten repetitions) of the word ‘living’. The
279 first and second series consisted of a normal pronunciation (‘living’), but recorded with two recording
280 devices (the built-in microphone of a laptop, and the built-in microphone of a smartphone), the third series
281 consisted of a pronunciation in which the intonation was changed (‘living?’), and the fourth series consisted
282 of a pronunciation in which the relative duration of the syllables was changed (‘li_ving’).

RESULTS

283 The correlation between the native-likeness ratings and the acoustic distances computed using our acoustic
284 method is $r = -0.71$ ($p < 0.0001$), and therefore accounts for about half of the variance in the native-
285 likeness ratings ($r^2 = 0.50$). **Figure 2** visualizes this correlation in a scatter plot. The acoustic distance
286 measure tends to underestimate the native-likeness (overestimate distances) when the speech samples are
287 rated as being very native-like.

288 Compared to the transcription-based method of [Wieling et al. \(2014\)](#), who used the Levenshtein distance
289 incorporating automatically determined linguistically-sensible segment distances, and reported a correlation
290 of $r = -0.77$, the performance of our measure is significantly lower (using the modified z -statistic of
291 [Steiger \(1980\)](#): $z = 2.10$, $p < 0.05$).

292 Impact of reference speakers

293 As the set of reference speakers might affect the correlation, we evaluated the impact of reducing the set
294 of reference speakers. The results are shown in **Table 1** and show that the correlation remains comparable,
295 irrespective of the (size of the) reference set (i.e. $-0.68 \leq r \leq -0.72$). To assess whether language
296 variation within the set of reference speakers might be important, we computed the acoustic distances
297 using as our reference set ($N = 14$) only the native English speakers who originated from the western

298 half of the U.S. and the English-speaking part of Canada. These areas are characterized by less dialect
299 variation compared to the eastern half of the U.S. (Boberg, 2010). Again, this did not substantially affect
300 the correlation, as it remained similar ($r = -0.70$).

301 Impact of segmentation and normalization

302 Two simplified (baseline) measures, each missing a single component of our acoustic measure, were
303 created to assess how segmentation and cepstral mean and variance normalization of the speech samples
304 contribute to acoustic distances that are more similar to human judgments of native-likeness. The results of
305 this experiment is shown in Table 2. It is clear that not using the normalization approach is much more
306 detrimental than not segmenting, but that the best results are obtained when doing both. The modified
307 z -statistic of Steiger (1980) indicates that our acoustic method significantly outperforms either of the two
308 simpler methods ($z = 4.11$, $p < 0.0001$).

309 Understanding the acoustic distance measure

310 We fitted a multiple linear regression model to determine whether the acoustic distance measure and the
311 transcription-based distance measure captured distinctive aspects of pronunciation. We also assessed
312 the influence of the number of mispronunciations. The coefficients and associated statistics of the
313 predictors used are shown in Table 3. The results show that the transcription-based distances and acoustic
314 distances both contribute significantly to the model fit ($p < 0.05$). This is not the case for the amount of
315 mispronunciations per speaker in the target data set ($p > 0.05$). The presented model accounts for 65% of
316 the variation in the human judgments of native-likeness ($r^2 = 0.65$). Only using the transcription-based
317 distance measure accounted for 60% of the variation. Consequently, our acoustic measure also seems to
318 capture information which is not present in phonetic transcriptions.

319 The results in Table 4 show that our acoustic measure can capture both intonation and timing differences
320 as these lead to larger distances than comparing individual repetitions of the same word pronounced by the
321 same speaker. However, it also shows that when recording the normal pronunciation by two microphones
322 simultaneously, the acoustic distances between the two simultaneous recordings are higher than zero,
323 whereas the pronunciation is in fact identical. Note that the acoustic distance when comparing the 10
324 normal pronunciations is also not zero, due to small deviations in the pronunciations.

325 Another indication of how well our acoustic measure captures segmental information is shown by the
326 significant positive correlation of $r = 0.68$ ($p < 0.0001$) between the formant-based acoustic vowel
327 differences and the computed acoustic differences between the hVd-words. Figure 3 shows these relative
328 vowel distances by using a multidimensional scaling visualization of the formant-based vowel differences
329 (visualizing all variation) and the DTW-based vowel differences (visualizing 47% of the variation in the
330 original differences).

DISCUSSION

331 We have created an acoustic-only approach for calculating pronunciation distances between utterances
332 of the same word by different speakers. We have evaluated the measure by calculating how different the
333 speech of non-native speakers of American-English is from native American-English speakers, and by
334 comparing our computed results to human judgments of native-likeness. While our method is somewhat
335 outperformed ($r = -0.71$ vs. $r = -0.77$) by the transcription-based method introduced by Wieling et al.
336 (2014), our measure does not require phonetic transcriptions, whose production is time consuming and
337 prone to errors. Given that our method is fully automatic, the trade-off in performance may be worthwhile.

338 Word segmentation and especially speaker-based cepstral mean and variance normalization of the
339 MFCC speech representations were important in creating an adequate acoustic-based distance measure.

340 These results show the importance of pre-processing continuous speech samples, as the comparison of
341 pronunciations in speech samples is most reliable when it is based on comparable and normalized segments
342 of speech that we obtain from word-level forced alignment.

343 The multiple regression model showed that the acoustic distance explained variance not accounted for by
344 the transcription-based distance measure. Particularly, our further experiments showed that our measure
345 is both sensitive to timing and intonation differences. However, the measure is also sensitive to different
346 recording devices, which is undesirable and may partly explain why the method is outperformed by the
347 transcription-based method. While the MFCC feature representation with cepstral mean and variance
348 normalization attempts to minimise non-linguistic confounds, it is only partly successful, as a computational
349 representation of general phonetic information remains a difficult issue in speech processing technology
350 (Gemmeke et al., 2011).

351 Consequently, future work should investigate whether other acoustic (pre-processing) techniques may
352 improve our acoustic measure. For example, contextual acoustic encoding techniques related to word
353 embeddings like *wav2vec* and *vq-wav2vec* may highlight acoustic details that are linguistically relevant
354 (Baeovski et al., 2019; Schneider et al., 2019). Additionally, generating a shared phonetic space through
355 which two speech samples may be compared (Ryant and Liberman, 2016) may be useful. Nevertheless, our
356 work serves as a useful and promising starting point for a fully automatic acoustic pronunciation distance
357 measure.

AUTHOR CONTRIBUTIONS

358 MB and MW conceptualized the research. MB and CR designed and conducted the experiments. MB
359 and CR performed data analysis and data visualization. MB wrote the first version of the manuscript. All
360 authors contributed to manuscript revision, read, and approved the submitted version.

CONFLICT OF INTEREST STATEMENT

361 The authors confirm that the research was conducted in the absence of any commercial or financial
362 relationship that could lead to a potential conflict of interest.

FUNDING DISCLOSURES

363 The research discussed in this paper has been funded by the Center for Groningen Language and Culture,
364 Faculty of Arts, University of Groningen. CR is supported by NSF grant DGE-1321851.

ACKNOWLEDGEMENTS

365 The authors would like to express their gratitude towards Hedwig Sekeres for her contribution visualizing
366 the formant-based vowel distances and assessing the amount the mispronunciations in the data set.

REFERENCES

- 367 Adank, P., Smits, R., and Van Hout, R. (2004). A comparison of vowel normalization procedures for
368 language variation research. *The Journal of the Acoustical Society of America* 116, 3099–3107
- 369 Arslan, L. M. and Hansen, J. H. (1996). Language accent classification in american english. *Speech*
370 *Communication* 18, 353–367
- 371 Arslan, L. M. and Hansen, J. H. (1997). A study of temporal features and frequency characteristics in
372 american english foreign accent. *The Journal of the Acoustical Society of America* 102, 28–40
- 373 Asher, J. J. and García, R. (1969). The optimal age to learn a foreign language. *The Modern Language*
374 *Journal* 53, 334–341

- 375 Auckenthaler, R., Carey, M., and Lloyd-Thomas, H. (2000). Score normalization for text-independent
376 speaker verification systems. *Digital Signal Processing* 10, 42–54
- 377 Baevski, A., Schneider, S., and Auli, M. (2019). vq-wav2vec: Self-supervised learning of discrete speech
378 representations. *arXiv preprint arXiv:1910.05453*
- 379 Bahlmann, C. and Burkhardt, H. (2004). The writer independent online handwriting recognition system
380 frog on hand and cluster generative statistical dynamic time warping. *IEEE Transactions on Pattern
381 Analysis and Machine Intelligence* 26, 299–310
- 382 Bailey, G. (2016). Automatic detection of sociolinguistic variation using forced alignment. In *University
383 of Pennsylvania Working Papers in Linguistics: Selected Papers from New Ways of Analyzing Variation
384 (NWAV 44)* (York), 10–20
- 385 Boberg, C. (2010). *The English language in Canada: Status, history and comparative analysis* (Cambridge
386 University Press)
- 387 Bucholtz, M. (2007). Variation in transcription. *Discourse Studies* 9, 784–808
- 388 Chen, N. F., Wee, D., Tong, R., Ma, B., and Li, H. (2016). Large-scale characterization of non-native
389 mandarin chinese spoken by speakers of european origin: Analysis on icall. *Speech Communication* 84,
390 46–56
- 391 Cucchiaroni, C. (1996). Assessing transcription agreement: methodological aspects. *Clinical Linguistics &
392 Phonetics* 10, 131–155
- 393 Danielsson, P.-E. (1980). Euclidean distance mapping. *Computer Graphics and image processing* 14,
394 227–248
- 395 Dave, N. (2013). Feature extraction methods lpc, plp and mfcc in speech recognition. *International journal
396 for advance research in engineering and technology* 1, 1–4
- 397 Davis, S. and Mermelstein, P. (1980). Comparison of parametric representations for monosyllabic word
398 recognition in continuously spoken sentences. *IEEE transactions on acoustics, speech, and signal
399 processing* 28, 357–366
- 400 Deller Jr, J. R., Proakis, J. G., and Hansen, J. H. (1993). *Discrete time processing of speech signals*
401 (Prentice Hall PTR)
- 402 Derwing, T. M. and Munro, M. J. (2009). Putting accent in its place: Rethinking obstacles to communication.
403 *Language teaching* 42, 476–490
- 404 Duckworth, M., Allen, G., Hardcastle, W., and Ball, M. (1990). Extensions to the international phonetic
405 alphabet for the transcription of atypical speech. *Clinical Linguistics & Phonetics* 4, 273–280
- 406 Efrat, A., Fan, Q., and Venkatasubramanian, S. (2007). Curve matching, time warping, and light fields:
407 New algorithms for computing similarity between curves. *Journal of Mathematical Imaging and Vision*
408 27, 203–216
- 409 Flege, J. E. (1981). The phonological basis of foreign accent: A hypothesis. *Tesol Quarterly* 15, 443–455
- 410 Flege, J. E. (1988). Factors affecting degree of perceived foreign accent in english sentences. *The Journal
411 of the Acoustical Society of America* 84, 70–79
- 412 Furui, S. (1981). Comparison of speaker recognition methods using statistical features and dynamic
413 features. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 29, 342–350
- 414 Galbally, J. and Galbally, D. (2015). A pattern recognition approach based on dtw for automatic transient
415 identification in nuclear power plants. *Annals of Nuclear Energy* 81, 287–300
- 416 Ganapathy, S., Pelecanos, J., and Omar, M. K. (2011). Feature normalization for speaker verification in
417 room reverberation. In *2011 IEEE International Conference on Acoustics, Speech and Signal Processing
418 (ICASSP) (IEEE)*, 4836–4839

- 419 Gemmeke, J. F., Virtanen, T., and Hurmalainen, A. (2011). Exemplar-based sparse representations for noise
420 robust automatic speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*
421 19, 2067–2080
- 422 Giorgino, T. et al. (2009). Computing and visualizing dynamic time warping alignments in r: the dtw
423 package. *Journal of statistical Software* 31, 1–24
- 424 Goldman, J.-P. (2011). Easyalign: an automatic phonetic alignment tool under praat. In *Proceedings*
425 *of the Annual Conference of the International Speech Communication Association, INTERSPEECH*.
426 3233–3236
- 427 Gooskens, C. and Heeringa, W. (2004). Perceptive evaluation of levenshtein dialect distance measurements
428 using norwegian dialect data. *Language variation and change* 16, 189–207
- 429 Goslin, J., Duffy, H., and Floccia, C. (2012). An erp investigation of regional and foreign accent processing.
430 *Brain and language* 122, 92–102
- 431 Haeb-Umbach, R. (1999). Investigations on inter-speaker variability in the feature space. In *1999 IEEE*
432 *International Conference on Acoustics, Speech, and Signal Processing. Proceedings. ICASSP99 (Cat.*
433 *No. 99CH36258)* (IEEE), vol. 1, 397–400
- 434 Hakkani-Tür, D., Riccardi, G., and Gorin, A. (2002). Active learning for automatic speech recognition.
435 In *2002 IEEE International Conference on Acoustics, Speech, and Signal Processing* (IEEE), vol. 4,
436 IV–3904
- 437 Heeringa, W., Johnson, K., and Gooskens, C. (2009). Measuring norwegian dialect distances using acoustic
438 features. *Speech Communication* 51, 167–183
- 439 Heeringa, W. J. (2004). *Measuring dialect pronunciation differences using Levenshtein distance*. Ph.D.
440 thesis, Citeseer
- 441 Jeszenszky, P., Stoeckle, P., Glaser, E., and Weibel, R. (2017). Exploring global and local patterns in the
442 correlation of geographic distances and morphosyntactic variation in swiss german. *Journal of Linguistic*
443 *Geography* 5, 86–108
- 444 Keogh, E. J. and Pazzani, M. J. (2001). Derivative dynamic time warping. In *Proceedings of the 2001*
445 *SIAM international conference on data mining* (SIAM), 1–11
- 446 Kessler, B. (1995). Computational dialectology in irish gaelic. In *Proceedings of the seventh conference on*
447 *European chapter of the Association for Computational Linguistics* (Morgan Kaufmann Publishers Inc.),
448 60–66
- 449 Koster, C. J. and Koet, T. (1993). The evaluation of accent in the english of dutchmen. *Language learning*
450 43, 69–92
- 451 Leather, J. (1983). Second-language pronunciation learning and teaching. *Language Teaching* 16, 198–219
- 452 Livescu, K. and Glass, J. (2000). Lexical modeling of non-native speech for automatic speech recognition.
453 In *2000 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat.*
454 *No. 00CH37100)* (IEEE), vol. 3, 1683–1686
- 455 Lu, X., Matsuda, S., Unoki, M., Shimizu, T., and Nakamura, S. (2009). Temporal contrast normalization
456 and edge-preserved smoothing on temporal modulation structure for robust speech recognition. In *2009*
457 *IEEE International Conference on Acoustics, Speech and Signal Processing* (IEEE), 4573–4576
- 458 Magen, H. S. (1998). The perception of foreign-accented speech. *Journal of phonetics* 26, 381–400
- 459 Mermelstein, P. (1976). Distance measures for speech recognition, psychological and instrumental. *Pattern*
460 *recognition and artificial intelligence* 116, 374–388
- 461 Molau, S., Hilger, F., and Ney, H. (2003). Feature space normalization in adverse acoustic
462 conditions. In *2003 IEEE International Conference on Acoustics, Speech, and Signal Processing,*
463 *2003. Proceedings.(ICASSP'03)*. (IEEE), vol. 1, I–I

- 464 Muda, L., Begam, M., and Elamvazuthi, I. (2010). Voice recognition algorithms using mel frequency
465 cepstral coefficient (mfcc) and dynamic time warping (dtw) techniques. *arXiv preprint arXiv:1003.4083*
- 466 Munro, M. J. (1995). Nonsegmental factors in foreign accent: Ratings of filtered speech. *Studies in Second*
467 *Language Acquisition* 17, 17–34
- 468 Munro, M. J. and Derwing, T. M. (2001). Modeling perceptions of the accentedness and comprehensibility
469 of l2 speech the role of speaking rate. *Studies in second language acquisition* 23, 451–468
- 470 Nerbonne, J. and Heeringa, W. (1997). Measuring dialect distance phonetically. In *Computational*
471 *Phonology. Third Meeting of the ACL Special Interest Group in Computational Phonology*. (Association
472 for Computational Linguistics (ACL)), 11–18
- 473 Novotney, S. and Callison-Burch, C. (2010). Cheap, fast and good enough: Automatic speech recognition
474 with non-expert transcription. In *Human Language Technologies: The 2010 Annual Conference*
475 *of the North American Chapter of the Association for Computational Linguistics* (Association for
476 Computational Linguistics), 207–215
- 477 Oppenheim, A. V. and Schaffer, R. W. (2004). From frequency to quefrequency: A history of the cepstrum.
478 *IEEE signal processing Magazine* 21, 95–106
- 479 Picone, J. W. (1993). Signal modeling techniques in speech recognition. *Proceedings of the IEEE* 81,
480 1215–1247
- 481 Piske, T., MacKay, I. R., and Flege, J. E. (2001). Factors affecting degree of foreign accent in an l2: A
482 review. *Journal of phonetics* 29, 191–215
- 483 Rao, K. S. and Manjunath, K. (2017). *Speech recognition using articulatory and excitation source features*
484 (Springer)
- 485 Raut, S. P. and Shah, D. S. N. (2015). Voice biometric system for speaker authentication. *International*
486 *Journal of Computer Applications* 975, 8887
- 487 Richter, C. (2017). Accent. Unpublished manuscript
- 488 Ryant, N. and Liberman, M. (2016). Large-scale analysis of spanish/s/-lenition using audiobooks. In
489 *Proceedings of Meetings on Acoustics 22ICA* (ASA), vol. 28, 060005
- 490 Sakoe, H., Chiba, S., Waibel, A., and Lee, K. (1990). Dynamic programming algorithm optimization for
491 spoken word recognition. *Readings in speech recognition* 159, 224
- 492 Schneider, S., Baeviski, A., Collobert, R., and Auli, M. (2019). wav2vec: Unsupervised pre-training for
493 speech recognition. *arXiv preprint arXiv:1904.05862*
- 494 Seltzer, M. L., Raj, B., and Stern, R. M. (2004). A bayesian classifier for spectrographic mask estimation
495 for missing feature speech recognition. *Speech Communication* 43, 379–393
- 496 Shafik, A., Elhalafawy, S. M., Diab, S., Sallam, B. M., and El-Samie, F. A. (2009). A wavelet based
497 approach for speaker identification from degraded speech. *International Journal of Communication*
498 *Networks and Information Security* 1, 52
- 499 Sluijter, A. M. and Van Heuven, V. J. (1996). Spectral balance as an acoustic correlate of linguistic stress.
500 *The Journal of the Acoustical society of America* 100, 2471–2485
- 501 Steiger, J. H. (1980). Tests for comparing elements of a correlation matrix. *Psychological bulletin* 87, 245
- 502 Stevens, S. S., Volkman, J., and Newman, E. B. (1937). A scale for the measurement of the psychological
503 magnitude pitch. *The Journal of the Acoustical Society of America* 8, 185–190
- 504 Torgerson, W. S. (1952). Multidimensional scaling: I. theory and method. *Psychometrika* 17, 401–419
- 505 Tsakalidis, S. and Byrne, W. (2005). Acoustic training from heterogeneous data sources: Experiments
506 in mandarin conversational telephone speech transcription. In *Proceedings.(ICASSP'05). IEEE*
507 *International Conference on Acoustics, Speech, and Signal Processing, 2005.* (IEEE), vol. 1, I–461

- 508 Vergin, R. and O'Shaughnessy, D. (1995). Pre-emphasis and speech recognition. In *Proceedings 1995*
509 *Canadian Conference on Electrical and Computer Engineering*. vol. 2, 1062–1065 vol.2. doi:10.1109/
510 CCECE.1995.526613
- 511 [Dataset] Weinberger, S. (2015). Speech accent archive
- 512 Weinberger, S. H. and Kunath, S. A. (2011). The speech accent archive: towards a typology of english
513 accents. *Language and Computers-Studies in Practical Linguistics* 73, 265
- 514 Wieling, M., Bloem, J., Mignella, K., Timmermeister, M., and Nerbonne, J. (2014). Measuring foreign
515 accent strength in english: Validating levenshtein distance as a measure. *Language Dynamics and*
516 *Change* 4, 253–269
- 517 Wieling, M., Margaretha, E., and Nerbonne, J. (2012). Inducing a measure of phonetic similarity from
518 pronunciation variation. *Journal of Phonetics* 40, 307–314
- 519 Wieling, M., Nerbonne, J., and Baayen, R. H. (2011). Quantitative social dialectology: Explaining
520 linguistic variation geographically and socially. *PloS one* 6
- 521 Young, S. J. and Young, S. (1993). *The HTK hidden Markov model toolkit: Design and philosophy*
522 (University of Cambridge, Department of Engineering Cambridge, England)
- 523 Yuan, J. and Liberman, M. (2008). Speaker identification on the scotus corpus. *Journal of the Acoustical*
524 *Society of America* 123, 3878
- 525 Zheng, F., Zhang, G., and Song, Z. (2001). Comparison of different implementations of mfcc. *Journal of*
526 *Computer science and Technology* 16, 582–589
- 527 Zhu, Q. and Alwan, A. (2000). On the use of variable frame rate analysis in speech recognition. In *2000*
528 *IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No.*
529 *00CH37100)* (IEEE), vol. 3, 1783–1786

Table 1. Pearson correlation coefficients r between the acoustic distances and human judgments of native-likeness depending on the size of the reference data set. All correlations are significant at the $p < 0.0001$ level.

Amount of reference speakers	r
10	-0.68
25	-0.71
50	-0.70
75	-0.72

Table 2. Pearson correlation coefficients r of acoustic distances compared to human judgments of native-likeness, using different methods to compute the acoustic distances. All correlations are significant at the $p < 0.0001$ level.

Model	r
Baseline 1 (only segmentation)	-0.27
Baseline 2 (only normalization)	-0.63
Acoustic measure (segmentation and normalization)	-0.71

Table 3. Coefficients of a multiple regression model predicting human judgments of native-likeness.

	Estimate	Std. Error	t -value	p -value
Intercept	24.19	2.68	9.04	< 0.001
Transcription-based distances	-379.30	34.26	-11.07	< 0.001
Acoustic-based distances	-2.79	0.44	-6.35	< 0.001
Amount of mispronunciations	0.01	0.03	0.26	0.795

Table 4. Averaged acoustic distances and standard errors of four variants of the word 'living'.

	Compared to normal pronunciation
Normal pronunciation	4.35 (0.50)
Normal pronunciation (different recording device)	6.94 (0.15)
Rising intonation	7.12 (0.13)
Lengthened first syllable	6.65 (0.13)

530 **Figure 1.** Diagram visualizing the features used in our acoustic distance algorithm.

531 **Figure 2.** Native-likeness ratings as a function of the computed acoustic distances ($r = -0.71$).

532 **Figure 3.** MDS plots visualizing the acoustic vowel distances (left) and the formant-based vowel distances
 533 (right). Individual pronunciations are shown in light gray, whereas the averages per vowel are shown in
 534 black.



