1 **A cognitively grounded measure of pronunciation distance**

2 Martijn Wieling[a*], John Nerbonne[b], Jelke Bloem[c], Charlotte Gooskens[b], Wilbert

3 Heeringa[b] and R. Harald Baayen[a,d]

4 [a]Department of Quantitative Linguistics, University of Tübingen, Tübingen, Germany [b]Center for

5 Language and Cognition Groningen, University of Groningen, Groningen, The Netherlands

6 [c]Amsterdam Center of Language and Communication, University of Amsterdam, Amsterdam, The

7 Netherlands [d]Department of Linguistics, University of Alberta, Edmonton, Alberta, Canada

8 *wieling@gmail.com

9

10 **Abstract**

11 In this study we develop pronunciation distances based on naive discriminative learning (NDL).

12 Measures of pronunciation distance are used in several subfields of linguistics, including

13 psycholinguistics, dialectology and typology. In contrast to the commonly used Levenshtein algorithm,

14 NDL is grounded in cognitive theory of competitive reinforcement learning and is able to generate

15 asymmetrical pronunciation distances. In a first study, we validated the NDL-based pronunciation

16 distances by comparing them to a large set of native-likeness ratings given by native American English

17 speakers when presented with accented English speech. In a second study, the NDL-based

18 pronunciation distances were validated on the basis of perceptual dialect distances of Norwegian

19 speakers. Results indicated that the NDL-based pronunciation distances matched perceptual distances

20 reasonably well with correlations ranging between 0.7 and 0.8. While the correlations were comparable

21 to those obtained using the Levenshtein distance, the NDL-based approach is more flexible as it is also

22 able to incorporate acoustic information other than sound segments.

23

27

## Introduction

Obtaining a suitable distance measure between two pronunciations is important, not only for dialectologists who are interested in finding the relationship between different dialects (e.g., [1]), but also for sociolinguists investigating the effect of political borders on vernacular speech [2], language researchers investigating the typological and genealogical relationships among the world's languages (e.g., [3]), applied linguists attempting to gauge the degree of comprehensibility among related languages [4], and researchers measuring the atypicality of the speech of the bearers of cochlear implants [5]. Furthermore, having a distance measure between word pronunciations enables quantitative analyses in which the integrated effect of geography and sociolinguistic factors can be investigated (e.g., [6]). Standard sociolinguistic analyses focus on whether specific categorical differences are present in the speech of people from different social groups. By using a *measure* of pronunciation difference, we allow more powerful numerical analysis techniques to be used. For these analyses to be meaningful, however, the measurements of pronunciation distance need to match perceptual distances as closely as possible.

There are various computational methods to measure word or pronunciation distance (or similarity), of which the Levenshtein distance has been the most popular [1,7,8,9,10]. The Levenshtein distance determines the pronunciation distance between two transcribed strings by calculating the number of substitutions, insertions and deletions to transform one string into the other [11]. For example, the Levenshtein distance between two accented pronunciations of the word Wednesday, [wɛnzdeɪ] and [wɛnəsde] is 3 as illustrated by the alignment in Table 1.

53   A clear drawback of this variant of the Levenshtein distance is that it does not

54   distinguish the substitution of similar sounds (such as [o] and [u]) from more different

55   sounds (such as [o] and [i]). Consequently, effort has been made to integrate more

56   sensitive segment distances in the Levenshtein distance algorithm [1,12]. As manually

57   determining sensitive segment distances is time-consuming and language-dependent,

58   Wieling and colleagues [13] developed an automatic method to determine sensitive

59   segment distances. Their method calculated the pointwise mutual information

60   between two segments, assigning lower distances between segments which aligned

61   relatively frequently and higher distances between segments which aligned relatively

62   infrequently. Results indicated that the obtained segment distances were acoustically

63   sensible and resulted in improved alignments [14]. Applying the adapted method to

64   the example alignment shown above yields the associated costs shown in Table 2.

65

66   While Levenshtein distances correlate well ($r = 0.67$) with perceptual dialect distances

67   between Norwegian dialects [15], there is no cognitive basis to link the Levenshtein

68   distance to perceptual distances (but see [16] for an attempt to adapt the Levenshtein

69   algorithm in line with theories about spoken word recognition). This is also

70   exemplified by the fact that the Levenshtein distance is symmetrical (i.e. the distance

71   between speaker A and B is the same as the other way around), while perceptual

72   dialect distances may also show an asymmetrical pattern [15].

73

74   As exposure to language shapes expectations and affects what is judged similar to

75   one's own pronunciation and what is different, we turn to one of the most influential

76   theories about animal and human (discrimination) learning: the model of Rescorla and

77   Wagner [17]. The basic assumption of this model is that a learner predicts an outcome

3

78    (e.g., the meaning of a word) based on the set of available cues (e.g., the sounds of a

79    word). Depending on the correctness of the prediction, the association strengths

80    between the outcome and the cues are adjusted so that future prediction accuracy

81    improves. Concretely, if an outcome is present together with a certain cue, its

82    association strength increases, while the association strength between an absent

83    outcome and that cue decreases. When an outcome is found together with multiple

84    cues (i.e. when there is cue competition), the adjustments are more conservative

85    (depending on the number of cues). The learning theory of Rescorla and Wagner is

86    formalized in a set of recurrence equations which specify the association strength $V_i^{t+1}$

87    of cue $C_i$ with outcome $O$ at time $t+1$ as $V_i^{t+1} = V_i^t + \Delta V_i^t$, where the change in

88    association strength $\Delta V_i^t$ is defined as:

89  
$$\Delta V_i^t = \begin{cases} 0 & \text{if ABSENT}(C_i,t) \\ \alpha_i \beta_1 (\lambda - \displaystyle\sum_{PRESENT(C_j,t)} V_j) & \text{if PRESENT}(C_i,t) \ \& \ \text{PRESENT}(O,t) \\ \alpha_i \beta_2 (0 - \displaystyle\sum_{PRESENT(C_j,t)} V_j) & \text{if PRESENT}(C_i,t) \ \& \ \text{ABSENT}(O,t) \end{cases}$$

90

91    In this definition, $\text{PRESENT}(X,t)$ denotes the presence of cue $X$ at time $t$ and

92    $\text{ABSENT}(X,t)$ its absence at time $t$. Whenever the cue occurs without the outcome

93    being present, the association strength is decreased, whereas it is increased when both

94    the cue and outcome are present. The adjustment of the association strength depends

95    on the number of cues present together with the outcome. The standard settings for the

96    parameters are $\lambda = 1$, all $\alpha$'s equal, and $\beta_1 = \beta_2$.

97

98    The Rescorla-Wagner model has been used to explain findings in animal learning and

99    cognitive psychology [18] and more recently, Ramscar and colleagues [19,20,21]

100 have successfully used this model in the context of children's language acquisition.

101 For example, Ramscar and colleagues [21] showed that the Rescorla-Wagner model

102 clearly predicted that exposure to regular plurals (such as *rats*) decreases children's

103 tendency to over-regularize irregular plurals (such as *mouses*) at a certain stage in

104 their development.

105

106 Danks [22] proposed parameter-free equilibrium equations (i.e. where $V_i^{t+1} = V_i^t$) for

107 the recurrence equations presented above: $\Pr(O \mid C_i) - \sum_{j=0}^{n} \Pr(C_j \mid C_i)V_j = 0$, where

108 $\Pr(C_j \mid C_i)$ represents the conditional probability of cue $C_j$ given cue $C_i$, and

109 $\Pr(O \mid C_i)$ the conditional probability of outcome $O$ given cue $C_i$. Consequently, it

110 is possible to directly calculate the association strength between cues and outcomes in

111 the stable (i.e. adult) state where further learning does not substantially change the

112 association weights. Baayen and colleagues [23] have proposed an extension to

113 estimate multiple outcomes in parallel. Their 'naive discriminative learning' (NDL)

114 approach (implementing the Danks equations [22]) lends itself for efficient

115 computation and is readily available via their R package 'ndl'. More details about the

116 underlying computations can also be found in [23].

117

118 After all association strengths of the adult state are determined, the activation (i.e.

119 activation strength) of an outcome given a set of cues can be calculated by summing

120 the corresponding association strengths. Especially these activations are important for

121 prediction. For example, Baayen and colleagues [23] found that the estimated

122 activation of words correlated well with experimental reaction times to those words.

123

124 Here we propose to use naive discriminative learning to determine pronunciation

125 distances. The intuition behind our approach is that a speaker of a certain dialect or

126 language variety is predominantly exposed to speakers who speak similarly, and this

127 input shapes the network of association strengths between cues (in our case,

128 sequences of three sound segments representing the pronunciation, i.e. substrings of

129 the phonetic transcription) and outcomes (in our case, the meaning of the pronounced

130 word) for the speaker. The use of sequences of three segments, so-called trigrams,

131 allows the measure to become sensitive to the adjustments sounds undergo in the

132 context of other sounds, and trigrams have been experimented with in dialectology

133 before [24]. (For comparison, we will also report results when using unigram and

134 bigram cues.) By exposing the speaker to a new pronunciation (in the form of its

135 associated cues) we can measure how well the speaker is likely to understand that

136 pronunciation by inspecting the activation strength of the corresponding outcome. The

137 activation strength of the outcome will depend on the association strengths between

138 the outcome and the cues involved in the pronunciation. If only cues are present

139 which have a high association strength with the outcome, the activation of the

140 outcome will be high, whereas the activation of the outcome will be somewhat lower

141 if one of the cues has a low association strength with the outcome. By calculating the

142 activation strength difference for two different pronunciations of the same word, we

143 obtain a (gradual) measure of pronunciation distance. For example, the word 'with'

144 would be highly activated when a native English listener hears [wɪθ]. However, when

145 a Mandarin speaker would incorrectly pronounce 'with' as [wɪz], this would result in

146 a somewhat lower activation.

147

148 Of course, using an adult state with fixed association weights between cues and

149 outcomes is a clear simplification. Language change is a continuous process and the

150 experience of a listener (i.e. the association weights between cues and outcomes) will

151 obviously be affected by this. However, as the new language experience only makes

152 up a small part of the total language experience of a listener, the effect of the past

153 experience is most important in determining the association weights. As a

154 consequence, and in line with the results of Labov's ([25]: Ch. 4) Cross-Dialectal

155 Comprehension (CDC) studies (which evaluated how well American English speakers

156 understand speakers from their own and other regions), our model will yield lower

157 meaning activations (i.e. more misunderstandings) when sound change is in progress

158 (i.e. the original sound segments will have a higher association strength with the

159 meaning than the new sound segments). In similar fashion, our model predicts higher

160 meaning activations for pronunciations closer to one's own pronunciation variant (i.e.

161 the "local advantage"). We also emphasize that our model is able to capture

162 differences in understandability per word (as each word has its own frequency of

163 occurrence) – which might explain Labov's finding that certain sounds are not always

164 correctly identified, even if they are characteristic of local speakers ([25]: pp. 84-85).

165 Furthermore, the model we propose is general, as it does not focus on a selection of

166 linguistic features (such as vowels), but takes into account all sound (sequences) in

167 determining the understandability of a certain pronunciation.

168

169 Besides being grounded in cognitive theory of competitive reinforcement learning, a

170 clear benefit of this approach is that the pronunciation distances obtained do not need

171 to be symmetrical, as they depend on the association strengths between cues and

172 outcomes, which are different for every speaker. This is illustrated in Section 2.2

173 below.

174

175 To evaluate the effectiveness of this approach, we conducted two experiments. The

176 first experiment focused on investigating foreignness ratings given by native

177 American English (AE) speakers when judging accented English speech, while the

178 second experiment focused on the asymmetric perceptual distances of Norwegian

179 dialect speakers.

180

181 As we noted in the introduction, the Levenshtein distance has been applied to

182 pronunciation transcriptions to assay the degree to which non-local pronunciations

183 sound "different" from local ones (in dialectology, see [1]), but also to predict the

184 comprehensibility of other language varieties (in applied sociolinguistics, see [4]).

185 Since pronunciations may sound non-native or non-local without suffering in

186 comprehensibility, one might suspect that the two notions are not the same, even if

187 they are clearly related. In the present paper we construct a model of an artificial

188 listener to discriminate well enough between words given sound trigrams, which is

189 essentially a comprehension task. But we shall evaluate the same model on how well

190 it predicts human judgments of how similar the speech is to one's own pronunciation

191 (i.e. how native-like foreign accents sound, or how close a pronunciation is to one's

192 own dialect). To the degree to which these experiments succeed, we may conclude

193 that the degree of comprehensibility is largely the same as the degree of nativeness (or

194 localness).

195

196

## Materials and Methods

### 1. Accented English speech

#### 1.1. Material: the Speech Accent Archive

The Speech Accent archive [26] is digitally available at http://accent.gmu.edu and contains a large sample of speech samples in English from people with various language backgrounds. Each speaker read the same paragraph of 69 words (55 of which are unique) in English:

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

All speech samples were transcribed by three phonetically trained transcribers (consensus was reached in the few cases where the transcriptions differed; [26]) according to the International Phonetic Alphabet (IPA). The transcriptions include diacritics, and the associated audio files are available. For this study, we extracted 395 transcribed speech samples and their audio from the Speech Accent Archive. The total number of native U.S.-born English speakers in this dataset was 115. The remaining 280 speech samples belonged to speakers with a different native language or who were born outside of the United States.

#### 1.2. Obtaining NDL-based pronunciation distances

221 For every transcribed pronunciation, we extracted all possible sets of sequences of

222 three sound segments (diacritics were ignored, and a separate segment was added to

223 mark word boundaries) as cues. To model a native AE listener, we randomly selected

224 about half (i.e. 58) of the native AE speakers. We used their pronunciations to

225 generate the pronunciation cues, and paired these with meanings as outcomes (i.e. the

226 pronunciation trigrams were linked to the corresponding meanings). We used only

227 half of the native speakers for the listener model in order to prevent overfitting, i.e.

228 learning the peculiarities of the speakers rather than the features of native American

229 English. The pronunciation of the other half of the speakers is used to represent

230 average American English speech to which the pronunciation of individual speakers is

231 compared. (While we could have used the speech of a single speaker for the listener model

232 and the speech of another individual speaker to represent native American English speech,

233 this would have biased the model to the specific dialectal variants of these speakers.) As the

234 association strength between cues and outcomes depends on the frequency with which

235 they co-occur, we extracted word frequency information from the Google N-Gram

236 Corpus [27]. The total frequency of each meaning outcome was equally divided

237 among all different pronunciations associated with it. For example, if the frequency of

238 the word 'frog' equals 580,000, the frequency of each of the 58 pronunciations was

239 set to 10,000. We then estimated the weights of the model using the 'ndl' package in

240 R (version 0.2.10) which implements the Danks equations [23] introduced above. The

241 resulting network of association strengths between pronunciation cues and meaning

242 outcomes represents a native AE listener. As an example, Table 3 shows part of the

243 input used for estimating the weights and Table 4 shows the association strengths

244 obtained after the weights have been estimated (i.e. the 'adult' association weights of

245 a native AE listener).

246

247 It is clear from Table 4 that the cues found together with a certain outcome generally

248 have a positive value. The more likely it is the cue is found together with the

249 associated outcome (and, crucially, not with other outcomes), the higher the

250 association strength between the two will be.

251

252 Given the table of association strengths representing a simulated native AE listener, it

253 is straightforward to determine the activations of each outcome for a certain

254 pronunciation (converted to cues) by summing the association strengths between the

255 cues in the pronunciation and the outcome. The top half of Table 5 shows that the

256 pronunciations of native AE speakers strongly activate the corresponding outcome

257 (the values are equal or very close to the maximum of 1).

258

259 Of course, we can also use the association strengths (of the simulated native AE

260 listener) to calculate the activations for accented speech. The bottom part of Table 5

261 clearly shows that accented speech results in lower activations (and thus reduced

262 understanding), compared to the pronunciations of native AE speakers (shown in the

263 top part of Table 5). In some cases, a foreign speaker might use a cue which would

264 never be used by a native AE speaker (such as '#xə' in Table 5). As these cues were

265 not encountered during the estimation of the model, no association strengths have

266 been set for those cues and, consequently, their values do not contribute to the

267 activation of the outcome.

268

269 To determine pronunciation distances with respect to native American English, we

270 exposed our model of a native AE speaker to both native American English speech as

271     well as accented English speech and investigated the activation differences of the

272     meaning outcomes. We used the following procedure:

273

274     1. For each of the native American English speakers not considered when

275         constructing the listener model (i.e. the remaining 57 native AE speakers), we

276         calculated the activation of the listener model for each of the 55 different

277         meaning outcomes (i.e. all unique words in our dataset). Whenever an

278         outcome occurred more than once (such as 'we', which occurs twice in the

279         paragraph of text), we averaged the activations associated with the

280         corresponding pronunciations (i.e. the associated cues). For each outcome, we

281         subsequently averaged the activations across all 57 speakers. This is our

282         baseline and can be interpreted as the activations (for 55 individual meanings)

283         of our native AE listener model when being exposed to the speech of an

284         *average* native AE speaker.

285     2. For each individual speaker (mostly non-native, see below), we obtained the

286         activations of our native AE listener model for each of the 55 meanings.

287         Again, whenever an outcome occurred more than once, we averaged the

288         activations associated with the corresponding pronunciations.

289     3. For each individual speaker, we calculated the activation difference compared

290         to the baseline for all 55 meanings separately. We then averaged these

291         activation differences across the 55 meanings. This resulted in a single value

292         for each speaker and represents the NDL-based pronunciation distance with

293         respect to an average native AE speaker.

294

295    As the specific sample of speakers used for estimating the native American English

296    listener model may influence the results, we repeated the random sampling procedure

297    (in which 58 speakers were selected whose pronunciations were used to estimate the

298    listener model) 100 times to generate 100 slightly different native AE listener models.

299    Obviously, this also resulted in a change of the remaining 57 speakers who were used

300    to represent an average AE speaker (see step 1, above). Consequently, we obtained

301    100 (slightly different) NDL-based pronunciation distances for each individual

302    speaker compared to an average AE speaker.

303

304    **1.3.    Validating automatically obtained foreignness ratings**

305    We evaluated the computed pronunciation distances by comparing them to human

306    native-likeness ratings. For this purpose, we developed an online questionnaire for

307    native U.S. English speakers. In the questionnaire, participants were presented with a

308    randomly ordered subset of 50 speech samples from the Speech Accent Archive. We

309    did not include all speech samples, as our goal was to obtain multiple native-likeness-

310    judgments per sample. For each speech sample, participants had to indicate how

311    native-like each speech sample was. This question was answered using a 7-point

312    Likert scale (ranging from 1: very foreign sounding to 7: native AE speaker).

313    Participants were not required to rate all samples, but could rate any number of

314    samples.

315

316    Of course, more advanced methods are possible to measure native-likeness, such as

317    indirect measures which assess the understandability of the accented pronunciations in

318    a certain context (cf. [25: Ch. 4]). However, as our dataset was limited to a small fixed

13

319 paragraph of text, we used a simple rating approach which, nevertheless, resulted in

320 consistent ratings (see results, below).

321

322 Via e-mail and social media we asked colleagues and friends to forward the online

323 questionnaire to people they knew to be native AE speakers. In addition, the online

324 questionnaire was advertised on Language Log by Mark Liberman. Especially that

325 announcement led to an enormous amount of responses. As a consequence, we

326 replaced the initial set of 50 speech samples five times with a new set to increase the

327 number of speech samples for which we could obtain native-likeness ratings. As there

328 was some overlap in the native AE speech samples present in each set (used to

329 calibrate the ratings), the total number of unique samples presented for rating was

330 286, of which 280 were samples from speakers who were not born in the U.S.

331

332 **2. Norwegian dialects**

333 **2.1. Material**

334 The Norwegian dialect material is taken from the study of Gooskens and Heeringa

335 [15], who perceptually evaluated the Levenshtein distance on the basis of IPA

336 transcribed audio recordings of 15 Norwegian dialect speakers reading the fable "The

337 North Wind and the Sun" (containing 58 unique words). The original dataset was

338 created by Jørn Almberg and Kristian Skarbø and is available at

339 http://www.ling.hf.ntnu.no/nos. The transcriptions (including diacritics) were made by

340 the same person, ensuring consistency. Perceptual distances (reported in Table 1 of

341 [15]) were obtained by asking 15 groups of high school pupils (in the corresponding

342 dialect areas) to rate all 15 dialectal audio samples on a scale from 1 (similar to own

14

343   dialect) to 10 (not similar to own dialect). Perceptual dialect distances were then

344   calculated by averaging these ratings per group.

345

346   **2.2.   Methods**

347   Following the same procedure as described in Section 1.2, we converted the

348   pronunciations for each of the 15 speakers in our sample to cues consisting of three

349   sequential sound segments (diacritics were ignored, and a separate segment was added

350   to mark word boundaries). The word frequencies were extracted from a Norwegian

351   word   frequency   list   (on   the   basis   of   subtitles   and   obtained   from

352   http://invokeit.wordpress.com/frequency-word-lists).

353

354   To determine pronunciation distance between dialects $D_i$ and $D_j$ from the perspective

355   of a listener of dialect $D_i$, we used the following procedure:

356

357   1. We estimated the NDL model (i.e. resulting in a specific weight matrix

358      associating cues with outcomes) using the cues on the basis of the

359      pronunciations from the speaker of dialect $D_i$. This model can be seen as

360      representing an experienced listener ($L_i$) of dialect $D_i$.

361   2. We expose $L_i$ to the cues on the basis of the pronunciations from dialect $D_i$ and

362      measure the activation of each of the corresponding 58 meaning outcomes.

363      (Because we only had a single speaker in our sample for each dialect, we

364      could not use separate pronunciations for estimating the listener model and

365      representing the speaker.). Whenever an outcome occurred more than once

366      (some words were repeated), we averaged the activations associated with the

367      corresponding pronunciations (i.e. the associated cues). These activations are

15

368      used as the baseline, and can be interpreted as the activations (for the 58

369      individual meanings) of $L_i$ when being exposed to speech of its own dialect.

370    3. We expose $L_i$ to the cues on the basis of the pronunciations of another dialect

371      $D_j$ and measure the (averaged, when a word occurred more than once)

372      activation of each of the corresponding 58 meaning outcomes.

373    4. For all 58 individual meaning outcomes, we calculated the difference between

374      the activations of $L_i$ for $D_j$ and the baseline $D_i$ and average these 58 differences

375      to get a single value representing the NDL-based pronunciation distance

376      between $D_i$ and $D_j$ (from the perspective of $L_i$).

377

378  The above procedure is repeated for all combinations of $D_i$ and $D_j$ resulting in 210

379  NDL-based pronunciation distances (15 x 15, but the 15 diagonal values are excluded

380  as they are always equal to 0). Table 6 shows these distances for a set of three

381  Norwegian dialects. Note that the NDL-based pronunciation distances between these

382  dialects are clearly asymmetric. The dialect of Bjugn is closer to the dialect of Bergen

383  from the perspective of Bergen (0.545) than the dialect of Bergen is from the

384  perspective of Bjugn (0.559).

385

386  To evaluate these distances, we correlated them with the corresponding perceptual

387  distances (obtained from [15]).

388  **Results**

389  **1. Results for accented English speech**

390  A total of 1143 native American English participants filled in the questionnaire (658

391  men: 57.6%, and 485 women: 42.4%). Participants were born all over the United

392  States, with the exception of the state of Nevada. Most people came from California

16

393   (151: 13.2%), New York (115: 10.1%), Massachusetts (68: 5.9%), Ohio (66: 5.8%),

394   Illinois (64: 5.6%), Texas (55: 4.8%), and Pennsylvania (54: 4.7%). The average age

395   of the participants was 36.2 years (SD: 13.9) and every participant rated on average

396   41 samples (SD: 14.0). Every sample was rated by at least 50 participants and the

397   judgments were consistent (Cronbach's alpha: 0.853).

398

399   To determine how well our NDL-based pronunciation distances on the basis of

400   trigram cues matched the native-likeness ratings, we calculated the Pearson

401   correlation $r$ between the averaged ratings and the NDL-based pronunciation

402   distances for the 286 speakers. Since we had 100 sets of NDL-based pronunciation

403   distances (based on 100 different random samplings of the native American English

404   speakers used to estimate the model), we averaged the corresponding correlation

405   coefficients, yielding an average correlation of $r = -0.72$ ($p < 0.001$). Note that the

406   direction of the correlations is negative as the participants indicated how *native-like*

407   each sample was, while the NDL-based pronunciation distance indicates how foreign

408   a sample is. As a scatter plot clearly revealed a logarithmic relationship (see Figure 1),

409   we log-transformed the NDL-based pronunciation distances, increasing the correlation

410   to $r = -0.80$ ($p < 0.001$). The logarithmic relationship suggests that people are

411   relatively sensitive to small differences in pronunciation in judging native-likeness,

412   but as soon as the differences have reached a certain magnitude (i.e. in our case an

413   NDL-based pronunciation distance of about 0.2) they hardly distinguish them

414   anymore. The sensitivity to small differences is also illustrated by the (slight) increase

415   in performance when trigram cues are used which incorporate diacritics. In that case,

416   the correlation strength increases to $r = -0.75$ ($r = -0.82$ for the log-transformed NDL-

417   based pronunciation distances). These results are comparable with the performance of

418    the Levenshtein distance when applied to this dataset ($r = -0.81$, $p < 0.001$ for the log-

419    transformed Levenshtein distance; unpublished data). In fact, the Levenshtein

420    distances and the NDL-based pronunciation distances also correlate highly, $r = 0.89$

421    ($p < 0.001$).

422

423    We should note that this correlation is close to how well individual raters agree with

424    the average native-likeness ratings (on average: $r = .84$, $p < .0001$). Consequently, the

425    NDL-based method is almost as good as a human rater, despite ignoring

426    suprasegmental pronunciation differences (such as intonation).

427

428    Figure 1 also shows that pronunciations which are perceived as native (i.e. having a

429    rating very close to 7), may correspond to NDL-based pronunciation distances greater

430    than 0. In this case, the NDL-based method classifies certain native-like features as

431    being non-native. This may be caused by our relatively small sample of only 58

432    speakers whose pronunciations were used to model the native AE listener. Real

433    listeners have much more experience with their native language, and therefore can

434    more reliably distinguish native-like from foreign cues.

435

436    The aforementioned results are all based on using trigram cues. When using unigram

437    cues instead, the correlation between the perceptual native-likeness ratings and the

438    NDL-based pronunciation distances dropped to $r = -0.54$ (log-transformed: $r = -0.57$).

439    When using bigram cues, the performance was almost on par with using trigram cues

440    ($r = -0.69$, log-transformed: $r = -0.79$). Using unigram and/or bigram cues together

441    with trigram cues did not affect performance, as these simpler cues are not

442    discriminative in the presence of trigram cues.

443

**2. Results for Norwegian dialects**

The correlation between the NDL-based pronunciation distances and the perceptual

distances was $r = 0.68$ ($p < 0.001$), which is comparable to the correlation Gooskens

and Heeringa [15] reported on the basis of the Levenshtein distance (i.e. $r = 0.67$).

Similar to the first study, log-transforming the NDL-based pronunciation distances

increased the correlation strength to $r = 0.72$ ($p < 0.001$). In line with the results for

the accent data, the Levenshtein distances and the NDL-based pronunciation distances

correlate highly, $r = 0.89$ ($p < 0.001$).

452

The aforementioned results are all based on using trigram cues. Using unigram cues

instead of trigram cues severely reduced performance ($r = 0.10$, log-transformed: $r =$

$0.31$), whereas using bigram cues was almost as good as using trigram cues ($r = 0.67$,

log-transformed: $r = 0.71$). Similar as before, adding unigram and/or bigram cues to

the trigram cues did not really improve performance. In contrast to the accent data,

incorporating diacritics in the cues also did not help; the correlation then dropped to $r$

$= 0.65$ (log-transformed: $r = 0.66$). This is likely caused by the relatively small

dataset.

461

**Discussion**

In the present paper we have shown that pronunciation distances derived from naive

discriminative learning match perceptual accent and dialect distances quite well.

While the results were on par with those on the basis of the Levenshtein distance, the

advantage of the present approach is that it is grounded in cognitive theory of

comprehension based on fundamental principles of human discrimination learning.

468    Furthermore, the Levenshtein distance is theoretically less suitable for modeling the

469    degrees of difference in the perception of non-local and non-native speech because it

470    is a true distance, i.e. always symmetric, while perceptions of similarity may also be

471    asymmetric [15]. The NDL-based approach naturally generates asymmetrical

472    distances.

473

474    We noted above that the task of recognizing words based on phonetic cues is

475    essentially a comprehensibility task. A second contribution of the present paper is

476    therefore to demonstrate that models constructed to comprehend local speech

477    automatically assign scores of non-nativeness (or of non-localness among dialects) in

478    a way that models native speakers judgments.

479

480    One may wonder why the NDL-based method only slightly improved upon the results

481    of the Levenshtein distance for the Norwegian dataset, especially since that dataset is

482    characterized by asymmetric perceptual distances. We note here that the 15 NDL

483    models (one for each listener) are only based on the pronunciation of a single speaker.

484    Consequently, it does not take into account the variation within each dialect (taken

485    into account by listeners living in the dialect area), which would have allowed for

486    more precise estimates of the association weights. A general limitation is that

487    Gooskens and Heeringa [15] already indicated that intonation is one of the most

488    important characteristics in Norwegian dialects, and no such cues have been used here

489    (as these were not available to us), thereby limiting the ability to detect relevant

490    asymmetries. Nerbonne and Heeringa ([28]: 563-564), on the other hand, speculate

491    that there is a limit to the accuracy of validating pronunciation difference measures on

492    the basis of aggregate judgments of varietal distance. If one supposes that poorer

493    measures are noisier – but not more biased – than better ones, then the noise will

494    simply be eliminated in examining large aggregates. If this is right, we cannot expect

495    to change mean differences by adopting more accurate measurements. They suggest

496    that improved validation will therefore have to focus on smaller units such as

497    individual words.

498

499    While we have not explored this in the present paper, another important advantage of

500    the NDL approach is that cues are not only restricted to phonetic segments. Cues with

501    respect to pronunciation speed or other acoustic characteristics (such as intonation)

502    can be readily integrated in an NDL model (e.g., linking cues representing different

503    intonation patterns to the individual meanings). A problem of the NDL method,

504    however, is that it only accepts discrete cues. A continuous measurement therefore

505    needs to be discretized to separate cues, and this introduces a subjective element in an

506    otherwise parameter-free procedure.

507

508    As our datasets only consisted of a few dozen words, our model was highly simplified

509    compared to the cognitive model of a human listener who will have access to

510    thousands of words. It is nevertheless promising that pronunciation distances on the

511    basis of our simplified models match perceptual distances at least as well as current

512    gold standards.

513

517 used in this study. We also thank Michael Ramscar for interesting discussions about

518 the ideas outlined in this paper.

519

520 **References**

521 1. Heeringa W (2004) Measuring Dialect Pronunciation Differences using

522     Levenshtein Distance. PhD thesis, Rijksuniversiteit Groningen.

523 2. Valls E, Wieling M, Nerbonne J (2013) Linguistic advergence and divergence

524     in Northwestern Catalan: A dialectometric investigation of dialect leveling

525     and border effects. LLC: Journal of Digital Scholarship in the Humanities

526     28(1): 119-146.

527 3. Bakker D, Müller A, Velupillai V, Wichmann S, Brown CH, et al. (2009)

528     Adding typology to lexicostatistics: A combined approach to language

529     classification. Linguistic Typology 13(1): 169-181.

530 4. Beijering K, Gooskens C, Heeringa W (2008) Predicting intelligibility and

531     perceived linguistic distances by means of the Levenshtein algorithm.

532     Linguistics in the Netherlands 15: 13-24.

533 5. Sanders NC, Chin SB (2009) Phonological distance measures. Journal of

534     Quantitative Linguistics 43: 96-114.

535 6. Wieling M, Nerbonne J, Baayen RH (2011) Quantitative social dialectology:

536     Explaining linguistic variation geographically and socially. PLOS ONE 6(9):

537     e23613. doi:10.1371/journal.pone.0023613.

538 7. Kessler B (1995) Computational dialectology in Irish Gaelic. In: Proceedings

539     of the Seventh Conference on European Chapter of the Association for

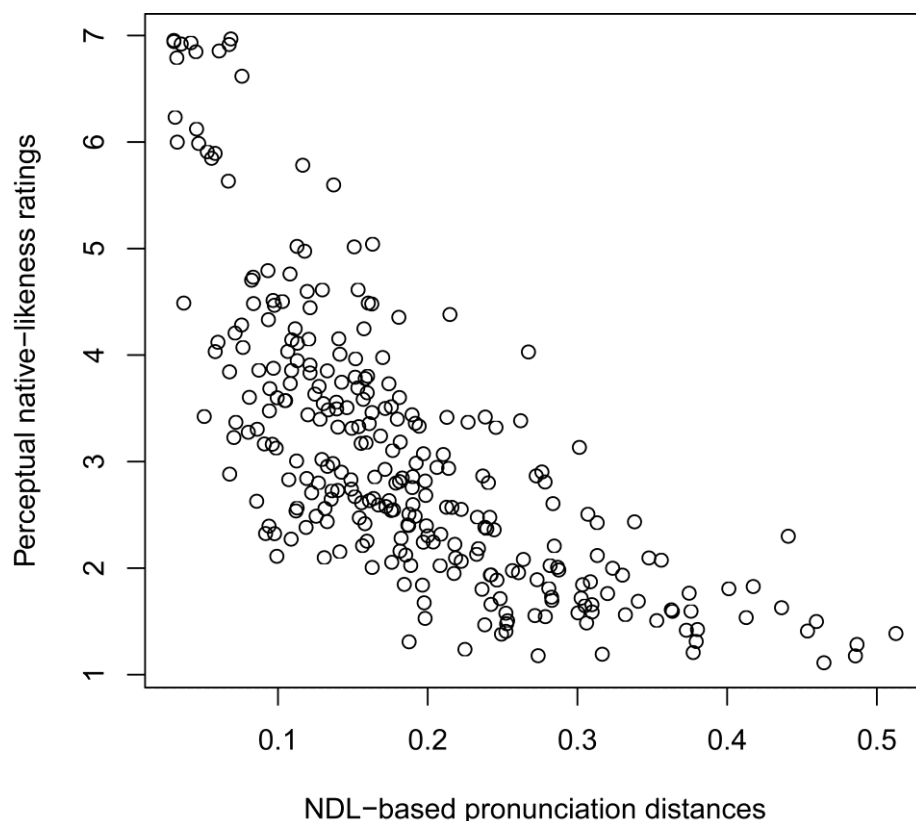540     Computational Linguistics. pp. 60-66.

541 8. Nerbonne J, Heeringa W (1997) Measuring dialect distance phonetically In:

542   Coleman J, editor. Workshop on Computational Phonology. Madrid: Special

543   Interest Group of the Association for Computational Linguistics. pp. 11-18.

544 9. Wichmann S, Holman EW, Bakker D, Brown CH (2010) Evaluating linguistic

545   distance measures. Physica A 389: 3632-3639.

546 10. Wieling M, Heeringa W, Nerbonne J (2007) An aggregate analysis of

547   pronunciation in the Goeman-Taeldeman-van Reenen-Project data. Taal en

548   Tongval 59(1): 84-116.

549 11. Levenshtein V (1965) Binary codes capable of correcting deletions, insertions

550   and reversals. Doklady Akademii Nauk SSSR 163: 845-848. In Russian.

551 12. Heeringa W, Braun A (2003) The use of the Almeida-Braun system in the

552   measurement of Dutch dialect distances. Computers and the Humanities

553   37(3): 257-271.

554 13. Wieling M, Prokić J, Nerbonne J (2009) Evaluating the pairwise alignment of

555   pronunciations. In: Borin L, Lendvai P, editors. Proceedings of the EACL

556   2009 Workshop on Language Technology and Resources for Cultural

557   Heritage, Social Sciences, Humanities, and Education. pp. 26-34.

558 14. Wieling M, Margaretha E, Nerbonne J (2012) Inducing a measure of phonetic

559   similarity from dialect variation. Journal of Phonetics 40(2): 307-314.

560 15. Gooskens C, Heeringa W (2004) Perceptive evaluation of Levenshtein dialect

561   distance measurements using Norwegian dialect data. Language Variation and

562   Change 16(3): 189-207.

563 16. Wieling M, Nerbonne J (2007). Dialect pronunciation comparison and spoken

564   word recognition. In: Osenova P et al., editors. Proceedings of the RANLP

565   Workshop on Computational Phonology. pp. 71-78.

566    17. Rescorla RA, Wagner AR (1972) A theory of Pavlovian conditioning:

567        Variations in the effectiveness of reinforcement and nonreinforcement. In:

568        Black AH, Prokasy WF, editors. Classical conditioning II: Current research

569        and theory. New York: Appleton-Century-Crofts. pp. 64-99.

570    18. Siegel SG, Allan LG (1996) The widespread influence of the Rescorla-Wagner

571        model. Psychonomic Bulletin and Review 3(3): 314-321.

572    19. Ramscar M, Yarlett D, Dye M, Denny K, Thorpe K (2010). The effects of

573        feature-label-order and their implications for symbolic learning. Cognitive

574        Science 34(6): 909-957.

575    20. Ramscar M, Dye M, Popick HM, O'Donnell-McCarthy F (2011) The Enigma

576        of number: Why children find the meanings of even small number words hard

577        to learn and how we can help them do better. PLOS ONE 6: e22501.

578        doi:10.1371/journal.pone.0022501.

579    21. Ramscar M, Dye M, McCauley S (2013) Error and expectation in language

580        learning: The curious absence of 'mouses' in adult speech. Language: In

581        press.

582    22. Danks D (2003). Equilibria of the Rescorla–Wagner model. Journal of

583        Mathematical Psychology 47: 109-121.

584    23. Baayen RH, Milin P, Filipovic Durdevic D, Hendrix P, Marelli M (2011) An

585        amorphous model for morphological processing in visual comprehension

586        based on naive discriminative learning. Psychological Review 118: 438-482.

587    24. Heeringa W, Kleiweg P, Gooskens C, Nerbonne J (2006). Evaluation of string

588        distance algorithms for dialectology. In: Nerbonne J, Hinrichs E, editors.

589        Linguistic Distances. Sydney: COLING/ACL. pp. 51-62.

590    25. Labov, W. (2010) Principles of Linguistic Change, Cognitive and Cultural

591         Factors, Vol. 3. Malden: Wiley-Blackwell.

592    26. Weinberger, SH, Kunath SA (2011) The Speech Accent Archive: Towards a

593         typology of English accents. Language and Computers 73: 265-281.

594    27. Brants T, Franz A (2009) Web 1T 5-gram, 10 European languages. Version 1.

595         Philadelphia: Linguistic Data Consortium.

596    28. Nerbonne J, Heeringa W (2010) Measuring dialect differences. In: Auer P,

597         Schmidt JE, editors. Language and Space: Theories and Methods. Berlin:

598         Mouton De Gruyter. pp. 550-566.

599

600   **Figure**



601

602   **Figure 1.** Logarithmic relationship between NDL-based pronunciation distances and

603   perceptual distances.

604 **Tables**

605 **Table 1.** Basic Levenshtein distance alignment.

| w | ε | n |   | z | d | e | ɪ |
|---|---|---|---|---|---|---|---|
| w | ε | n | ə | s | d | e |   |
|   |   |   | 1 | 1 |   | 1 |   |

606

607 **Table 2.** Levenshtein distance alignment with sensitive sound distances.

| w | ε | n |     | z | d | e | ɪ |
|---|---|---|-----|---|---|---|---|
| w | ε | n | ə   | s | d | e |   |
|   |   |   | 0.031 | 0.020 |   | 0.030 |   |

608 **Table 3.** Part of the table used for estimating the association strengths. The '#' marks

609 the word boundary.

| Speaker | Outcome | Pronunciation | Cues | Frequency |
|---------|---------|---------------|------|-----------|
| english23 | with | [wɪθ] | #wɪ, wɪθ, ɪθ# | 28,169,384 |
| english167 | with | [wɪð] | #wɪ, wɪð, ɪð# | 28,169,384 |
| english23 | her | [həɹ] | #hə, həɹ, əɹ# | 852,131 |
| english167 | her | [ɚ] | #ɚ# | 852,131 |

610

611

612 **Table 4.** The association strengths for the cues and outcomes in Table 1 for our

613 simulated native AE listener after these have been estimated on the basis of the input

614 of 58 randomly selected native AE speakers.

| Cue | Association strength for 'with' | Association strength for 'her' |
|-----|--------------------------------|-------------------------------|
| #wɪ | 0.2519 | 0.0000 |
| wɪθ | 0.3738 | 0.0000 |
| ɪθ# | 0.3738 | 0.0000 |
| wɪð | 0.3741 | 0.0000 |
| ɪð# | 0.3741 | 0.0000 |
| #hə | 0.0000 | 0.4973 |
| həɹ | 0.0000 | 0.2433 |
| əɹ# | 0.0000 | 0.2594 |
| #ɚ# | 0.0000 | 1.0000 |

26

615 **Table 5.** The activations of different outcomes on the basis of the association

616 strengths between the cues and outcomes for our simulated native AE listener (shown

617 in Table 2).

| Speaker | Outcome | Pronunciation | Cues | Activation of outcome |
|---|---|---|---|---|
| english23 | with | [wɪθ] | #wɪ, wɪθ, ɪθ# | 0.9995 |
| english167 | with | [wɪð] | #wɪ, wɪð, ɪð# | 1.0000 |
| english23 | her | [həɹ] | #hə, həɹ, əɹ# | 1.0000 |
| english167 | her | [ɚ] | #ɚ# | 1.0000 |
| mandarin10 | with | [wɪz] | #wɪ, wɪz, ɪz# | 0.2519 |
| serbian10 | her | [xəɹ] | #xə, xəɹ, əɹ# | 0.2594 |

618

619

620 **Table 6.** Part of the NDL-based Norwegian dialect pronunciation distances.

| | Bergen | Bjugn | Bodø |
|---|---|---|---|
| Bergen | X | 0.545 | 0.584 |
| Bjugn | 0.559 | X | 0.319 |
| Bodø | 0.574 | 0.314 | X |

621