

Computational modelling of segmental and prosodic levels of analysis for capturing variation across Arabic dialects

Georgina Brown^{1*}, Sam Hellmuth²

¹Department of Linguistics and English Language, Lancaster University, United Kingdom

²Department of Language and Linguistic Science, University of York, United Kingdom

¹g.brown5@lancaster.ac.uk

²sam.hellmuth@york.ac.uk

* Correspondence:

Georgina Brown

g.brown5@lancaster.ac.uk

Keywords: Arabic dialects, accent, intonation, automatic accent recognition, Support Vector Machines.

Abstract

Dialect variation spans different linguistic levels of analysis. Two examples include the typical phonetic realisations produced and the typical range of intonational choices made by individuals belonging to a given dialect group. Taking the modelling principles of a specific automatic accent recognition system, the work here characterises and observes the variation that exists within these two levels of analysis among eight Arabic dialects. Using a method that has previously shown promising performance on English accent varieties, we first model the segmental level of analysis from recordings of Arabic speakers to capture the variation in the phonetic realisations of the vowels and consonants. In doing so, we show how powerful this model can be in distinguishing between Arabic dialects. This paper then shows how this modelling approach can be adapted to instead characterise prosodic variation among these same dialects from the same speech recordings. This allows us to inspect the relative power of the segmental and prosodic levels of analysis in separating the Arabic dialects. This work opens up the possibility of using these modelling frameworks to study the extent and nature of phonetic and prosodic variation across speech corpora.

1 Introduction

Many recent approaches to automatic accent recognition have depended heavily on machine learning techniques, falling in line with trends across the breadth of speech technology (Najafian, et. al., 2018; Shon et. al., 2018). Usually though, these approaches do not yield accent recognition rates that are

33 comparable with the low error rates we see in related areas like automatic speaker recognition
34 (Snyder et. al., 2017). Additionally, these approaches demand enormous, and therefore often
35 unattainable, datasets to develop working systems. One way of overcoming the need for very large
36 datasets in automatic accent recognition is to be selective in its development and inform the system
37 of the specific features it should use to model speakers' accents. The York ACCDIST-based
38 automatic accent recognition system (Brown, 2015; Brown and Wormald, 2017) is an example of a
39 system that takes this more targeted approach. Based on the ACCDIST metric (Huckvale, 2004,
40 2007), Y-ACCDIST models encapsulate only a subset of features that are expected to represent a
41 speaker's production of the phoneme inventory. In doing so, Y-ACCDIST has a lowered reliance on
42 machine learning techniques that would otherwise involve the extraction of many features from right
43 across the speech sample, which would then be used to derive a subset that is estimated to comprise
44 the most useful features for the task at hand. As implemented to date, Y-ACCDIST targets the
45 phonetic realisations of the individual vowels and consonant segments in the language and compares
46 one speaker's set of realisations with the corresponding sets of other speakers. This comparison
47 gauges which group of speakers (grouped by accent) the speaker is most similar to. The first
48 experiments in this paper demonstrate the performance of this "segmental" version of the Y-
49 ACCDIST system on speech recordings taken from speakers of eight Arabic dialects. These
50 experiments simultaneously show its use as an automatic dialect classification system and as a way
51 of observing variation among accents and dialects.

52
53 While attempting to isolate the segmental level has its advantages (as it is the level of analysis that is
54 expected to be most valuable to dialect classification), we are aware that there are other potentially
55 useful features within the speech signal that this approach overlooks. There is growing evidence of
56 accent- or dialect-specific intonation patterns in a number of languages. For example, computational
57 analysis of data from the Intonational Variation in English (IViE) project in seven different British
58 English varieties, showed differences in the shape and distribution of f0 contours across dialects
59 (Grabe, Kochanski, & Coleman, 2007). A key contribution of this paper is to ascertain whether the
60 modelling procedure in the standard segmental form of the Y-ACCDIST system can also be applied
61 to the prosodic level of analysis. This will then enable us to compare the contribution of segmental
62 and prosodic cues to a specific dialect classification task, while removing other potentially distracting
63 information embedded within the speech signal.

64

65 The dataset that has been used in the experiments presented in this work is the *Intonational Variation*
66 *in Arabic* (IVAr) corpus (Hellmuth and Almbark, 2019). There are other, larger, speech corpora
67 available that would allow for research to be conducted on different Arabic dialects. The Multi-Genre
68 Broadcast (MGB-5) challenge dataset (Ali et. al., 2019) is one such example which consists of
69 hundreds of hours of data from 17 countries, a subset of which has been labelled for dialect group by
70 human annotators. Despite MGB-5's appealing size, there are a number of reasons why the IVAr
71 corpus is better suited to the present study. Firstly, Ali et. al., (2019) concede that there will be
72 labelling errors as a result of their dataset construction method. Much of the metadata is often led by
73 the country of the YouTube channels, for example, from which the speech data have been identified.
74 Dialect labels may therefore be estimations at times, bringing in noise to any dialect research. The
75 IVAr corpus metadata, on the other hand, are extremely controlled and reliable, allowing us to draw
76 more robust findings. Secondly, the speech samples in the MGB-5 dataset are generally too short.
77 MGB-5 speech samples are categorised according to their durations: *short* (<5s), *medium* (5-20s) and
78 *long* (>20s). This distribution of sample durations is insufficient for the methods implemented in the
79 present study, where, ideally, we would be using at least one minute of speech per speaker. Thirdly,
80 the IVAr corpus was collected in such a way that elicited speech for the purpose of prosodic research
81 (i.e. a carefully selected and informed set of sentences and speech tasks that prompt intonation
82 patterns of interest). The present study would not be possible without such control in the data
83 construction. Lastly, the IVAr corpus has already had a substantial amount of prosodic analysis
84 conducted on it (Hellmuth, 2018; Hellmuth, to appear). This enables us to interpret the performance
85 of the modelling procedure in the context of prosodic analysis that has been conducted using more
86 traditional analytical methods. These kinds of analytical procedures have typically involved manual
87 qualitative labelling of samples of data using a system of prosodic annotation such as the Tones and
88 Break Indices (ToBI) system (Beckman & Elam, 1997; Beckman, Hirschberg, & Shattuck-Hufnagel,
89 2005) or more recent systems proposed for use across languages (Hualde & Prieto, 2016), as well as
90 quantitative approaches such as visualisation and statistical analysis of f0 contour shapes (Hellmuth,
91 2018).

92
93 Recently, a more innovative way of capturing prosodic variation has been proposed. Elvira-García et.
94 al. (2018) introduced the *ProDis* dialectometric tool for measuring prosodic distances between
95 linguistic varieties based on acoustic measurements. *ProDis* involves logging the correlations
96 between the pitch contours of specified sentences produced by speakers, and then comparing these
97 correlations among a speaker set representing a range of languages. This provides a dialectometric

98 method that aims to reveal prosodic similarities and differences between linguistic varieties. The
99 authors motivate their work by pointing out that efforts have been made to measure dialect and
100 language differences by making phonological or lexical comparisons, but that we lack an equivalent
101 that makes use of prosodic information. Their demonstration of using *ProDis* shows its application to
102 a subset of AMPER (Atlas Multimédia Prosodique de l’Espace Roman) (Contini and Romano, 2002),
103 which is an international effort to capture data that represents a full range of Romance linguistic
104 varieties. Within their work, they applied the *ProDis* tool to 7 dialects from across 5 Romance
105 languages. Using *ProDis*, Elvira-García et. al. were able to perform cluster analyses and associated
106 data visualisations on these data, followed by some qualitative evaluation. For example, they
107 produced a dendrogram of their *ProDis* data representations. One of their clusters was neatly made
108 up of varieties that are largely spoken in Sardinia, and they were able to provide an accompanying
109 example of the characteristic intonation contour shape of yes/no-questions produced by speakers of
110 those varieties.

111

112 Similarly, one version of the Y-ACCDIST system has been presented as another way to quantify
113 differences among accent varieties, by measuring and modelling phonetic realisational differences of
114 segments, demonstrated in Brown and Wormald (2017). Like Elvira-García et. al.’s study above,
115 Brown and Wormald were able to draw observations from a dendrogram of Y-ACCDIST
116 representations of different speakers in a speech dataset. In their work, they looked at the accent
117 differences between Punjabi-English and Anglo-English speakers in Bradford and Leicester in
118 England. One of the pertinent patterns to emerge was that there were some clusters that grouped the
119 speakers according to the community centre they attended, which perhaps went beyond the types of
120 grouping that the authors originally expected. As well as the cluster analyses, Brown and Wormald
121 were also able to perform some feature selection analyses (using the Y-ACCDIST models as a
122 framework of features) which indicated the vowels and consonants that were estimated to separate
123 the accent varieties in the dataset. This analysis pointed towards the GOAT vowel and /ɪ/ as features
124 that discriminated these accent varieties, which corresponded with some of the more traditional
125 acoustic analysis conducted in Wormald (2016).

126

127 Another ACCDIST-based system was demonstrated to observe accent variation among a larger
128 number of accents from across the British Isles in Ferragne and Pellegrino (2010), which also took
129 advantage of the variation in phonetic realisations. In their study, Ferragne and Pellegrino took
130 controlled wordlist data and created an ACCDIST-based model of the vowel systems of 261 speakers

131 who represented 13 accents from the Accents of the British Isles (ABI) corpus (D’Arcy et. al., 2004).
132 They also found that these models yielded linguistically explicable patterns in visualisations of the
133 data. For example, they found a very neat split in a cluster analysis between the Scottish, Irish and
134 English accent varieties in the corpus.

135
136 By implementing a Y-ACCDIST-based framework to model speakers’ intonational inventories, in
137 this work we apply a similar modelling procedure to that presented in Elvira-García et. al. (2018).
138 However, by implementing a framework that has also been used to capture segmental phonetic
139 realisational differences between different accent varieties, we can draw comparisons between how
140 prosodic information and segmental information distinguish linguistic varieties under investigation.
141 Additionally, by modelling numerous speakers per dialect group, we have an opportunity to train a
142 dialect classification system on the prosodic information alone to be able to observe how much this
143 single level of analysis could contribute to an accent or dialect classification task. Although the
144 dataset used to demonstrate *ProDis* in Elvira-García et. al. was very large, the number of speakers per
145 variety was very small (less than 5), and so did not provide the opportunity for an experiment of the
146 kind presented here.

147
148 Until the current work, Y-ACCDIST had only been tested on datasets of speech in English. We first
149 demonstrate its performance in distinguishing between dialects of Arabic in its original segmental
150 configuration (i.e. targeting the phonetic realisations of different segments), and we show results on
151 both controlled read speech and spontaneous speech. We then move on to explore the Y-ACCDIST-
152 based framework for modelling the prosodic variation among accents, allowing us to compare the
153 different value that segmental and prosodic levels of speech analysis bring to the dialect recognition
154 task. We also delve into the inner workings of the machine learning within the system to determine
155 whether we can identify particularly useful features within the segmental and prosodic models that
156 can discriminate the Arabic dialects. All the analysis tasks conducted for this study are interpreted in
157 the context of the existing prosodic analysis conducted on these same data.

158 In summary, this paper addresses the following broad objectives:

- 159 • to observe the Y-ACCDIST system’s recognition performance on Arabic dialect varieties and
160 interpret the results in the context of existing linguistic analyses of the data.
- 161 • to compare the performance of the Y-ACCDIST system on read speech and spontaneous
162 speech on the same dialect classification task.

- 163 • to transfer Y-ACCDIST’s modelling technique from the segmental level of analysis to the
164 prosodic level and compare dialect classification performance between these two levels of
165 analysis.

166

167 **2 Arabic Dialects**

168 **2.1 Overview of Arabic dialects**

169 Arabic is one of the world’s largest languages, spoken as a native language by at least 300 million
170 speakers (Owens, 2013), yet consisting of a diverse array of spoken vernaculars which vary from
171 each other at all levels of linguistic analysis – from phonetics and phonology to morphosyntax and
172 lexis (Retsö, 2013). There is a clear divide between western ‘maghreb’ dialects spoken in North
173 Africa and eastern ‘mashreq’ dialects spoken elsewhere (Behnstedt & Woidich, 2013), such that
174 human listeners can distinguish these two broad groups based solely on prosodic information (Barkat,
175 Ohala, & Pellegrino, 1999). A commonly used geographical approach to grouping Arabic dialects,
176 based on shared linguistic features within groups, results in the following five-way grouping, from
177 west to east (Versteegh, 2014): i) dialects of North Africa (including Morocco, Algeria, Libya and
178 Tunisia); Egyptian dialects (including Egypt and Sudan); Levantine dialects (including Jordan,
179 Lebanon, Syria and Palestine); Mesopotamian dialects (including Iraq); and dialects of the
180 Gulf/Arabian Peninsula (including Saudi Arabia, Kuwait, Bahrain, Qatar, Oman and Yemen). This
181 five-way split has been widely implemented in computational approaches to the Arabic dialect
182 classification task (e.g. Biadisy et. al., 2009). Nevertheless, the degree of dialectal variation within
183 each of these five groups is considerable, with additional important dialectal discontinuities due to
184 historical contact and migration, social categories and lifestyle (with a common broad divide between
185 dialects which are sedentary/urban versus nomadic/rural in origin) as well as religious or sectarian
186 affiliation (Behnstedt & Woidich, 2013). As a result of these cross-cutting factors contributing to
187 dialectal variation, Arabic is frequently described as a ‘mosaic’ of dialects. ‘Successful’ dialect
188 classification for Arabic would ideally be able to tackle different degrees of granularity, both between
189 and within the broad regional groupings that are usually taken as targets.

190 **2.2 Automatic classification of Arabic dialects**

191 As indicated in the Introduction, many approaches to automatic dialect identification have depended
192 heavily on machine learning approaches, usually inspired by the techniques tested for Language
193 Identification (LID). These approaches have demanded vast quantities of data for training. Biadisy et.

194 al. (2009) applied a Phone Recognition followed by Language Modelling (PRLM) approach to
195 Arabic dialect classification, which was first introduced by Zissman (1996) for the purpose of LID.
196 As the name suggests, PRLM starts by feeding a speech sample through a phone recognition system
197 to establish an estimated sequence of phones in the sample. This estimated sequence is then
198 compared against the phone sequences and distributions computed for the different linguistic
199 varieties in the reference system (i.e. the training data). PRLM therefore depends on the different
200 varieties we are distinguishing between to have phone sequences and distributions that are separable.
201 For LID, this seems to achieve reasonable performance, but as the varieties we are distinguishing
202 between become more and more similar (i.e. dialects and then accents), this approach is expected to
203 become less effective. In their work, Biadisy et. al. reported that the PRLM approach achieved 81.6%
204 accuracy for an identification task involving speakers of five Arabic dialect groups (using the
205 commonly used grouping described in Section 2.1 above).

206
207 The PRLM approach is the more traditional one for this sort of task. Researchers have since applied
208 classifiers based on neural networks to the problem of Arabic dialect recognition (Najafian, et. al.,
209 2018; Shon, et. al., 2018). These works follow in the footsteps of developments in speaker
210 recognition research, where a new method of modelling the variation among different speakers in the
211 form of “embeddings” was proposed, in an effort to improve on the performance of i-vector-based
212 systems (Snyder et. al., 2017). Such methods demand vast amounts of training data (ideally,
213 hundreds of speech samples per dialect group). Both of the studies mentioned above which apply the
214 neural network based approach to Arabic dialect identification used the Multi-Genre Broadcast 3
215 (MGB-3) dataset, which offers 63.6 hours of training data across the five main Arabic dialect groups.
216 Shon et. al. (2018) achieved 73% accuracy using a neural network based system, outperforming the i-
217 vector systems they compared on the same task.

218
219 In this paper, our experiments will be conducted on a corpus of speech recordings taken from 96
220 speakers spanning 8 Arabic dialect categories. We therefore present ourselves with a dialect
221 classification problem which has a fraction of the data to train a system on. In addition, we assume
222 that this is a more difficult problem in that we have increased the level of similarity between dialects
223 by having 8 dialect categories, rather than 5 broader ones. The Y-ACCDIST-based method we are
224 employing is much better suited to a dataset of this size and nature (as demonstrated in Brown
225 (2016)).

226

227 2.3 The IVAr Corpus

228 The core Intonational Variation in Arabic (IVAr) corpus contains recordings from 12 speakers each
229 in 8 spoken dialects of Arabic (96 speakers in total), collected on location in North Africa and the
230 Middle East (Hellmuth & Almbark, 2019)¹. IVAr provides at least one dataset from each regional
231 dialect group, with more than one dataset for the more linguistically diverse regional groups
232 (Levantine/Gulf/North Africa). The corpus thus provides for an eight-way dialect classification task,
233 across the geographically defined dialects listed in Table 1.

235 Table 1. Dialects represented in the Intonational Variation in Arabic Corpus
236

Code	Dialect	Recording location	Regional group
moca	Moroccan Arabic (Casablanca)	Casablanca, Morocco	North Africa
tuns	Tunisian Arabic (Tunis)	Tunis, Tunisia	
egca	Egyptian Arabic (Cairo)	Cairo, Egypt	Egyptian
joka	Jordanian Arabic (Karak)	Karak, Jordan	Levantine
syda	Syrian Arabic (Damascus)	Amman, Jordan	
irba	Iraqi Arabic (Muslim Baghdadi)	Amman, Jordan	Mesopotamian
kwur	Kuwaiti Arabic (Urban)	Kuwait City, Kuwait	Gulf/Arabian Peninsula
ombu	Gulf Arabic (Buraimi)	Buraimi, Oman	

237
238

239 Use of IVAr allows us to demonstrate the dialect identification task at a more granular level than is
240 typical in the field, since most other work on dialect identification for Arabic attempts at most a five-
241 way regional classification (due, in turn, to the fact that most large Arabic corpora provide datasets
242 defined at a regional level only).

243

244 The corpus contains speech elicited in a range of speech styles, from scripted read speech to
245 unscripted semi-spontaneous speech. The scripted materials were presented to participants printed in
246 Arabic script, using the informal spelling conventions of each local dialect (rather than following the
247 norms of standard Arabic); in this paper we use data elicited by means of a scripted dialogue (sd)
248 performed as a role play between pairs of speakers and a monologue narrative folk tale (sto). The
249 spontaneous speech data used in this paper comprise a monologue folk tale retold from memory (ret),
250 an information-gap map task performed in dialogue between pairs of speakers (map), and free

¹ The full corpus comprises 10 datasets across 8 dialects; that is, for one of the 8 dialects, Moroccan Arabic, there are two additional datasets: one with bilingual speakers of Moroccan Arabic and Tashlhiyt Berber aged 18-35 (mobi), and one with Moroccan Arabic speakers aged 40-60 (moco). These two additional datasets are not investigated in the present study.

251 conversation between pairs of speakers (fco). Further information about the instruments used to elicit
252 the data is available at ivar.york.ac.uk/.

253

254 The participants in each location were recruited through a local fieldwork representative, typically
255 through an educational institute such as a university or private language school. Participants ages
256 ranged from 18-35 years. All recordings took place in the city in which participants were resident,
257 and recruitment was carefully monitored to ensure participants were speakers of the target dialect and
258 had been raised in the target city. The only exception was speakers of Syrian and Iraqi Arabic, who
259 were recruited in Amman, Jordan due to the prevailing security situation in Syria and Iraq at the time
260 of recording. Detailed participant metadata is provided with the published corpus. All participants
261 received an information sheet in Arabic and provided informed written consent prior to recording.

262

263 Participants were recorded in pairs using head-mounted Shure SM10A dynamic microphones directly
264 to .wav format on a Marantz PMD660/620 digital recorder at 44.1kHz 16 bit, with each speaker
265 recorded to a separate stereo channel which can be split to analyse speakers separately. Recording
266 sessions were run by a local fieldworker who was a native speaker of the same dialect. All of the
267 tasks, scripted and unscripted, were performed in a single recording session, with the same
268 interlocutor and under the same recording conditions.

269

270 The spontaneous speech data were orthographically transcribed by native speaker research assistants
271 using a romanised phonetically transparent transliteration system adapted for each dialect; these
272 transcriptions are available as part of the published corpus. For the read speech, the script used during
273 data collection was transcribed into the same transliteration system, and are also made available with
274 the corpus. For the present project we created a merged dictionary of all of the dialect-specific forms
275 used in transcripts for read and spontaneous speech across all dialects; a native speaker of Arabic
276 proficient in Modern Standard Arabic (MSA) created a transcription of each dialect-specific form
277 using a common MSA phone set to create the merged dictionary. This was based on the accepted
278 cognate sound in MSA of dialect-specific variants. For example, the name of the main character in
279 the folk tale retold from memory is variously produced in the dialects as [zuħa], [dzuħa] or [guħa]
280 <جحا> and appears in the merged dictionary as pronounced in MSA i.e. as [zuħa]. We intend on
281 publication of the present paper to make this merged dictionary available as an appendix to the main
282 published IVAr corpus.

283

284 As already discussed in the Introduction, larger speech datasets of Arabic dialects exist, such as
285 MGB-3 and MGB-5. However, such datasets have not been collected in a way that allows us to
286 explore the specific research questions in this paper that involve analysing prosodic variation as well
287 as segmental variation.

288
289

290 **3 The Y-ACCDIST System**

291 Y-ACCDIST is a text-dependent system, which requires a transcription to be processed alongside the
292 audio sample we are classifying. However, a text-dependent system here is defined as one that
293 requires a transcription, but the speech can be spontaneous (as discussed in Brown (2018)). In some
294 works, text-dependent systems only refer to those where the spoken content of the test samples and
295 the training samples match. This is one of the key features that separates Y-ACCDIST from other
296 ACCDIST-based recognisers found in Huckvale (2004, 2007) and Hanani et. al., (2013). The initial
297 experiments will allow us to compare the performance of this approach on the IVAr dataset on both
298 read speech (where the spoken content is matched across training and test data), and spontaneous
299 speech (where the spoken content does not match across speech samples).

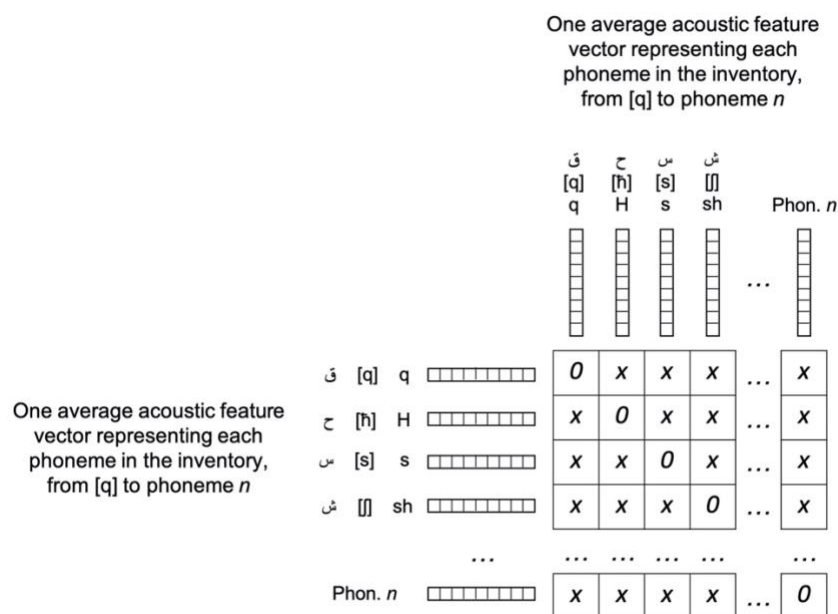
300

301 **3.1 System Description**

302

303 For each speaker in the IVAr dataset, we take a speech sample and a transcription and pass them
304 through a forced aligner (developed in-house using the Hidden Markov Model Toolkit (HTK)
305 (Young et al, 2009)) to estimate where each phone in the sequence is produced in the sample. Given
306 a speech recording and a phonemic transcription of that recording, the aligner extracts acoustic
307 features from across the speech sample and estimates where each phone is in the signal, i.e.
308 producing an estimated time alignment of the phone sequence. Some forced aligners, particularly
309 those that are widely available, have ready-trained acoustic models for a given language that may
310 provide multiple options for a phonemic transcription of a given word. The specific phone labels
311 attributed to a speech sample will therefore be partly determined by the acoustics of the segments in
312 the speech sample, and how they compare against the pre-trained acoustic models of the forced
313 aligner. For the present study, however, we created a bespoke lexicon containing all lexical items in
314 the analysed data subset (described in section 4.1), based on the phoneme inventory of Modern
315 Standard Arabic (MSA). To achieve this, we generated a cross-dialectal lexicon from the dialect-
316 specific transcripts made available with the IVAr corpus, which was manually edited by an Arabic
317 speaker to replace dialect-specific phoneme labels with MSA phoneme labels; for example, a dialect-

318 specific entry for the word ‘heart’ such as [galb] or [2alb] appears in the bespoke lexicon as [qalb].
319 We then used the IVAr dataset itself to train speaker-specific acoustic models for the MSA phoneme
320 categories in the bespoke lexicon, which the aligner used to estimate where each phoneme is in the
321 sample. We initialised the models by “flat-starting”; that is, we imposed evenly spaced notional
322 phoneme boundaries on the speech samples as a starting point. We then repeatedly applied an
323 Expectation-Maximization algorithm which iteratively adjusted the placement of these boundaries to
324 more accurately segment the sample according to phone segments. More reliable boundaries should
325 be reflected in the production of increasingly stable acoustic models during this process. Performing
326 forced alignment in this way was possible because we had enough speech per speaker to do so. This
327 allows us to impose just one set of MSA symbols on the range of different productions that different
328 speakers may produce. This lays the foundations for our method of dialect classification.
329
330 Using these estimated time boundaries between phones in the sequence, a vector of Mel Frequency
331 Cepstral Coefficients (MFCCs) (Davis & Mermelstein, 1980) was extracted at the midpoint to
332 acoustically represent each phone. The MFCCs used in this work consist of 12 coefficients. Larger
333 MFCC vectors have been trialled in past work (Brown, 2014), but 12 coefficients were shown to
334 provide sufficient information. An average MFCC was calculated for each phoneme category in
335 MSA from these midpoint acoustic features. The result of this is that we have the phoneme inventory
336 represented by average acoustic features (one per phoneme) for the speaker. By using midpoint
337 features, this approach overlooks temporal differences that might exist between dialects. This is a
338 factor to keep in mind when interpreting the results.
339
340 Using this set of averaged acoustic features, we calculated the Euclidean distance between all
341 phoneme-pair combinations that are possible within the phoneme inventory. This was achieved by
342 computing the Euclidean distances between all the possible pairs of average MFCC vectors that
343 represent each phoneme. We can organise this in a matrix (for clarity, this is illustrated below in
344 Figure 1).
345



346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367

Figure 1. A demonstration of how a speaker-specific matrix is calculated with the whole segmental inventory from [q] to phoneme n . The ‘0’/‘x’ symbols represent the Euclidean Distance for that pair.

The resulting set of Euclidean distances is expected to encapsulate the range of phonetic realisations that are associated with a speaker’s pronunciation system (or accent). This matrix of distances is our model of a speaker’s accent. Using British English accents as an example, typical speakers in Northern England will produce similarly realised vowels for FOOT and STRUT (both realised as [ʊ]), whereas typical speakers in the South of England will produce different vowel realisations (FOOT would be produced as [ʊ], while STRUT would be more likely to be produced as [ʌ]). A parallel example for Arabic would arise for consonants; an Arabic speaker from Egypt will typically realise the target sound < ق > [q] in the same way as target < ء > [ʔ], whereas an Arabic speaker from Morocco will more frequently produce these two target sounds ([q]~[ʔ]) as two separate categories. One of the Euclidean distances in the matrix is expected to reflect this accent-specific feature of the speaker. An entire matrix is therefore expected to contain numerous accent-specific features of this kind. Simultaneously, by computing intra-speaker distances in this way, we should eliminate other information embedded within the speech signal that does not necessarily assist in the accent classification task. For example, the distance between the FOOT and STRUT vowels for a typical Northern male speaker and a typical Northern female speaker should be equally small; similarly, the distance between targets [q]~[ʔ] will all be equally small for a typical Egyptian female speaker and a typical Egyptian male speaker.

368 We performed the above procedure on the speech samples and transcriptions of all our training
369 speakers. The resulting speaker-specific matrices are then fed as features into a Support Vector
370 Machine (SVM) classifier (Vapnik, 1998). It is also possible to make use of Deep Neural Networks
371 (DNNs) as classification mechanisms in these sorts of experiments. However, DNNs are much more
372 suited to extremely large datasets of thousands of data samples. SVMs also tend to require larger
373 datasets, but they are not as “data-hungry” as DNNs.

374

375 Within the SVM, which acts as multi-dimensional space, a *one-against-the-rest* rotation is
376 implemented for classification. In turn, each accent group of training speakers becomes the ‘one’,
377 while the speakers for all other accent categories are collapsed into a single group that form ‘the rest’.
378 An optimal hyperplane (i.e. a separating boundary within multidimensional space) is computed on
379 each rotation to achieve the best separation between these two groups of speakers². To classify an
380 unseen speaker, we form a matrix model for that speaker as described above, and this model is
381 presented to the SVM on each rotation. The accent category of the unseen speaker is determined by
382 the clearest margin it forms with the hyperplane in each of these rotations.

383 **4. Experiments**

384 A sequence of experiments was conducted in the commonly implemented *leave-one-out cross-*
385 *validation* setup, where each speaker became the test speaker, in turn, while the remaining speakers
386 in the dataset were used to train the Y-ACCDIST system. This was in an effort to maximise the
387 number of training speakers.

388

389 **4.1 Segmental Modelling**

390

391 The above process was conducted for read speech recordings from the speakers (where speakers were
392 asked to read the same scripted dialogue and story) and also spontaneous speech as a comparison of
393 performance on the two modes of speech. As we pointed out above, a transcription must accompany
394 the recordings. Most, but not all, speakers’ spontaneous speech samples have been orthographically
395 transcribed. For these experiments, we have therefore used data for both read speech and spontaneous
396 speech experiments from a subset of the speakers (reduced from 96 speakers to 86 speakers). This

² While SVMs can be very useful in classification problems, they are susceptible to ‘overfitting’, particularly on moderate-sized datasets. In these experiments, while overfitting is a risk, we have used a linear kernel, and have also set the regularization parameter to tolerate some errors during training. The controlled nature of the dataset also mitigates against overfitting as it provides less “noise” and therefore fewer overfitting opportunities.

397 results in an imbalance in the number of speakers for different dialect groups, though an even gender
 398 balance was retained within each group. Table 2 shows the number of speakers per dialect group in
 399 our analysed subset, along with the volume of data in minutes used in the experiment (with silences
 400 removed), by dialect, gender and speech style.

401

402 Table 2. Number of speakers per accent category in the data subset, with total duration (rounded up
 403 to the nearest whole minute) of speech data used in training and/or testing (silences removed).

Dialect Group	Code	Speakers		Scripted data (mins)			Unscripted data (mins)		
		Female	Male	Female	Male	Total	Female	Male	Total
Egyptian (Cairo)	egca	4	4	15	12	26	15	7	20
Iraqi (Muslim Baghdadi)	irba	6	6	15	13	28	21	15	36
Jordanian (Karak)	joka	6	6	15	15	30	21	17	37
Kuwaiti (Urban)	kwur	6	6	14	14	28	18	23	41
Moroccan (Casablanca)	moca	6	6	14	14	29	21	44	65
Gulf (Buraïmi, Oman)	ombu	6	6	16	15	31	18	12	30
Syrian (Damascus)	syda	3	3	17	15	32	12	17	29
Tunisian (Tunis)	tuns	6	6	14	13	27	23	21	44

404

405

406 Table 3 provides the means and standard deviations of the amount of speech (in seconds) per speaker
 407 used in model training and/or testing in this study.

408

409 Table 3. Mean/standard deviation of speech in seconds per speaker in the data subset by speech task.

410

	Speech task	Mean amount of speech per speaker (seconds)	Standard deviation per speaker (seconds)
Read Speech (scripted)	Story	72.09	10.74
	Read sentences	73.36	9.64
	Total (read)	145.45	16.54
Spontaneous Speech (unscripted)	Free conversation	77.28	46.72
	Map task	70.01	59.64
	Retold folk tale	65.11	17.26
	Total (spontaneous)	212.41	102.48

411

412

413 We present the overall results and their corresponding confusion matrices in the subsections below.

414

415 4.1.1 Read Speech

416 The read speech data used for these experiments come from a scripted role-play dialogue which was
 417 designed to elicit a number of different sentence types, including declarative statements (*dec*),

418 yes/no-questions (*ynq*), wh-questions (*whq*) and coordinated questions (*coo*, also known as
 419 alternative questions, of the form “is it X or Y?”). The sentences were designed to control the
 420 segmental content and prosodic structure of the last lexical item in each utterance, so that it contained
 421 mostly sonorant sounds (to facilitate pitch tracking) and the position of the stressed syllable was
 422 systematically varied over the last three syllables of the word. A set of sample yes/no-questions
 423 elicited in one dialect (here, Jordanian Arabic) are provided in Table 4.

424
 425 Table 4. Sample set of yes/no questions (in *joka*) elicited using the scripted dialogue.
 426

Code	Target sentence	
ynq1	ruħt l-nna:di l-‘jamani	<i>Did you go to the Club Yemeni?</i>
Ynq2	l-zawa:ʒ l-madani raħ jku:n fi-l-mabna l-‘baladi	<i>Will the civil wedding be in the <u>municipal</u> office?</i>
Ynq3	ga:balu baʕid ^o ʕan tʕari:g ‘ze:na	<i>Did they meet each other through <u>Zena</u>?</i>
ynq4	jaʕni raħ tzu:r ʕuxutha la‘ja:li	<i>Do you mean she will visit her sister <u>Layali</u>?</i>
ynq5	yaʕni tʕarrafit ʕale: fi-l-mat ^o ʕam illi fi-l-‘mo:l	<i>Do you mean they met in the restaurant in the <u>mall</u>?</i>
ynq6	wa:lid nabi:l raħ jku:n maw‘zu:d	<i>Will Nabil’s parents be <u>present</u>?</i>

427
 428
 429 For the story task participants read a monologue narrative folk-tale ‘Guha and the banana seller’,
 430 adapted from a story in Abdel-Massih (2011) and adjusted to contain appropriate lexical and
 431 grammatical forms for each target dialects. The story is typically realised by speakers in 40-45
 432 prosodic phrases or breath groups. As noted above all scripted material was presented in Arabic
 433 script using local spelling conventions.

434
 435 Although considerable effort went into making the reading material as comparable as possible across
 436 dialects, the scripts read by speakers across dialects did not necessarily match word-for-word. This is
 437 because certain words are simply not shared across dialects so there is some lexical and grammatical
 438 variation across speech samples. We acknowledge that this may have weighted the result to some
 439 extent in that a small number of the phones were produced in specific phonological environments and
 440 this varies according to dialect. However, we do not expect this to be the leading factor in
 441 determining accent as we are cutting out individual phonemes and taking acoustic values from the
 442 midpoints of these segments.

443
 444 Using the read speech data, the Y-ACCDIST system achieved 95.3% correct. The accompanying
 445 confusion matrix is presented in Table 5.

446

447 Table 5. Confusion matrix for Arabic dialect classification task using the segmental models on the
 448 full read speech dataset (scripted dialogue/story).

449

		Predicted Labels							TOTAL	
		egca	irba	joka	kwur	moca	ombu	syda		Tuns
True Labels	egca	12 (100%)	0	0	0	0	0	0	0	12
	irba	0	11 (91.7%)	1 (8.3%)	0	0	0	0	0	12
	joka	0	0	8 (100%)	0	0	0	0	0	8
	kwur	0	0	0	12 (100%)	0	0	0	0	12
	moca	0	0	0	0	12 (100%)	0	0	0	12
	ombu	0	0	0	0	0	12 (100%)	0	0	12
	syda	0	0	3 (50%)	0	0	0	3 (50%)	0	6
	tuns	0	0	0	0	0	0	0	12 (100%)	12

450

451 The least successful result in this experiment was for the Syrian group of 6 speakers, 3 of whom were
 452 identified as Jordanian. We note that all of the Syrian speakers were resident in Jordan at time of
 453 recording.

454

455 4.1.2 Spontaneous Speech

456

457 Using the spontaneous speech data for modelling, the Y-ACCDIST system achieves 77.9% correct
 458 (67/86). The accompanying confusion matrix is presented in Table 6.

459

460 Table 6. Confusion matrix for Arabic dialect classification task using the segmental models on
 461 spontaneous speech.

		Predicted Labels							TOTAL	
		Egca	irba	joka	kwur	moca	ombu	syda		tuns
True Labels	egca	6 (75%)	0	0	0	1 (12.5%)	0	1 (12.5%)	0	8
	irba	1	8 (66.7%)	1 (8.3%)	2 (16.7%)	0	0	0	0	12
	joka	0	1 (8.3%)	8 (66.7%)	1 (8.3%)	0	1 (8.3%)	1 (8.3%)	0	12
	kwur	0	1 (8.3%)	0	9 (75%)	0	2 (16.7%)	0	0	12
	moca	0	0	0	0	12 (100%)	0	0	0	12
	ombu	0	0	1 (8.3%)	3 (25%)	0	8 (66.7%)	0	0	12

	syda	1 (16.7%)	0	1 (16.7%)	0	0	0	4 (66.7%)	0	6
	tuns	0	0	0	0	0	0	0	12 (100%)	12

462

463 We should also note that for the spontaneous speech condition, speakers did not necessarily produce
464 the same quantity of speech (durations of speech per speaker generally ranged from 4 to 8 minutes),
465 and so there is variability across the dataset in this respect.

466

467 From the above two results, we can get an indication of the detriment to performance that content-
468 mismatched data has. This is because the different phone tokens are produced in different
469 environments which is likely to introduce an additional element of variability that is not present in
470 the read speech condition. We should also bear in mind that there is a smaller number of speakers
471 available for training the system for some dialect categories which is also likely to impact on the
472 result.

473

474 **4.2. Prosodic Modelling**

475

476 As discussed above, the Y-ACCDIST-based approach has allowed us to isolate the segmental level
477 of analysis and ignore other information embedded within the acoustic signal that might distract
478 away from cues useful to the accent recognition task. In this part of the study, we aim to transfer the
479 principles of the Y-ACCDIST modelling procedure to the prosodic level of analysis to see whether
480 we can confirm previous prosodic analysis of these same data, which indicated that there are prosodic
481 patterns that are typical of speakers of one or more Arabic dialects but different from patterns
482 observed in a parallel context in one or more other dialects (Hellmuth 2018).

483

484 **4.2.1 Organisation of Prosodic Data**

485

486 The read speech data from the IVAr scripted dialogue include the sentence types presented in Table
487 7, elicited because these may be characterised by different prosodic patterns between sentence types
488 within one dialect, and/or by different prosodic patterns between dialects within one sentence type.

489

490 Table 7. Sentence types elicited using the IVAr corpus role-play scripted dialogue.

491

Code	Sentence type	
dec	declarative	response to an open question (e.g. 'what's new?')
whq	wh-question	question using wh-word such as who or what
ynq	yes/no-question	polar question inviting a yes or no answer

coo	coordinated question (or alternative question)	question between two alternatives (e.g. ‘is it X or Y?’)
inf	information focus	statement produced in response to a wh-question
con	contrastive focus	statement produced in response to a yes/no-question
idf	identification focus	statement produced in response to a coordinated question

492

493 For this reason, it is only these sentences extracted from the scripted dialogue (sd) that are being used
494 in these experiments that compare the prosodic Y-ACCDIST system with the segmental Y-
495 ACCDIST system, with the read story (sto) data removed from the dataset. A set of results for each
496 of these system configurations are therefore presented within this section, where each has been
497 trained and tested on only the sentences extracted from the scripted dialogue.

498

499 **4.2.2 Integration of Prosodic Data into the Y-ACCDIST System**

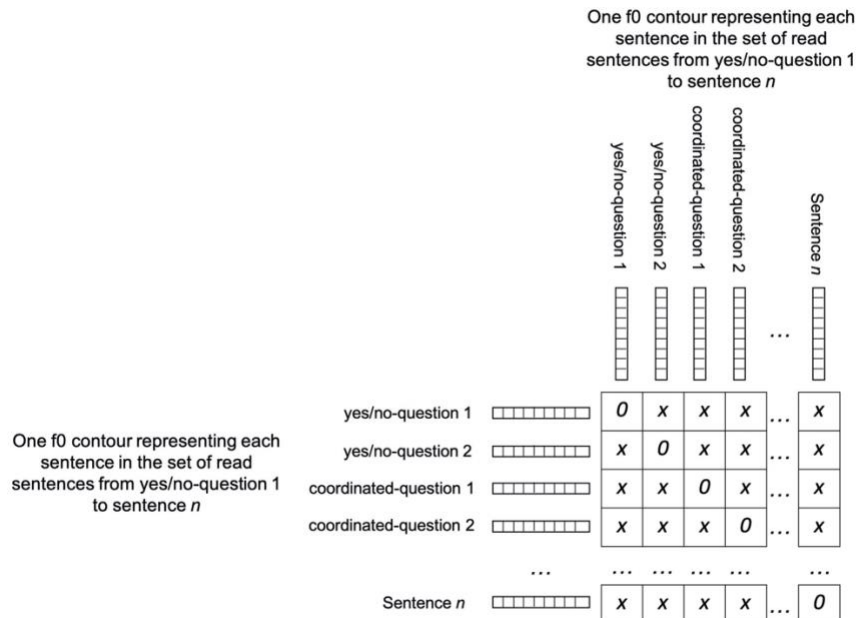
500

501 To provide the system with prosodic information, we calculated Euclidean distances between the f0
502 contours of all the possible pairs of read sentences available for each speaker. It is this collection of
503 Euclidean distances between f0 contours that is expected to characterise the intonational patterns of a
504 speaker. While it may not be immediately obvious what sorts of Euclidean distances are likely to
505 occur between these f0 contours, it is expected that logging the similarities and differences between
506 f0 contours in this way will express any systematic similarities and differences in intonational
507 patterns within and between dialects. For example, for Syrian speakers, who more frequently use a
508 rising contour in declarative sentences than is observed in other Arabic dialects (Hellmuth, 2020), the
509 f0 contour between a particular declarative sentence X and a particular polar interrogative sentence Y
510 might be reasonably expected to be more similar to one another than for a speaker of another dialect.

511

512 We used the read portion of the corpus in which all speakers produced more or less the same
513 sentences. F0 contours were extracted by marking out 50 equally distributed points throughout an
514 utterance and extracting the f0 at those points in the signal. Of course, at the points in the signal
515 where there is no voicing, extracting f0 was not possible. This therefore reduced these vectors down
516 to a size which was slightly smaller than 50, and the same reduction was performed on the f0 vector
517 that it was being compared against. The result of this was a Y-ACCDIST matrix that reflected the
518 intonation realisations of the “prosodic contour inventory” that the dataset allowed. Like the default
519 segmental configuration explained in Section 3.1, this modelling method has the advantage of
520 normalising against factors such as gender. By making intra-speaker calculations in this way, the
521 model only characterises the shapes of a speaker’s f0 contours, and so, regardless of whether the
522 speaker has a relatively high or low f0 range on average, any dialect-specific contour shapes should

523 be expressed in the matrix. Figure 2 provides an illustration of the prosodic modelling procedure that
 524 can be compared with the segmental modelling procedure.



525

526

527 Figure 2. A diagram to demonstrate how a matrix is formed using f0 contours, rather than acoustic
 528 feature vectors, with the set of sentence types in the IVAr dataset, from yes/no-question 1 to sentence
 529 n . The '0'/'x' symbols represent the Euclidean Distance for that pair.

530

531 One key difference between these prosodic models and the segmental models is that there is no
 532 averaging of vectors in the construction of the matrices. While it is expected that intonation is
 533 affected by sentence type in Arabic, it is also affected by a range of other discourse factors, such as
 534 information structure (topic, givenness and focus) (Krifka 2008) as well as the interactional context
 535 (Walker 2014). The read speech sentences in the corpus were elicited at different points throughout a
 536 scripted dialogue, and thus appear in subtly different discourse contexts with varying information
 537 structure. This meant that it would be artificial to try to model an average f0 contour for each whole
 538 sentence type. We have therefore treated each individual sentence that was produced as a single
 539 category in the construction of the speaker-specific matrices. Each individual sentence was elicited in
 540 the same discourse context (i.e. position in the scripted dialogue) from each speaker in each dialect.
 541 Overall, this means that we are restricted to performing these experiments on a dataset where
 542 speakers produce the same spoken content. We return to this point further below in the Discussion.

543

544 Using these prosodic matrices to represent each speaker, we followed the same experimental
 545 procedure as for the segmental experiments described above to achieve a recognition rate and
 546 confusion matrix. The recognition rate we achieved in this configuration was 52.1% correct. The
 547 confusion matrix for this task is shown in Table 8.

548
 549 Table 8. Confusion matrix of Y-ACCDIST’s performance on the IVAr dataset using the prosodic
 550 models on the read speech data subset (scripted dialogue only).

		Predicted Labels							
		egca	irba	joka	kwur	moca	ombu	syda	tuns
True Labels	egca	8 (66.7%)	2 (16.7%)	0	0	0	0	1 (8.3%)	1 (8.3%)
	irba	1 (8.3%)	6 (50%)	3 (25%)	0	0	1 (8.3%)	1 (8.3%)	0
	joka	1 (8.3%)	1 (8.3%)	2 (16.7%)	3 (25%)	0	2 (16.7%)	2 (16.7%)	1 (8.3%)
	kwur	0	0	1 (8.3%)	9 (75%)	1 (8.3%)	0	0	1 (8.3%)
	moca	0	0	1 (8.3%)	1 (8.3%)	7 (58.3%)	2 (16.7%)	1 (8.3%)	0
	ombu	0	0	1 (8.3%)	0	3 (25%)	5 (41.7%)	2 (16.7%)	1 (8.3%)
	syda	1 (8.3%)	1 (8.3%)	2 (16.7%)	0	2 (16.7%)	3 (25%)	2 (16.7%)	1 (8.3%)
	tuns	0	0	0	0	0	0	1 (8.3%)	11 (91.7%)

551
 552 We can compare these results with the segmental system’s results that were produced using the same
 553 subset of read data, which was 93.75% correct, and the confusion matrix for this is shown in Table 9.
 554
 555 Table 9. Confusion matrix of Y-ACCDIST’s performance on the IVAr dataset using the segmental
 556 models on the read speech data subset (scripted dialogue only; this is the same dataset that was used
 557 to build and train the prosodic system).

		Predicted Labels								TOTAL
		egca	irba	joka	kwur	moca	ombu	syda	tuns	
True Labels	egca	12 (100%)	0	0	0	0	0	0	0	12
	irba	0	11 (91.7%)	0	0	0	0	1 (8.3%)	0	12
	joka	0	0	11 (91.7%)	1 (8.3%)	0	0	0	0	12
	kwur	0	0	0	11 (91.7%)	0	0	1 (8.3%)	0	12
	moca	0	0	0	0	12 (100%)	0	0	0	12
	ombu	0	0	0	0	0	12	0	0	12

							(100%)			
syda	0	0	3 (25%)	0	0	0	0	9 (75%)	0	12
tuns	0	0	0	0	0	0	0	0	12 (100%)	12

559

560 For the results from the prosodic system, although accuracy in the classification task varies
561 considerably, all dialects are recognised by prosodic contour alone at above chance levels (if chance
562 is 12.5% correct). The four ‘best’ recognised dialects are tuns (91%), kwur (75%), egca (67%) and
563 moca (58%). The four ‘worst’ dialects are irba (50%) and ombu (42%), followed by syda and joca
564 (both at 16%). These best and worst groupings resemble those observed in the segmental
565 classification task on spontaneous speech data (Table 6): tuns and moca (100%), kwur and egca
566 (75%), then irba, ombu, syda and joca (all on 67%).

567

568

569 **5 Feature contributions to Arabic dialect classification**

570

571 One obvious area of interest is identifying the specific linguistic units (i.e. phonemes or sentences)
572 that are contributing most to distinguishing between the dialect varieties. This section presents an
573 attempt to access this information within the inner workings of the systems. It builds on a similar
574 attempt to achieve this in Brown and Wormald (2017), which simply applied ANOVA to Y-
575 ACCDIST models of different British English speakers to reveal which phoneme-pairs were
576 estimated to distinguish between four accent varieties. The work here makes use of the machine
577 learning mechanisms implemented in this study to help identify any linguistic units or categories
578 which might be key features in separating the varieties.

579

580 For both the segmental and prosodic systems, SVMs were used as the classification mechanism.
581 SVMs assign different weights to the features of the models or representations they are learning.
582 These weights help with the separation of the groups in the SVM which, in turn, should help to
583 achieve better classification results. The weights can also be used for a feature selection process
584 called *Recursive Feature Elimination* where they are used to rank the features by their weights, and
585 then the weakest features are removed (either iteratively or as a specified amount, n). By removing
586 features expected to be less useful to the task, it is thought that the classification performance of a
587 system could be improved. One option is to run experiments that iteratively remove the weaker
588 features, and observe what the effect is on classification performance. However, a more efficient, and

589 perhaps more direct, way of accessing information about the estimated value of the different features
590 in a feature set is to look at the full set of weights that the SVM has assigned.

591

592 In the present case, the Euclidean distances that make up the Y-ACCDIST matrices are assigned
593 different weights, according to those that are estimated to discriminate the different dialect groups
594 among the training data. Because we have applied a linear kernel in the SVM in these experiments, it
595 is possible to look more closely at the weights that the different Y-ACCDIST features are assigned
596 by the SVM, and then use them to make estimations around which features are most useful to the
597 task of discriminating Arabic dialects. This was carried out for both the segmental system and the
598 prosodic system to assess whether this method allows us to pick out any particular phonemes or
599 sentence types that are particularly useful in distinguishing between the dialects. We were keen to use
600 as much data as possible in order to observe the most reliable indications of sociophonetic variation
601 within this dataset, so we chose to include the read speech from both the scripted dialogue and story
602 tasks in this analysis using the segmental system.

603

604 Using the Y-ACCDIST modelling method, however, this process of drawing on SVM weights to
605 observe individual feature contribution is not wholly straightforward. The speaker representations
606 that we feed into the SVM are values that are computed between pairs of phonemes or pairs of
607 sentence types (i.e. the values represented by “x” in Figures 1 and 2), rather than a single value
608 mapping directly on to an individual phoneme or sentence type. While the modelling method has
609 been shown to be very strong, these pairs are very difficult to disentangle to be able to observe the
610 contribution of individual phonemes and sentence types. Clear and obvious patterns may therefore
611 not emerge. Nevertheless, it is still of interest to see whether this method yields any insight into
612 feature contributions and so we observe the values that we can in this section. To estimate feature
613 contribution, we have accumulated all of the absolute weight values assigned to all the pairs of
614 features that a single feature belongs to, and we reflect these values in boxplots.

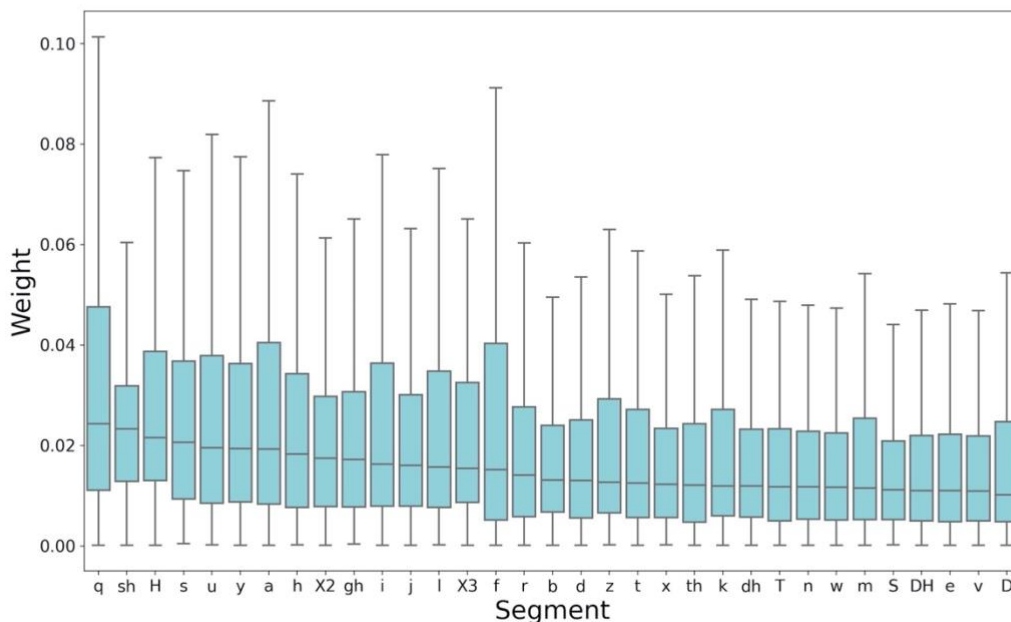
615

616 **5.1 Segmental feature contributions**

617

618 Figure 3 shows this tentative measure of which phonemes appear to contribute the most weight to the
619 task of distinguishing the eight dialect groups.

620



621

622

623 Figure 3. Boxplots to represent the distributions of feature weights associated with each phoneme
 624 segment³.

625

626 The segments have been ordered according to the median, where we find those segments that are
 627 estimated to have the greatest contribution overall to the left, and the segments that are estimated to
 628 make the least contribution are positioned to the right. To reiterate, because of the pairwise nature of
 629 the modelling method, the visual evidence of an individual segment's contribution to a classification
 630 task is somewhat diluted. We should also keep in mind that segments are represented by midpoints
 631 and so segments that differ in terms of temporal characteristics (rather than quality characteristics)
 632 are less likely to emerge in this analysis.

633

634 The highest ranked phoneme in the feature weights boxplot is (q) /q/, matching systematic variation
 635 in the realisation of this sound across Arabic dialects (Al-Essa 2019). Similarly, the relatively high
 636 ranking of (j) /dʒ/ matches the status of that sound as a **known** locus of variation between dialects.
 637 For both these sounds, variation between dialects in their realisation is well-documented in the

³ A key for the symbols used can be found here: <https://reshare.ukdataservice.ac.uk/852878/15/transliteration.pdf>

638 research literature, and they frequently appear as a variable in variationist sociolinguistic studies of
639 individual dialects.

640
641 In contrast, most of the other highly ranked sounds in Figure 3, by feature weight value, are **classes**
642 **of sound** which have received little attention in the research literature on Arabic sociolinguistic or
643 dialectal variation. These include both fricatives, notably (sh) /ʃ/, (s) /s/ and (H) /ħ/, and vowels /a/
644 and /u/ (with /i/ not far behind). Variation across Arabic dialects in the gradient phonetic realisation
645 of fricatives and vowels is under-researched and as a result not yet fully documented, but those **few**
646 studies that exist are nevertheless consistent with the patterns seen here. **For vowels, Alghamdi**
647 **(1998) reports a complex pattern of differences in values of the first formant (reflecting vowel height)**
648 **in data from Saudi, Egyptian and Sudanese speakers, for [a, a:, i, i:, u, u:]. Al-Tamimi & Ferragne**
649 **(2005) similarly report a difference in the size of the i~a~u vowel space between Moroccan and**
650 **Jordanian Arabic (with comparison also to French); furthermore, they show that Principal**
651 **Component Analysis on a simple measure of vowel space size (using between-vowel-vectors for the**
652 **first and second formants) yielded 88% correct classification of the three languages. For fricatives,**
653 **Alsabhi et al (2020) report a main effect of *dialect* in models of standard acoustic measures of overall**
654 **spectral shape (centre of gravity and peak Hz) for /s/ and /s^h/, in experimental data elicited alongside**
655 **the IVAr corpus from the same speakers and dialect groups examined in the present study.**

656
657 In addition, we note that dialectal variation in realisation of (q) and (j) in Arabic can be characterised
658 as sociolinguistically salient: the variation is above the level of awareness among speakers and may
659 serve as a stereotype of particular dialects (Ateek & Rasinger, 2018). To our knowledge no studies
660 have systematically investigated the relative sociolinguistic salience of different linguistic features in
661 Arabic dialects. **However, these feature weights suggest that there may be gradient sociophonetic**
662 **features related to fine-grained phonetic realisation of fricatives and/or vowels, which may be below**
663 **the level of phonological awareness among speakers and thus not perceived as dialect stereotypes,**
664 **but which nevertheless contribute to automatic dialect classification.**

665
666 We have already indicated that what draw can be drawn from the feature weights for this particular
667 modelling method is rather limited. Having said this, there are patterns emerging that align with some
668 expectations based on previous research, as well as patterns that have perhaps previously gone
669 virtually unnoticed. The method has therefore focussed attention on potential features of interest and
670 motivated future research objectives in relation to Arabic dialects.

671

672 **5.2 Prosodic feature contributions**

673

674 We also produced the equivalent boxplot visualisations for the prosodic system models to determine
675 whether any sentence types were particularly influential in distinguishing between the dialects. Our
676 conclusion here is that there seems to be very little to report, as we found very flat and invariable
677 distributions among the different features. This will in part be due to the particular modelling method
678 (as we have already said, a pairwise modelling approach is not the best foundation for reporting the
679 value of individual segments). This will also be due to the fact that these features are not particularly
680 powerful dialect discriminators, as the classification results have already demonstrated.

681 The combination of the lower classification result and the invariable boxplots indicates that we can
682 expect to find a lot of variability among the intonation contours within the dialect groups. As noted
683 earlier, however (section 4.2.2) this is exactly what we would expect, since prosodic contour
684 realisation varies not only according to semantic categories such as question versus statement, but
685 also due to the information structure and the wider discourse and interactional context. The lack of
686 feature weight information thus supports the methodological choice to have Y-ACCDIST use
687 individual sentences (realised at the same position in the dialogue sequence) as the unit of analysis.

688

689

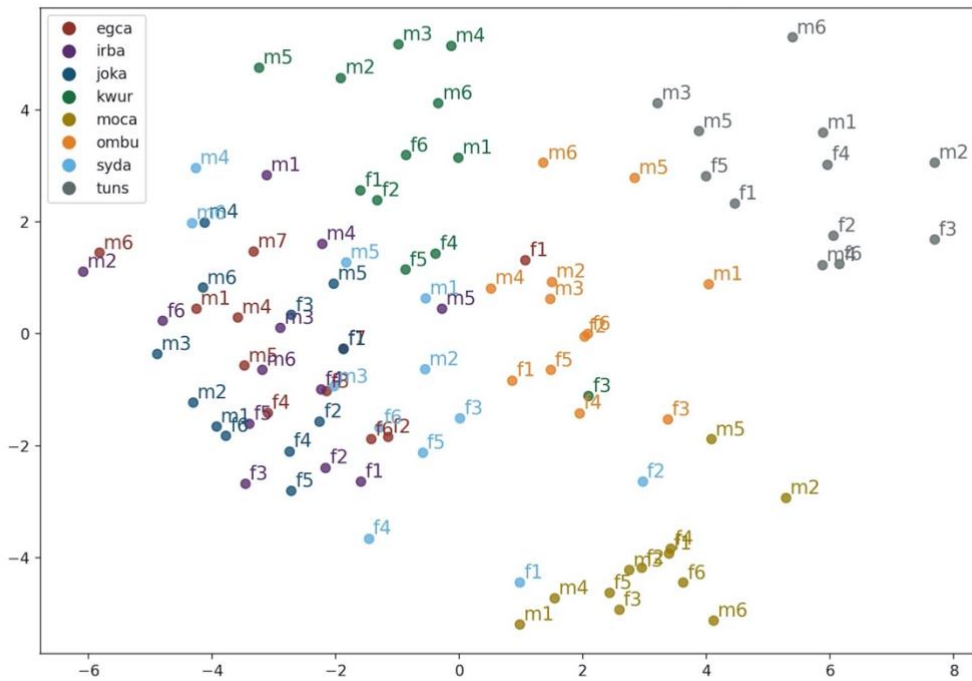
690 **6 Comparison of visualisations of segmental and prosodic models**

691

692 It is also possible to compare visualisations of the two modelling methods for the IVAr speakers in
693 the read speech scripted dialogue data subset. Having modelled the speakers in this data subset, we
694 performed multi-dimensional scaling (MDS) on the data, once under the segmental configuration and
695 once under the prosodic configuration. This allows us to observe any interesting clusters of speakers
696 for each of the levels of analysis in isolation. These are presented in Figure 4 and Figure 5
697 respectively.

698

699

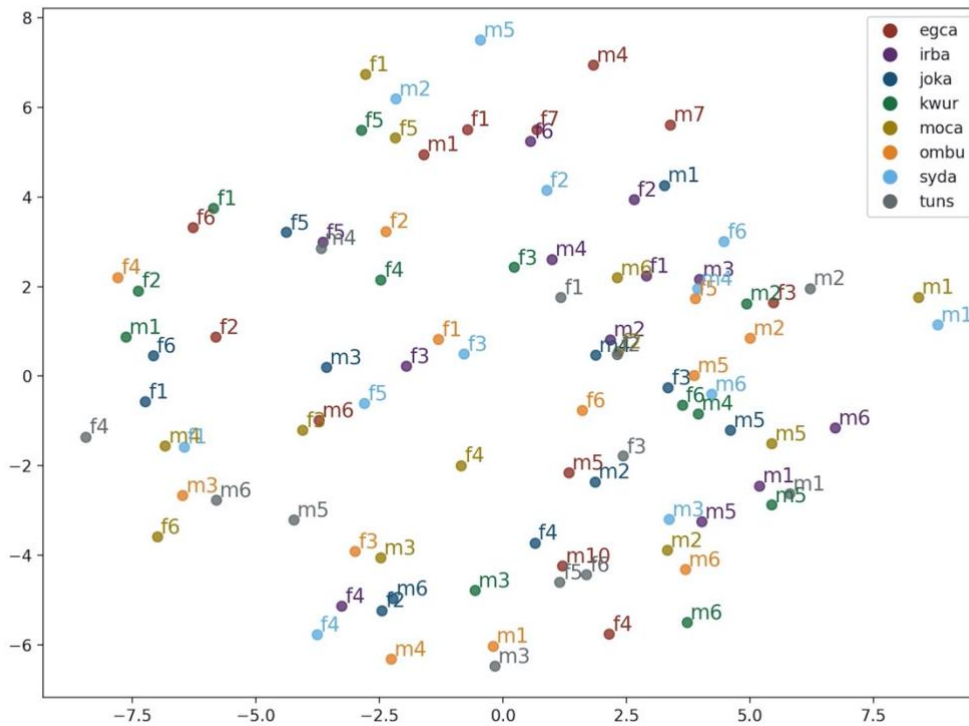


700

701 Figure 4. Multidimensional Scaling of the IVAr dataset based on segmental Y-ACCDIST modelling
 702 of speakers

703

704



705

706 Figure 5. Multidimensional Scaling of the IVAr dataset based on prosodic Y-ACCDIST modelling of
 707 speakers.

708

709 Figure 4 shows clear groupings for most of the individual dialect groups, showing that the segmental
 710 level of analysis is a good unifier of speakers of the same dialect. This is consistent with the very
 711 high classification result that this version of the system achieved. The clustering reflects both the
 712 geographical spread of dialects and their position in the Arabic dialect continuum. The clusters of
 713 speakers from Egypt, Iraq and the Levant (Jordan and Syria) overlap somewhat, towards the left of
 714 the plot, but corresponding to their more central geography and position in the middle of the dialect
 715 continuum. The Gulf dialects (Kuwait/Oman) are distinct from each other, matching their positions at
 716 extreme north and southern ends of the Gulf Arabic dialect group, but both equally separate and
 717 distinct from the central dialects. Similarly, the North African dialects (Tunisia/Morocco) are clearly
 718 separated from each other, again matching their geographical and dialectal separation within the
 719 Maghreb group, but are both equally separate and distinct from the central dialects, and placed at a
 720 greater distance from the central group than the Gulf dialects, reflecting the clear east/west divide
 721 noted in section 2.1.

722 Figure 5 for the prosodic model shows a less clear clustering of individual dialect groups. The
723 relatively tight cluster of Tunisian speakers just right of centre of Figure 5 perhaps aligns with the
724 very high classification rate achieved for Tunisian speakers in the prosodic experiments, and
725 interestingly, the Egyptian speakers seem to form a somewhat more consistent cluster compared to
726 the other dialect groups. Such a clustering for these two dialects would be consistent with some of the
727 more distinctive prosodic patterning in these dialects found in previous work by the second author,
728 namely a distinctive rise-plateau contour in yes/no-questions in Tunisian Arabic (Hellmuth, 2018)
729 and overall higher frequency in the distribution of prosodic peaks in Egyptian Arabic compared to
730 other dialects (Hellmuth 2007, 2020).

731 **7 Discussion**

732 This paper has demonstrated approaches to analysing dialect variation that take into account whole
733 collections of features, rather than just focussing on a single feature and seeing how it varies across
734 different dialects. These approaches are reliant on there being an “inventory” of categories for the
735 models to work with. In the case of the segmental system, this is the phoneme inventory, and in the
736 case of the prosodic system, this is a range of different sentences produced at different points in a
737 scripted dialogue. These aggregate approaches therefore currently only offer a broad-brush account
738 of the variation in a dialect dataset on either one of these levels of analysis, rather than a detailed
739 account of exactly which feature is discriminating the different varieties.

740
741 In the case of the prosodic version of the system, we have presented the modelling approach only on
742 a subset of read speech data in which we could guarantee balance and control in the different
743 sentence types that we used. The corpus was originally designed for prosodic research to be
744 conducted on it and so other analysis on the prosodic variation in this dataset had already been done
745 which opened up the opportunity to corroborate results or to even uncover surprising findings.
746 Having tested the approach on these very controlled data and having discovered that it appears to
747 have some value in capturing the variation and distinguishing between the dialects, it is natural to
748 now consider how it could be transferred to spontaneous speech data. Given a dataset of spontaneous
749 speech recordings that have been tagged for key features such as sentence type and information
750 structure, we could evaluate the prospect of using this approach on spontaneous recordings. In
751 addition, it could be that a larger dataset than the one used in this work would be required to achieve
752 a more stable representation of the specific variation that exists in Arabic prosody.

753

754 There has been some other work that has actively sought to integrate prosody’s potential role in the
755 automatic classification of Arabic dialects. Biadisy and Hirschberg (2009) modelled spontaneous
756 speech utterances using a combination of features relevant for intonation and rhythm. They captured
757 a selection of pitch and rhythm values to represent whole utterances (e.g. capturing the variation in
758 pitch and vocalic proportion of an utterance). These were termed more “global” measurements. They
759 then went on to characterising utterances with “sequential prosodic features”, which logged various
760 characteristics of the pitch and intensity contours of utterances. On a broad four-way Arabic dialect
761 classification task, the “global” features alone achieved 60% accuracy, whereas when more
762 sophisticated sequential features were combined with them, they achieved 72% accuracy. Between
763 Biadisy and Hirschberg’s study and the present one, there are many differences to do with the size and
764 the dimensions of the datasets used, but it may be of interest in future to compare these two methods
765 like-for-like. One key difference is that Biadisy and Hirschberg applied their prosodic modelling
766 method to spontaneous speech, a natural next step for the Y-ACCDIST modelling method
767 implemented in the present study.

768

769 Although the Y-ACCDIST modelling approaches themselves are very adaptable and can feasibly be
770 used on large datasets of speech recordings, there is some manual preprocessing of the data (i.e.
771 either broad transcription or tagging) that is required before modelling, classification and
772 visualisation can take place. It could be possible to overcome this preprocessing by either
773 automatically transcribing or tagging a corpus, but this will inevitably introduce errors. Work on this
774 less labour-intensive version is currently ongoing.

775

776

777 **8 Conclusion**

778

779 In this paper we have focussed on automatically classifying speakers of Arabic into different dialect
780 groups and considered the systems’ outputs in the context of an interest in the variation among
781 Arabic dialects. We have demonstrated how these kinds of system can both reinforce what we know
782 about a set of linguistic varieties, but also how it could possibly illuminate new questions to pursue
783 around certain features. Previous work has shown that we can do this on one level of analysis, but
784 part of this work has demonstrated that sometimes performance might be too high for us to learn
785 about a set of dialects from the errors that a system makes. However, this paper has demonstrated that
786 it is possible to transfer similar modelling principles that have been used for one level of analysis
787 across to another. By isolating the segmental level of analysis and then the prosodic level in a similar

788 modelling framework, we can observe the contribution of each level of analysis to distinguishing
789 between the classification of a particular set of varieties. It is probably no surprise that the segmental
790 level outperforms the prosodic level in a simple dialect classification task, but the prosodic version of
791 the system showed a performance that sat well above the level we would expect if the system were
792 working by chance. This difference in performance between the two levels of analysis is likely to be
793 down to the fact that one forms models based on a full and well-established phoneme inventory and
794 the other makes use of a (partly arbitrary) list of target sentences. The former is both more fine-
795 grained and more controlled than the latter.

796
797 In the context of Arabic dialects, we were able to corroborate some of the findings surrounding the
798 prosodic system's outputs with past prosodic analyses conducted on the same data. The work here
799 was also able to indicate that Arabic has a wealth of sociophonetic variation to discover at the
800 segmental level, which is arguably under-explored in Arabic dialects. The detail of this segmental
801 variation cannot be accurately uncovered by the macro-level computational method implemented in
802 this work, but would require other more detailed methods to gain a richer understanding.

803

804 **Data availability**

805 The corpus that was used for this study can be found at the following reference:

806 Hellmuth, S., & Almbark, R. (2019). *Intonational Variation in Arabic Corpus (2011-2017)*. Retrieved from:
807 <http://reshare.ukdataservice.ac.uk/852878/>

808

809

810

811 **Author contributions**

812

813 GB took the lead in the technology development, computational methods and generating results. SH
814 took the lead in data collection, data processing and interpretation of results. The authors contributed
815 to the planning and writing of this article in equal measure.

816

817 **Acknowledgments**

818

819 The authors would like to thank Rana Alhusein Almbark for her valuable input towards the
820 beginning of the project when planning and executing the data processing required for this project.

821

822

823

824 **References**

- 825
- 826 Abdel-Massih, E.T. 2011. *An introduction to Egyptian Arabic*. Ann Arbor, University of Michigan.
- 827
- 828 Al-Essa, Aziza (2019) Phonological and morphological variation. In E. Al-Wer, & U. Horesh (Eds.),
829 *Routledge Handbook of Arabic Sociolinguistics*. Routledge. Pp151-168.
- 830
- 831 Al-Tamimi, J.-E., & Ferragne, E. (2005). Does vowel space size depend on language vowel
832 inventories? Evidence from two Arabic dialects and French. Paper presented at the 9th European
833 Conference on Speech, Communication and Technology (Interspeech 2005), Lisbon, Portugal.
- 834
- 835 Alghamdi, M. M. (1998). A spectrographic analysis of Arabic vowels: A cross-dialect study. *Journal*
836 *of King Saud University*, 10(1), 3-24.
- 837
- 838 Ali, A., Shon, S., Samih, Y., Mubarak, H., Abdelali, A., Glass, J., Renals, S. and Choukri, K. (2019).
839 The MGB-5 Challenge: Recognition and Dialect Identification of Dialectal Arabic Speech.
840 *Proceedings of Automatic Speech Recognition and Understanding*. (Sentosa, Singapore). 1026-1033.
- 841
- 842 Alsabhi, M., Bailey, G and Hellmuth, S. (2020) Cross-dialectal variation in phonetic realisation of
843 the emphatic contrast in Arabic fricatives. Paper presented at the 4th Arabic Linguistics Forum, 30th
844 June-2nd July 2020, Leeds.
- 845
- 846 Ateek, M., & Rasinger, S. M. (2018). Syrian or non-Syrian? Reflections on the use of LADO in the
847 UK. In I. M. Nick (Ed.), *Forensic Linguistics: Asylum-seekers, Refugees and Immigrants* (pp. 75-93):
848 Vernon Press.
- 849
- 850 Barkat, M., Ohala, J., & Pellegrino, F. (1999). Prosody as a distinctive feature for the discrimination
851 of Arabic dialects. *Eurospeech 99*, 395-398.
- 852
- 853 Beckman, M., & Elam, G. A. (1997). *Guidelines for TOBI Labelling (version 3.0 1997)*.
- 854
- 855 Beckman, M., Hirschberg, J., & Shattuck-Hufnagel, S. (2005). The original ToBI system and the
856 evolution of the ToBI framework. In S.-A. Jun (Ed.), *Prosodic Typology: The phonology of*
857 *intonation and phrasing* (pp. 9-54). Oxford: Oxford University Press.
- 858
- 859 Behnstedt, P., & Woidich, M. (2013). Dialectology. In J. Owens (Ed.), *The Oxford Handbook of*
860 *Arabic Linguistics* (pp. 300-325). Oxford: Oxford University Press.
- 861
- 862 Biadisy, F., & Hirschberg, J. (2009). Using Prosody and Phonotactics in Arabic Dialect Identification.
863 *Proceedings of Interspeech 2009 (Brighton, UK)*, 208-211.
- 864
- 865 Biadisy, F., Hirschberg, J., & Habash, N. (2009). Spoken Arabic Dialect Identification using
866 Phonotactic Modeling. *Proceedings of the EACL Workshop on Computational Approaches to Semitic*
867 *Languages. (Athens, Greece)*, 53-61.
- 868
- 869 Brown, G. (2014). Y-ACCDIST: An automatic accent recognition system for forensic applications.
870 Master Thesis, University of York, UK.
- 871

872 Brown, G. (2015). Automatic recognition of geographically-proximate accents using content-
873 controlled and content-mismatched speech data. *Proceedings of the 18th International Congress of*
874 *Phonetic Sciences. (Glasgow, UK)*. Paper 458.
875

876 Brown, G. (2016). Automatic accent recognition systems and the effects of data on performance.
877 *Proceedings of Odyssey: the speaker and language recognition workshop. (Bilbao, Spain)*. Paper 29.
878

879 Brown, G. (2018). Segmental content effects on text-dependent automatic accent recognition.
880 *Proceedings of Odyssey: the speaker and language recognition workshop. (Les Sables d'Olonne,*
881 *France)*. 9-15.
882

883 Brown, G., & Wormald, J. (2017). Automatic Sociophonetics: Exploring corpora with a forensic
884 accent recognition system. *Journal of the Acoustical Society of America*. 142. 422-433.
885

886 Contini, M. and Romano, A. (2002) Atlas Multimédia Prosodique de l'Espace Roman. URL
887 <http://dialecto.u-grenoble3.fr/AMPER/new.htm>.
888

889 D'Arcy, S., Russell, M., Browning, S. and Tomlinson, M. (2004). The Accents of the British Isles
890 (ABI) corpus. *Proceedings of Modélisations pour l'Identification des Langues. (Paris, France)*. 115-
891 119.
892

893 Davis, S., & Mermelstein, P. (1980). Comparison of parametric representations for monosyllabic
894 Word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech and*
895 *Signal Processing*. 28.
896

897 Elvira-García, W., Balocco, S., Roseano, P., & Fernández-Planas, A. M. J. S. C. (2018). ProDis: A
898 dialectometric tool for acoustic prosodic data. *Speech Communication*, 97, 9-18.
899

900 Ferragne, E., & Pellegrino, F. (2010). Vowel systems and accent similarity in the British Isles:
901 Exploiting multidimensional acoustic distances in phonetics. *Journal of Phonetics*. 38. 526-539.
902

903 Grabe, E., Kochanski, G., & Coleman, J. (2007). Connecting intonation labels to mathematical
904 descriptions of fundamental frequency. *Language and Speech*, 50(3), 281-310.
905

906 Hanani, A., Russell, M., & Carey, M. (2013). Human and computer recognition of regional accents
907 and ethnic groups from British English speech. *Computer Speech and Language*. 27. 59-74.
908

909 Hellmuth, S. (2007). The relationship between prosodic structure and pitch accent distribution:
910 Evidence from Egyptian Arabic. *The Linguistic Review*, 24(2-3), 289-314.
911

912 Hellmuth, S. (2018). *Variation in polar interrogative contours within and between Arabic dialects.*
913 *Proceedings of the 9th International Conference on Speech Prosody. (Poznań, Poland)*. 989-993.
914

915 Hellmuth, S. (2020). Contact and variation in Arabic intonation. In C. Lucas & S. Manfredi (Eds.),
916 *Arabic and contact-induced change: a handbook*. Berlin: Language Science Press.
917

918 Hellmuth, S. (to appear). *Intonation in spoken Arabic dialects*. Oxford: Oxford University Press.
919

920 Hellmuth, S., & Almbark, R. (2019). *Intonational Variation in Arabic Corpus (2011-2017)*.
921 Retrieved from: <http://reshare.ukdataservice.ac.uk/852878/>
922

923 Hualde, J., & Prieto, P. (2016). Towards an International Prosodic Alphabet (IprA). *Laboratory*
924 *Phonology: Journal of the Association for Laboratory Phonology*, 7(1).
925

926 Huckvale, M. (2004). ACCDIST: a metric for comparing speakers' accents. *Proceedings of the*
927 *International Conference on Spoken Language Processing. (Jeju, Korea)*. 29-32.
928

929 Huckvale, M. (2007). ACCDIST: an accent similarity metric for accent recognition and diagnosis. In
930 C. Müller (Ed.) *Lecture Notes in Computer Science: Speaker Classification*. Vol 2. Berlin
931 Heidelberg: Springer-Verlag. 258-274.
932

933 Krifka, M. (2008). Basic notions of information structure. *Acta Linguistica Hungarica*, 55(3-4), 243-
934 276.
935

936 Najafian, M., Khurana, S., Shon, S., Ali, A., & Glass, J. (2018). Exploiting convolutional neural
937 networks for phonotactic based dialect identification. *Proceedings of the IEEE International*
938 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*. (Calgary, Canada.) 5174-5178.
939

940 Owens, J. (2013). A house of sound structure, of marvelous form and proportion: An introduction. In
941 J. Owens (Ed.), *The Oxford Handbook of Arabic Linguistics*. (pp. 1-22). Oxford: Oxford University
942 Press.
943

944 Retsö, J. (2013). What is Arabic? In J. Owens (Ed.), *The Oxford Handbook of Arabic Linguistics* (pp.
945 433-450). Oxford: Oxford University Press.
946

947 Shon, S., Ali, A., & Glass, J. (2018). Convolutional neural network and language embeddings for
948 end-to-end dialect recognition. In *Proceedings of Odyssey: the speaker and language recognition*
949 *workshop*. Les Sables d'Olonne, France. 98-104.
950

951 Snyder, D., Ghahremani, P., Povey, D., Garcia-Romero, D., & Carmiel, Y. (2017). Deep neural
952 network embeddings for text-independent speaker verification. *Proceedings of Interspeech*.
953 (Stockholm, Sweden). 165-170.
954

955 Vapnik, V. (1998). *Statistical Learning Theory*. New York: Wiley.
956

957 Versteegh, K. (2014). *The Arabic Language*. Edinburgh University Press.
958

959 Walker, T. (2014). Form (does not equal) function: The independence of prosody and action.
960 *Research on Language and Social Interaction*, 47(1), 1-16.
961

962 Wormald, J. (2016). *Regional Variation in Punjabi-English*. PhD thesis. University of York, UK.
963

964 Young, S., Evermann, G., Gales, M., Hain, T., Kershaw, D., Liu, X., Moore, G., Odell, J., Ollason,
965 D., Povey, D., Valtchev, V., & Woodland, P. (2009). *The HTK Book for HTK Version 3.4*.
966 Cambridge University Engineering Department, Cambridge.
967

968 Zissman, M. (1996). Comparison of four approaches to automatic language identification of
969 telephone speech. *IEEE Transactions on Speech and Audio Processing*. 4. 31-44.