## Deep learning assessment of syllable affiliation of intervocalic consonants

Zirui Liu,[1,a] Yi Xu[1, b]

*[1] Speech, Hearing and Phonetic Sciences, University College London,*

*London, WC1N 1PJ, United Kingdom*

In English, a sentence like "He made out our intentions." could be misperceived as "He may doubt our intentions." because the coda /d/ sounds like it has become the onset of the next syllable. The nature and the occurrence condition of this resyllabification phenomenon are unclear, however.  Previous empirical studies mainly relied on listener judgment, limited acoustic evidence such as voice onset time (VOT) or average formant values to determine the occurrence of resyllabification. This study tested the hypothesis that resyllabification is a coarticulatory re-organisation that realigns the coda consonant with the vowel of the next syllable. We used deep learning in conjunction with dynamic time warping (DTW) to assess syllable affiliation of intervocalic consonants. The results suggest that convolutional and recurrent neural network (CNN-RNN) based models can detect cases of resyllabification using Mel-frequency spectrograms. DTW analysis shows that neural network inferred resyllabified sequences are acoustically more similar to their onset counterparts than their canonical productions. A binary classifier further

[a] zirui.liu.17@ucl.ac.uk
[b] yi.xu@ucl.ac.uk

1

24  suggests that similar to the genuine onsets, the inferred resyllabified coda

25  consonants are coarticulated with the following vowel. These results are

26  interpreted with an account of resyllabification as a speech-rate-dependent

27  coarticulatory reorganisation mechanism in speech.

## I.    INTRODUCTION

Despite the wide recognition of the syllable as a speech unit among speakers and researchers (Browman & Goldstein, 1992; Levelt, Roelofs & Meyer, 1999; MacNeilage, 1998), there have been doubts about the role of the syllable due to ambiguity associated with syllable boundaries. One situation where ambiguity is especially severe is in regard to the syllable affiliation of intervocalic consonants. For example, the phrase "escort us" in British English (/ɛs#kːɔt#əs/) can be syllabified as /ɛs#kːɔ#təs / in connected speech, according to observation of a noisy release during the word final /t/ (Levelt et al., 1999). The phenomenon is more formally known as resyllabification, which usually denotes a shift of syllabification of a coda consonant into the onset of the following vowel-initial syllable (Levelt et al., 1999; Schiller et al., 1997). For English, empirical work examining resyllabification goes back as early as 70 years ago, when Stetson used the kymograph to investigate CV and VC production at different speech rates (Stetson, 1951). He observed that in a sequence of syllables like /bi bi bi…/, the CV structure remains stable regardless of speech rate. In contrast, a sequence of VC syllables such as /ib ib ib…/, becomes very similar to /bi bi bi…/ when repeated at a fast rate, according to kymograph data, indicating that the coda /b/ is resyllabified as an onset consonant. The perceptual finding was consistent with articulatory patterns recorded by the kymograph. Stetson's findings were later replicated by Tuller and Kelso (1990, 1991), with glottal transillumination data, which showed that glottal

51 movements shifted drastically at a critical rate of speech, and perception of

52 the spoken sequences also shifted to be mostly identified as /ip ip ip.../.

53 In languages such as Spanish and French (Bermúdez-Otero, 2011, Gaskell

54 et al., 2002), resyllabification is recognised as a phonological process,

55 although there are cross dialect variations according to acoustic evidence

56 such as consonantal duration (Strycharczuk & Kohlberger, 2016). Due to

57 the lack of clear empirical evidence, the existence of resyllabification in

58 English is questioned (Shattuck-Hufnagel, 2011), as mentioned above.

59 Furthermore, the status of the syllable is called into question because of

60 boundary ambiguity due to resyllabification (Blevins, 2003; Steriade, 1999).

61 A major source of the difficulty of determining the syllabification status of

62 segments is that it is mainly based on the subjective judgment of listeners

63 (Ní Chiosáin et al., 2012; Content, 2001; Goslin & Frauenfelder, 2001;

64 Schiller et al., 1997). Even when acoustic measurements are taken, listener

65 judgments are still treated as the "ground truth" (de Jong et al., 2004;

66 Mullooly, 2003). But as found in de Jong et al. (2004), listeners agree with

67 each other well only in cases where a gap between the release of the coda

68 consonant and the beginning of voicing for the next vowel can be easily

69 detected. In the absence of apparent gaps, listener judgments become very

70 diverse. Those authors therefore suggested that the difference between the

71 coda and onset consonant is more closely related to how they are

72 *motorically optimised* in production in ways that are too subtle for most

73 listeners to detect.

3

74  What is needed is an alternative definition of resyllabification, that departs

75  from conventional definitions that are based on language-specific

76  phonotactics (what is phonologically legal), perceptual impression, and

77  language-specific acoustic properties (aspiration, voicing, etc.). In this

78  study, we consider an articulatory-acoustic definition that specifies the

79  affiliation of an intervocalic consonant based on an articulatory definition of

80  the syllable. And the definition of the syllable, as will be reviewed next, also

81  addresses coarticulation, another essential issue of speech articulation.

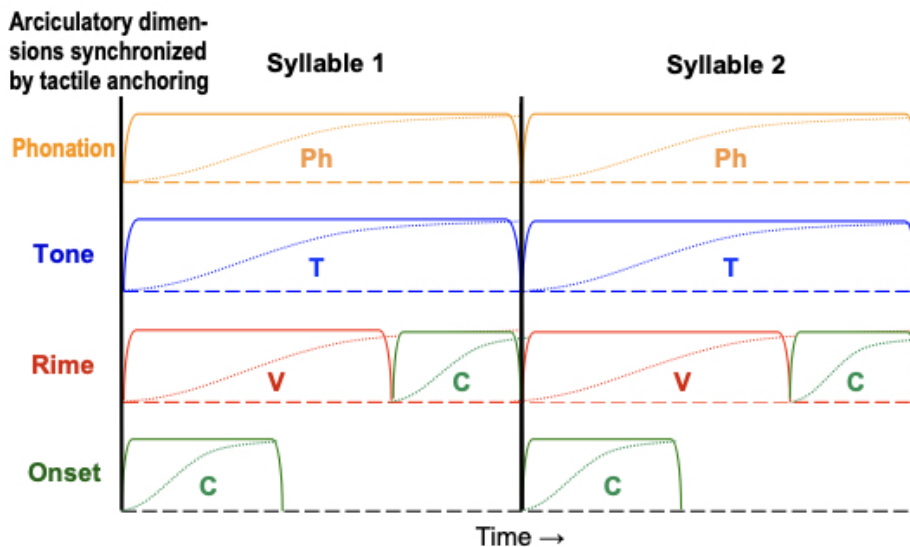82  **A. Resyllabification, coarticulation and the syllable**

83  Resyllabification is closely related to a well-documented asymmetry

84  between onset and coda consonants in both phonology and phonetics. For

85  languages that allow for coda consonants, codas are more vulnerable than

86  their onset counterparts, as they are more susceptible to deletion and

87  reduction (Barlow & Gierut, 1999; Xu, 1986, 2020). In contrast, onset

88  consonants are often inserted when the syllable is vowel initial, such as

89  glottal stop insertion (Birgit, 2001; Garellek, 2012), intrusive /r/s (Gick,

90  1999; Uffmann, 2007), and vowel hiatus breakers (Mudzingwa, 2013; Smith,

91  2001). In terms of canonical syllable structures, CV syllables are also more

92  common than both VC and CVC syllables in many languages (Clements &

93  Keyser, 1983; Levelt et al., 1999; Xu, 2020).

94  According to articulatory phonology, the vulnerability of codas is likely

95  related to an asymmetry in coarticulation within the syllable. That is, onset

4

96  consonants are coupled "in-phase" with the vowel, resulting in synchronous

97  activation between the vocalic and onset C gestures (Goldstein et al., 2006).

98  On the other hand, coda consonants are coupled "anti-phase" with the

99  vowel, which is a less stable mode of coordination. Resyllabification is

100 therefore "analysed as an abrupt transition to a more stable coordination

101 mode" that is likely to occur under increased speaking rate (Goldstein et al.,

102 2006:237).

103 An alternative account of resyllabification is provided by the

104 synchronisation model of the syllable (Xu, 2020), as shown in Fig. 1, which

105 shares some similarities with articulatory phonology but differs from it in

106 certain critical details. The model assumes that syllable is a mechanism for

107 eliminating most of the temporal degrees of freedom by synchronising

108 consonant, vowel and glottal movements at syllable onset (vertical lines),

109 whereby each movement (dotted lines) is to approach an underlying target

110 within its allocated time interval. The synchronisation makes the initial

111 consonant fully overlapped, hence coarticulated, with the initial portion of

112 the "following" vowel. In contrast, a coda consonant is articulated

113 sequentially after the vowel, because its closing movement directly conflicts

114 with the opening movement of the vowel (Xu & Liu, 2006). There are two

115 differences between this model and articulatory phonology that are directly

116 relevant for the current study. First, synchronisation is assumed to be a

117 fundamental design of the syllable (likely centrally controlled) rather than

118 emerging from the coupling of the gestural planning oscillators as in

5

119  articulatory phonology (Goldstein et al., 2006). Second, the sequential

120  articulation of coda consonant is not modelled in terms of phase relation

121  between C and V, because a) individual target approximation movements

122  are frequently allocated insufficient amount of times (Nakatani et al., 1981;

123  Xu & Wang, 2009), thus disallowing them to from complete movement

124  cycles (Xu & Prom-on, 2019), and b) syllables constantly vary their duration,

125  due to stress, phrasing and other linguistic factor, which makes it difficult

126  for syllable sequences, together with their constituent segments, to be

127  temporally periodic to make oscillation-based modelling possible.
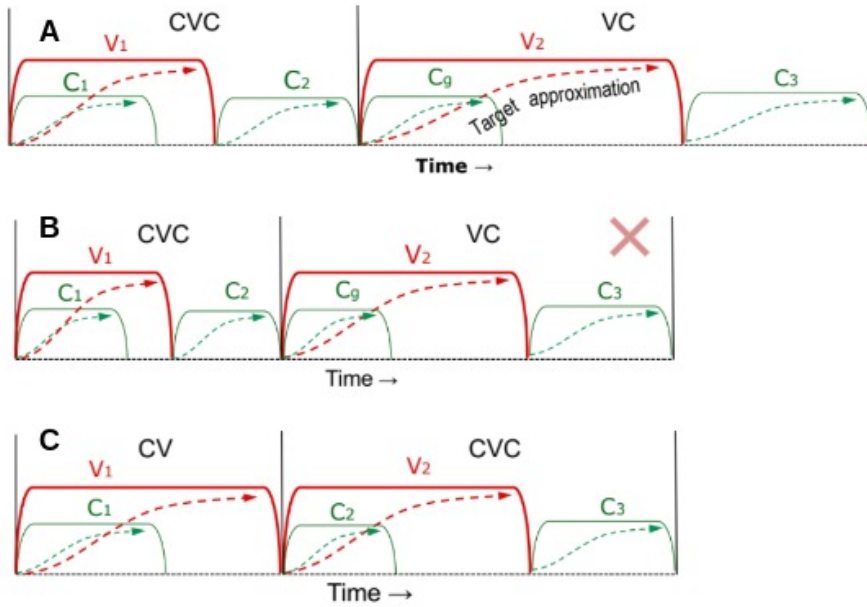
128



129  FIG. 1. The synchronisation model of the syllable (Xu, 2020).

130  According to the synchronisation model of the syllable, resyllabification is

131  due to a lack of articulation time, as schematised in Fig. 2, rather than due

132  to transition from anti-phase to in-phase articulatory coordination. In Fig.

133  2A, the coda consonant ($C_2$) occupies its own time interval because it is

6

134 sequentially articulated after the first vowel ($V_1$). Meanwhile, the second

135 syllable is not articulated as a true VC because it actually starts with a

136 glottal stop ($C_g$). Such glottal stops have been reported as frequently

137 occurring at slow speech rate (Birgit, 2001; de Jong, 2001), but would

138 disappear as speech rate reached a certain threshold, leading to a

139 perceptual shift from /VC#VC/ to /CV#CV/ (de Jong, 2001). As illustrated in

140 Fig. 2B, as speech rate increases, less time is allocated to the syllable,

141 which would require the duration for both $V_1$ and $C_2$ to be shortened to an

142 implausible extent (as indicated by the red cross). The increased time

143 pressure (Tiffany, 1980; Xu & Prom-on, 2019) may then lead to the

144 replacement of the glottal stop ($C_g$) with $C_2$ when speech rate approaches a

145 certain threshold (e.g. 350 ms per syllable (de Jong, 2001)). $C_2$ now

146 becomes the initial consonant of the second syllable, as shown in Fig. 2C.

147 This reorganisation gives $V_1$ more articulation time while preserving all the

148 segmental composition of the original syllables.

7

FIG. 2. Illustration of articulatory resyllabification based on the synchronisation model of the syllable.

Based on this account of resyllabification, two predictions can be made: 1) Due to similarity in articulatory structure, resyllabified codas spectrally resemble their onset counterparts more than their canonical form, and the opposite can be observed for the non NN-resyllabified ones. 2) Because a resyllabified coda is fully coarticulated with the vowel of the following syllable, there is similar amount of vowel information shared between the resyllabified onsets and the canonical onsets, but not between canonical codas and canonical onsets. These predictions can be tested on English by applying machine learning models on acoustic data.

8

**161** **B. Using deep neural networks with acoustic data to identify**

**162** **resyllabification**

**163** Given the difficulty of subjectively judging the occurrence of reyllabification

**164** (de Jong et al., 2004), an alternative is to obtain objective evidence by

**165** taking advantages of recent development in machine learning technology.

**166** This study therefore aims to determine the occurrences of resyllabification

**167** using deep learning models and dynamic time warping in combination with

**168** continuous acoustic data. The deep learning models used were inspired by

**169** state-of-the-art automatic speech recognition (ASR) networks (Amodei et al.,

**170** 2015). ASR systems without language models are error prone when

**171** detecting the canonical structure of resyllabified sequences (Adda-Decker et

**172** al., 2002; Mirzaei et al., 2018; Wu et al., 1997). For example, a sequence

**173** like "fade out" could be recognised as "Fay doubt" if the coda /d/ is

**174** resyllabified as the onset of the second syllable. We trained recognition

**175** networks on slow speech data with no resyllabification occurrences and

**176** used them to classify data from normal rate speech. The reason behind

**177** using data from the slow speech rate condition for training is to ensure that

**178** there are no resyllabified sequences in the training data. In other words, for

**179** the model to be able to misclassify a sequence as its onset counterpart due

**180** to resyllabification, it should not be trained with a resyllabified sequence

**181** labelled as its canonical version. The misclassified sequences in normal

**182** speech rate (i.e. "fade out" as "fay doubt") were further examined to shed

**183** some light on the articulatory structure of the syllable.

9

## II.  Methods

We trained a deep neural network classifier to identify word sequences such as "coo part" and "coop art". The utterances in the slow condition were used for training the classifiers. Then, we used the trained classifiers to classify the same utterances spoken in the normal rate recordings. A /CVC#VC/ sequence such as 'coop art' was categorised as resyllabified if the classifier "misclassified" it as its counterpart /CV#CVC/ sequence, i.e. 'coo part'. These neural network inferred resyllabified sequences are referred to as NN-resyllabified to avoid confusion between the cognitive process of syllable reorganisation and the inferred syllabification status by the classifier. Dynamic time warping was then used to investigate the spectral similarities between the NN-resyllabified sequences in the normal speaking rate and the sequences in the slow rate (e.g. NN-resyllabified "coop art" vs. slow "coo part" or NN-resyllabified "coop art" vs. non resyllabified slow "coop art"). Furthermore, to test prediction (2), we built binary neural network classifiers to categorise contrastive pairs such as "coop art" vs. "coop eat", whose training data only consisted of the intervocalic consonantal portions of the acoustic signal (e.g. aspiration for /p/). The closure interval was not included due to very little acoustic energy in the data, as /p/ is a voiceless stop. The results were compared between speech rates and syllable structures.

10

## A. Subjects

Eight subjects aged 20-40 participated in this study, whose first language was Southern Standard British English (6 female and 2 males). No speaking or hearing disorders were reported prior to recording. To ensure data quality, all potential participants had to submit a short recording on Gorilla. The experimenters then visually inspected the recordings in the computer program Praat (Boersma & Weenink, 2022). Only participants with an external microphone and sufficient recording quality took part in the study.

## B. Stimuli and data collection

Table I lists the word sequences used in this study. The stimuli include three groups of four sequences. For each group, the onset pair and coda pair match in terms of segments and differ in syllable structure, e.g. /CVC#VC/ vs. /CV#CVC/. This maximises the possibility that if the classifier misclassified a coda sequence as its onset counterpart, it is likely due to the shift in syllable structure, i.e. resyllabification.

TABLE I. Stimuli.

| Group | Onset | | Coda | |
|---|---|---|---|---|
| 1 | Lee steal | Lee stale | Least eel | Least ale |
| 2 | Do mart | Do meet | Doom art | Doom eat |
| 3 | Coo part | Coo Pete | Coop art | Coop eat |

11

221  Note that there exist differences other than syllabification between onset

222  and coda sequences, such as lexical, syntactic or prosodic properties. For

223  example, "doom art" is a noun/verb noun sequence, where as "do mart" is a

224  verb noun sequence. The neural network classifier could use information

225  such as syllabification, syntactic and lexical differences between the onset

226  and coda tokens. Therefore, it is important to minimise the *similarities*

227  between items such as "coo part" and "coop art" due to the following: If the

228  classifier misclassified "coop art" as "coo part", it is important to minimise

229  the possibility that the misclassification took place due to prosodic or lexical

230  similarity between the two, rather than coarticulation between the

231  intervocalic C and the second V. Therefore, within each onset and coda pair,

232  we use word combinations that differ in their morphosyntactic structure

233  (e.g. "Lee steal" vs. "least eel"). However, other unknown factors may still

234  result in similarities between the onset and coda pairs which could

235  contribute to misclassification. The current design can only assume that

236  when a coda sequence is misclassified as its onset counterpart, it is due to

237  similarity in coarticulation structure rather than other unknown factors.
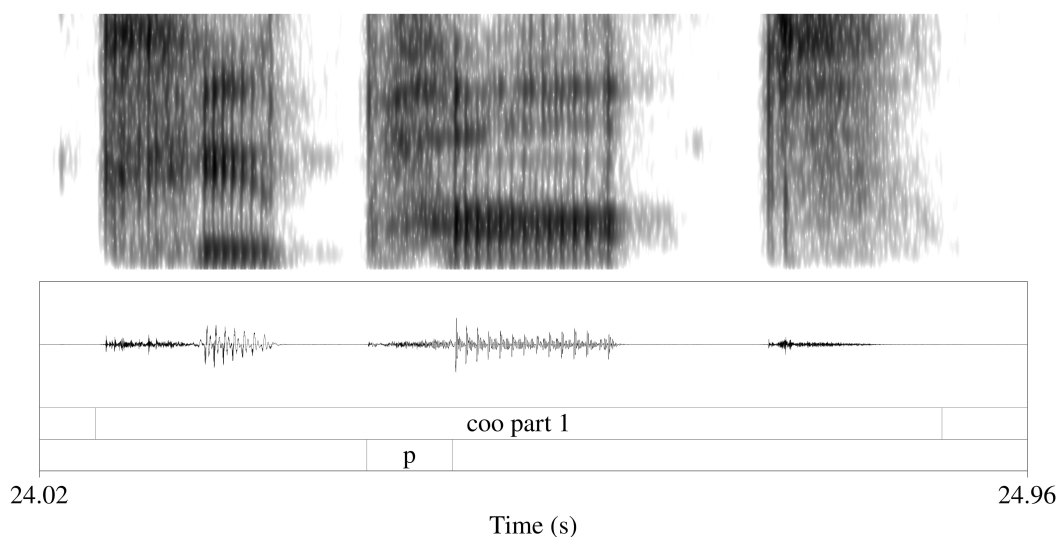
238  There is also a vowel minimal contrast in the second syllable for each

239  syllable structure condition in each group. The vowel contrast allows us to

240  examine the amount of coarticulation in the intervening consonant, by

241  assessing the performance of a binary classifier at predicting the second

242  vowel identity using only acoustic data from the annotated consonant

243  interval. Previous studies have used a minimal pair design and showed that

12

244 when a consonant is coarticulated with the upcoming vowel, acoustic

245 information associated with the vowel can be detected during the consonant

246 (Liu & Xu, 2021, Liu et al., 2022). Liu and Xu (2021) also show that the

247 entire cluster in /clusterV/ syllables in British English is coarticulated with

248 the vowel. Thus, a cluster triplet is included in the current study to

249 investigate whether the following vowel is coarticulated from the onset of

250 the consonant cluster.

251 Participants were instructed to say the word sequences in isolation in two

252 blocks of different speaking rates – first slow, then normal. For the slow

253 block, the speakers were instructed to articulate the words clearly and

254 fluently, at a slow pace. In the normal condition, speakers were informed to

255 speak at a faster pace in a colloquial style. There were no instructions on

256 what resyllabification was, or whether they should or should not resyllabify

257 anything. The stimuli were read aloud with 20 and 10 repetitions for the

258 randomised slow and normal blocks, respectively, yielding 360 tokens per

259 speaker ($12 \times 20 + 12 \times 10$). Around 3% of the data were excluded due to

260 background noise during recording.

261 The recording took place online over Zoom with the sampling rate of 32

262 kHz, with Zoom's original sound feature turned on, which preserved the

263 original recording quality by minimising the amount of audio enhancement.

264 All the participants used an external microphone during the experiment and

265 the recording quality was assessed by the researcher prior to the

266 experiment. For the resyllabification classifiers, the recordings were

267 annotated in either $[C_1V_1\#C_2V_2C_2]$ or $[C_1V_1C_1\#V_2C_2]$ format (subscripts

268 denote syllable position), with the first boundary being the start of acoustic

269 landmark of onset $C_1$ (e.g. lateral murmur for /l/), and the second boundary

270 being the end of acoustic landmark of the coda $C_2$. For the binary classifiers,

271 the consonantal intervals were segmented as the plosive aspiration for /p/,

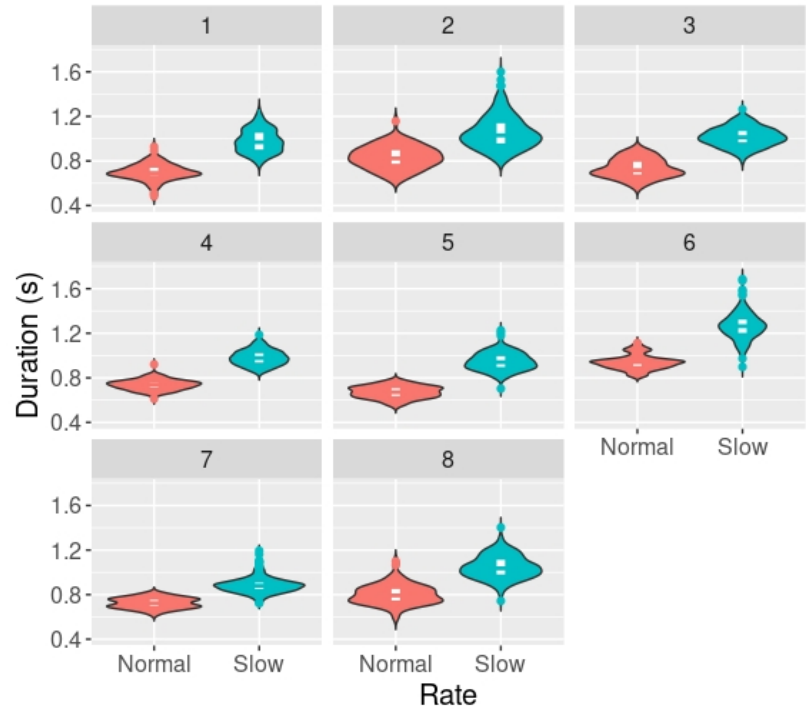272 nasal murmur for /m/ and frication for /s/. An example is shown in Fig. 3.



coo part 1

p

24.02      24.96

Time (s)

273

274 FIG. 3. An annotation example of "coo part" from one speaker, with the

275 vertical lines indicating the segmentations.

276 **C. Speech rate analysis**

277 As speech tempo can be speaker-specific due to difference in speaker

278 characteristic (Jacewicz et al., 2009), participants were free to speak at a

279 rate they deemed appropriate as slow or normal. For both the slow and

14

280 normal rate condition, participants were instructed to speak fluently (i.e.

281 without spontaneous pausing). No spontaneous pauses were identified in

282 the data during the annotation process. Therefore, speech rate in the

283 present study is analogous to articulation rate, which does not include

284 hesitation, pausing or emotional expressions. The duration values of

285 annotated tokens are presented in Fig. 4. As the figure shows, speech rate

286 was faster for the normal condition compared to the slow condition for all

287 speakers. On average, speakers produced 2.9 syllables per second for the

288 normal rate and 2 syllables per second for the slow rate. According to de

289 Jong (2001), resyllabification should take place when articulation rate

290 approaches 2.8 syllables per second.

291



292 FIG. 4. Annotated sequence duration for 8 speakers.

15

### D. Neural network classifier for identifying resyllabification
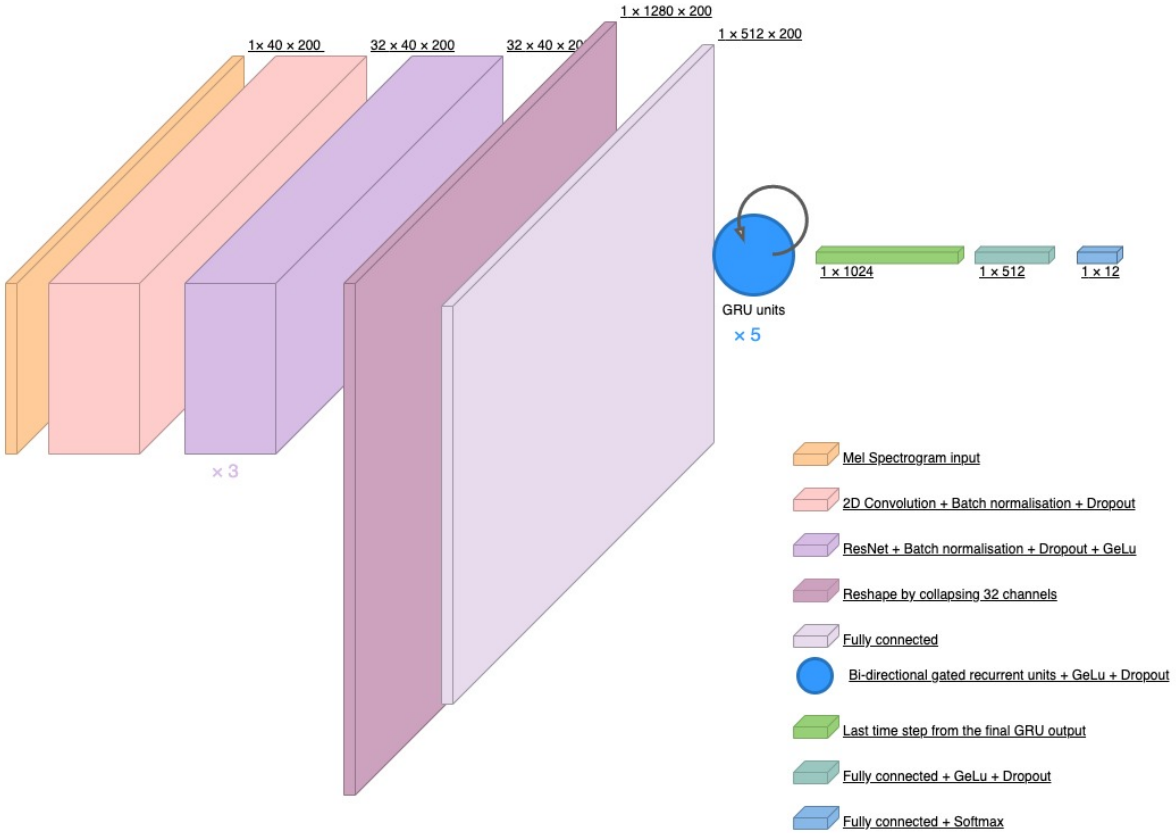
#### *1. Data preparation*

To ensure high accuracy, neural networks were trained for each speaker individually. The segmented word sequences from the slow condition were converted into mel-frequency spectrograms with 40 mel filter-banks with 25 ms as the window length and a hopping interval of 5 ms. We augmented the data to boost the amount of training data by using common augmentation techniques, such as speed augmentation, noise addition and frequency/time masking (Ko et al., 2015; Park et al., 2019). First, half of the tokens from the speaker were selected and sped up randomly between the factor of 0.3 to 0.9, by using the Audacity software with a custom Python script (Audacity Team, 2021). This resulted in 360 samples per speaker. Then, 15% of the resultant dataset were reserved as the testing set (N = 54), and 85% as the training set (N = 306)[1]. Note that the samples were randomised before data splitting. Since the original data are balanced between word classes, the train and test split should also contain approximately balanced data, resultant of the random sampling process. The training set was then further boosted by augmenting 30% with random Gaussian noise addition to the raw acoustic signal (Pervaiz et al., 2020), or frequency or time masking to the spectrograms (Park et al., 2019), yielding 398 samples for the training set. Not only does data augmentation improve model generalisation and performance, the sped-up samples also familiarise the model with shorter acoustic signal such as those in the normal speech rate condition. The

16

316  motivation for doing noise addition and masking boost after the speed boost

317  is to provide the benefit of these augmentation techniques for the sped-up

318  tokens as well rather than just the original slow sequences.

319  ### *2. Model architecture*

320  The model architecture is shown in Fig. 5[2], which was inspired by a

321  combination of Deep Speech and ResNet, developed by Baidu (Amodei et

322  al., 2015) and Microsoft (He et al., 2015), respectively. Each model was

323  trained for 120 epochs, unless the average accuracy across the last 5

324  epochs has reached the threshold of 98% for the testing set. For each

325  epoch, the spectrograms were padded to the same duration as the longest

326  sequence in the batch (N = 32), then fed into the neural network. Note that

327  Fig. 5 demonstrates the flow of data through the network by a batch size of

328  1. The spectrogram is first passed through a 2D convolutional layer (i.e.

329  convolutional neural network (CNN)), which had a 3 × 3 kernel with a

330  stride of 1, and 32 channels. The output from the 2D convolutional layer is

331  then passed through 3 residual blocks (He et al., 2015), the convolutional

332  layers in each residual block had a 5 × 5 kernel with a stride of 1. For both

333  the 2D convolutional and residual layers, padding was used to retain the

334  shape of the tensors. The motivation behind these two types of

335  convolutional layers is for the model to extract features such as dynamic

336  information of spectral energy between frequencies or time steps (e.g.

337  velocity of energy variation between time steps) (Luo et al., 2018; Sharma

338  et al., 2020). To preserve as much acoustic information as possible, no

17

pooling was used. The output from the residual layers was reshaped by collapsing the 32 channels, resulting in tensors with the shape of 1280 by n timesteps, which was further reduced by a fully connected layer with 512 units. Five layers of bi-directional Gated Recurrent Units (GRU) were then used to process the sequential acoustic features. Only the last timestep's output was used from the GRU. Finally, the output was fed into two fully connected layers with a final SoftMax activation which generated the 12-dimensional probability vector, one for each word sequence in Table I. Due to the complexity of the model, we used dropout as the regularisation technique to combat overfitting (Semeniuta et al., 2016). A dropout rate of 0.1 was used throughout the network (see Fig. 5 for dropout locations). Furthermore, batch normalisation was applied after each mini batch to stabilise learning, as well as provide some regularisation effect (Ioffe & Szegedy, 2015). The hyperparameters were tuned by using grid search with data from the pilot study. The hyperparameters used can be found in Table II in the Appendix section.

1x 40 x 200   32 x 40 x 200   32 x 40 x 20   1 x 1280 x 200   1 x 512 x 200

× 3

GRU units
× 5

1 x 1024    1 x 512    1 x 12

Mel Spectrogram input

2D Convolution + Batch normalisation + Dropout

ResNet + Batch normalisation + Dropout + GeLu

Reshape by collapsing 32 channels

Fully connected

Bi-directional gated recurrent units + GeLu + Dropout

Last time step from the final GRU output

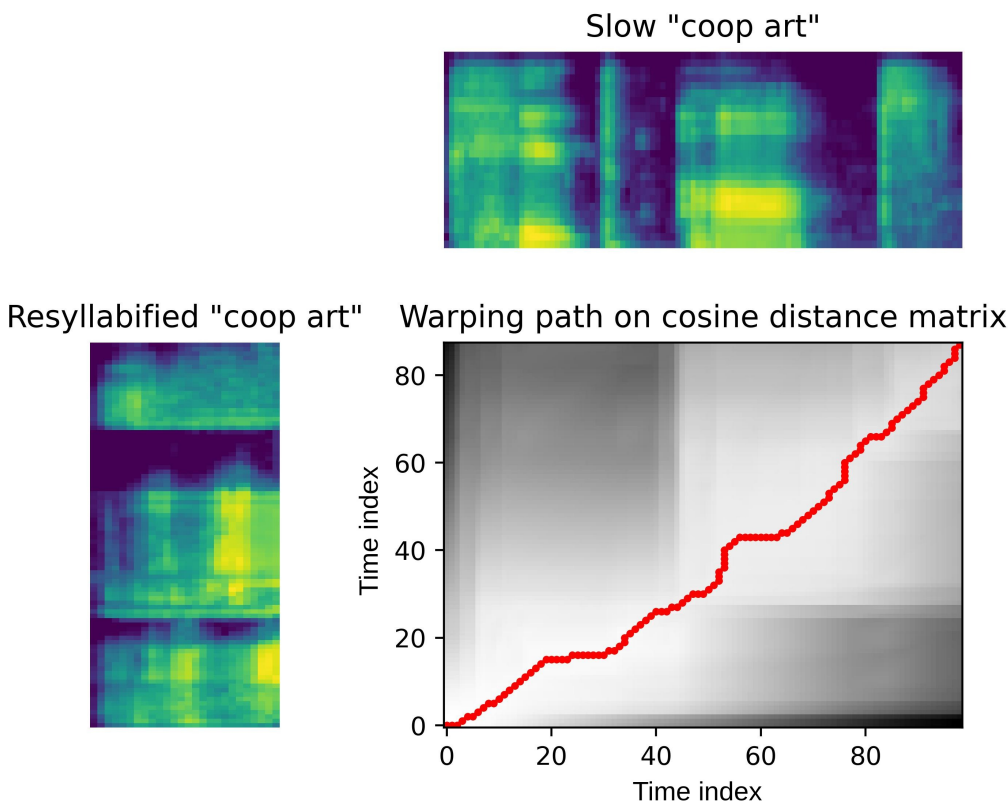Fully connected + GeLu + Dropout

Fully connected + Softmax

355

356  FIG. 5. Model architecture for the resyllabification classifier. The figure

357  shows the tensor dimensions for a batch size of 1. The box sizes reflect

358  tensor shapes, as annotated above each box. The depth, height and width of

359  the boxes are not to scale and is for illustration purposes only.

360  The trained models were used to classify tokens from the normal speech

361  rate condition for each speaker. If a coda sequence was misclassified as its

362  onset counterpart (e.g. "coop art" classified as "coo part"), we categorised it

363  as resyllabified.

19

**E. Dynamic time warping analysis**

Dynamic time warping (DTW) was used to measure how similar the NN-resyllabified and non NN-resyllabified tokens were in relation to the onset or coda conditions in the slow speech rate condition. DTW has been demonstrated to be effective at measuring similarity between sequences such as acoustic signals. For example, it has been widely used for speech recognition (Sakoe & Chiba, 1978; Zhang et al., 2014), as well as other applications such as bird song recognition (Kogan & Margoliash, 1998), speech segment clustering (Lerato & Niesler, 2019), and accent quantification (Bartelds et al., 2020). The DTW algorithm is illustrated in Fig. 6. First, a cost matrix is computed by measuring the distance between the feature vectors (in this case we used mel-spectrograms) between two sequences at each time step. We used cosine similarity for the calculation of distance, as it is not affected by the magnitude of spectral energy, i.e. frequency decibels (e.g. the same recording played at different volumes would measure 0 in cosine distance but not Euclidean distance). The lower right heatmap in Fig. 6 shows the cosine distance between the mel-frequency vectors in the two sequences at all time steps. DTW works by finding the path in the distance matrix that result in the lowest cumulative distance (i.e. cost). Therefore, the DTW distance between the two sequences in Fig. 6 is the sum of the distance values through the warping path shown by the red line.

20

Slow "coop art"

Resyllabified "coop art"    Warping path on cosine distance matrix

386

FIG. 6. Demonstration of the DTW algorithm. The dotted line shows the dynamic warping path. The spectrograms are mel-spectrograms of the tokens "coop art" (bottom left) and "coop art" (top). The pixel intensity in the lower right heatmap represent feature distances at each time step between the two spectrograms.

Using DTW, we can compute the similarity between word sequences, while minimising the effect of speech tempo. For this study, we calculated the distances between the NN-resyllabified as well the the non NN-resyllabified coda sequences and their onset and coda counterparts in the same group from the slow rate condition (e.g. NN-resyllabified "coop art" vs. slow "coo part", "coo Pete" or NN-resyllabified "coop art" vs. slow "coop art" and

21

398 "coop eat"). Note that since the vowel contrast is constant between the

399 distance comparisons, it should not confound the analysis.

400 The DTW analysis was used to compare the similarities between the NN-

401 resyllabified sequences and the onset and coda sequences in the slow

402 condition. In addition, a parallel DTW analysis was conducted for the non

403 NN-resyllabified (correctly classified normal rate coda sequences) to assess

404 whether they are more similar to their canonical form.

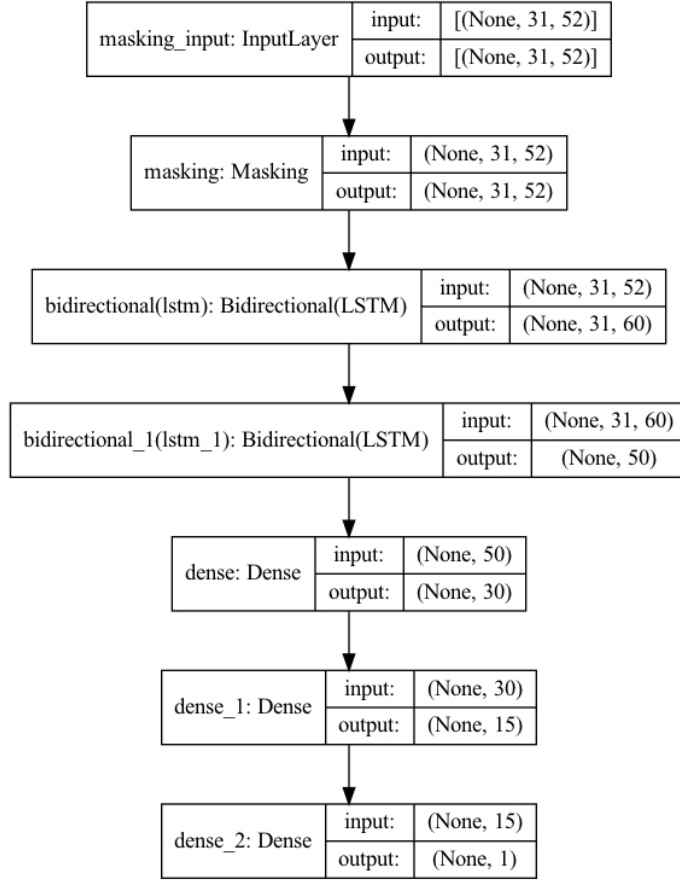405 **F. Detecting $V_2$ information in the intervocalic consonant**

406 As illustrated in Fig. 3, the researcher manually segmented the canonical

407 acoustic intervals from the intervocalic consonant or the first cluster

408 component (i.e. nasal murmur for /m/, aspiration for /p/ and frication for /s/),

409 which were used to investigate the articulatory alignment of the consonant

410 and the following vowel. The segmented intervals differ in terms of

411 articulatory meaning between groups, as aspiration correspond to the

412 consonantal release gesture and the other two correspond to consonantal

413 closures. This difference should have an impact on the amount of vowel

414 information detected in each group. Similar to methods used in Tilsen

415 (2020), Tilsen et al. (2021) and Liu and Xu (2021), to detect vowel

416 information in the segmented intervocalic C, we trained a simple recurrent

417 neural network to predict the second vowel identity between contrastive

418 pairs (e.g. NN-resyllabified "coop art" vs. NN-resyllabified "coop eat"). Liu

419 and Xu (2021) showed that for tautosyllabic $C_nV$, binary classifiers are able

22

420 to detect vowel information in the acoustic intervals of onset C, such as

421 during frication or lateral murmur.

422 For each minimal pair, tokens from all 8 speakers were used. From the

423 normal speech rate condition, only the NN-resyllabified tokens and the true

424 onset tokens were examined. According to results from the neural network

425 classifiers, not all coda tokens were NN-resyllabified, which gave rise to the

426 possibility of accuracy scores from the onset conditions being higher than

427 the NN-resyllabified codas, due to having significantly more training data.

428 For example, a speaker resyllabified 5 out of 10 repetitions of "coop art"

429 and "coop eat", which would result in 10 samples in total for the neural

430 network, whereas 20 samples are available for the onset condition (i.e. 10

431 repetitions of "coo part" and "coo Pete"). Therefore, we balanced the

432 sample sizes between the two conditions by randomly sub-sampling the

433 onset condition for each speaker to match the number of NN-resyllabified

434 ones. For instance, if a speaker resyllabified 5 out of 10 repetitions of "coop

435 eat", only 5 random selections of "coo Pete" were used from this speaker for

436 training the binary classifier.

437 The classifiers were bi-directional recurrent neural networks with Long

438 Short-Term Memory (LSTM) units (Soltau et al., 2016). The network details

439 are shown in Fig. 7. The hyperparameters were tuned with data from the

440 pilot study using grid search, and details can be found in Table III in the

441 Appendix section. The segmented tokens were converted into mel-

23

442 spectrogams with 26 filter banks, with 0.025 s as the window length and

443 0.005 s as the hop length. Before training, all the spectrograms were

444 padded to the same length as the longest one. As Fig. 7 shows, masking

445 was applied in the input layer, which tells the model to ignore the padded

446 duration. Due to the absence of CNN, we included delta coefficients (i.e.

447 first order differentials) to aid model performance, which resulted in a 52

448 dimensional vector at each time step. The data were split into training and

449 testing splits with the ratio of 8:2. We randomly shuffled the data for each

450 minimal pair and trained a model from scratch 80 times and reported the

451 accuracy distribution on the testing sets. The motivation behind examining

452 an accuracy distribution is to avoid the issue of accidental above chance

453 performance, which could arise with small datasets (Combrisson & Jerbi,

454 2015; Ojala & Garriga, 2009).

455

FIG. 7. Model architecture of the binary classifiers. The tensor shapes are denoted on the right of each box.

## 1. Bayesian analysis

To test the amount of vowel information in the acoustic signal, we used Bayesian analysis with beta likelihood to model the effect of syllable structure (i.e. onset vs. coda) on model accuracy. A conventional non-significant result cannot be used to validate a null hypothesis, as it only suggests a failure to reject it. The advantage of using Bayesian statistics is that it simply tells us which model is more supported by the evidence in the data, and the models do not need to be nested. The motivation behind using

25

466  beta regression is due to the nature of accuracy rate being bounded

467  between 0 and 1. Beta regression assumes that the data generating process

468  can be modelled by a beta distribution (Balakrishnan & Nevzorov, 2003),

469  where the distribution can be parameterised with the mean-precision (μ-φ)

470  parameters, where φ is analogous to the inverse of data dispersion. Since Y

471  ~Beta(μ, φ), beta regression presumes that the mean μ of the response

472  given the predictor X is linear on the logit transformed scale (Douma &

473  Weedon, 2019). In other words, in a beta regression model, the dependent

474  variable can be mapped from the bounded space [0, 1] to unbounded real

475  numbers with a link function (most commonly the logit function), where an

476  ordinary linear regression can be used to model the logit transformed data.

477  During Bayesian estimation of the posterior distribution of the model

478  parameters, the likelihood function with the μ-φ parameterisation is:

479
$$f(y;\mu,\phi)=\frac{\Gamma(\phi)}{\Gamma(\mu\phi)\Gamma((1-\mu)\phi)}y^{\mu\phi-1}(1-y)^{(1-\mu)\phi-1} \tag{1}$$

480  and:

481
$$\mu=logit^{-1}(X\beta) \tag{2}$$

482  Γ is the Gamma function, μ is the inverse logit transformed model

483  prediction, y is the observed data bounded between 0 and 1, and φ is the

484  precision parameter. Note that model predictions are mapped back to the

485  bounded space with the inverse logit function. Our accuracy data contains

486  values equal to one. Therefore, the one-inflated beta distribution is needed,

26

487 which produces a mixture density (Ospina & Ferrari, 2012). The likelihood

488 function using the one-inflated beta distribution incorporates a new

489 parameter α:

490
$$f(y;\alpha,\mu,\phi)=\begin{cases}(1-\alpha)f(y;\mu,\phi) & (0<y<1)\\ \alpha & (y=1)\end{cases}$$
(3)

491 To construct beta regression models with Bayesian analysis with the one-

492 inflated beta distribution for the likelihood function, we defined a custom

493 response distribution with the brms package in R[3]. Weakly informative

494 Gaussian priors ($\beta \sim N(0, 5^2)$) were used as the priors for the regression

495 coefficients. The half Cauchy distribution was used for φ ($\varphi \sim$ Cauchy[0,

496 $5^2$)), and the beta distribution for α ($\alpha \sim$ Beta(0.5, 8)). Note that model

497 coefficients do not need to be bounded in any way, as model output is

498 transformed with the inverse logit function into the bounded space.

499 Bayes Factors (BF) were used for model comparison (Dienes, 2016; Liu et

500 al., 2022; Stone, 2013). There is controversy regarding using BF to

501 substitute for null hypothesis testing (Gelman et al., 2013). However, BF is

502 used here to compare which model is more likely given the evidence (i.e.

503 the data), rather than the likelihood of the observed effect being due to

504 chance, as is the case in null hypothesis testing (Morey et al., 2016;

505 Wagenmakers et al., 2016). Other popular methods such as the Bayes leave-

506 one-out (LOO) analysis show limitations when the ground truth is consistent

507 with the null hypothesis. Gronau and Wagenmakers (2019) demonstrates

27

508 that when the number of observations consistent with the simpler model

509 (i.e. $H_0$) grows larger, LOO's support for it reaches an upper bound, and this

510 bound can sometimes be very modest. It was also shown that depending on

511 the prior distribution, as more $H_0$ consistent data is added, LOO's support

512 for $H_0$ can decrease. Therefore, to avoid potential bias towards the more

513 complex model, we use BF for model comparison.

514 If $BF_0$ (the BF indicating evidence for $H_0$ over $H_1$) is between 0 and 1/10, the

515 data strongly supports $H_1$ over $H_0$. Conversely, if $BF_0$ is larger than 10, there

516 is strong evidence for the null hypothesis (Jeffreys, 1961; Biel & Friedrich,

517 2018; Dienes, 2014; Harms & Lakens, 2018; Lakens et al., 2020;

518 Schönbrodt & Wagenmakers, 2018; Lee & Wagenmakers, 2014).

519 For each speech rate condition, a full model was constructed with the main

520 effects of syllable structure (onset vs. coda for the slow rate and onset vs.

521 NN-resyllabified coda for the normal rate) and group. The null model was

522 constructed with group as the only main effect. We also tested whether the

523 effect of syllable structure differed between item groups, by including an

524 interaction term.

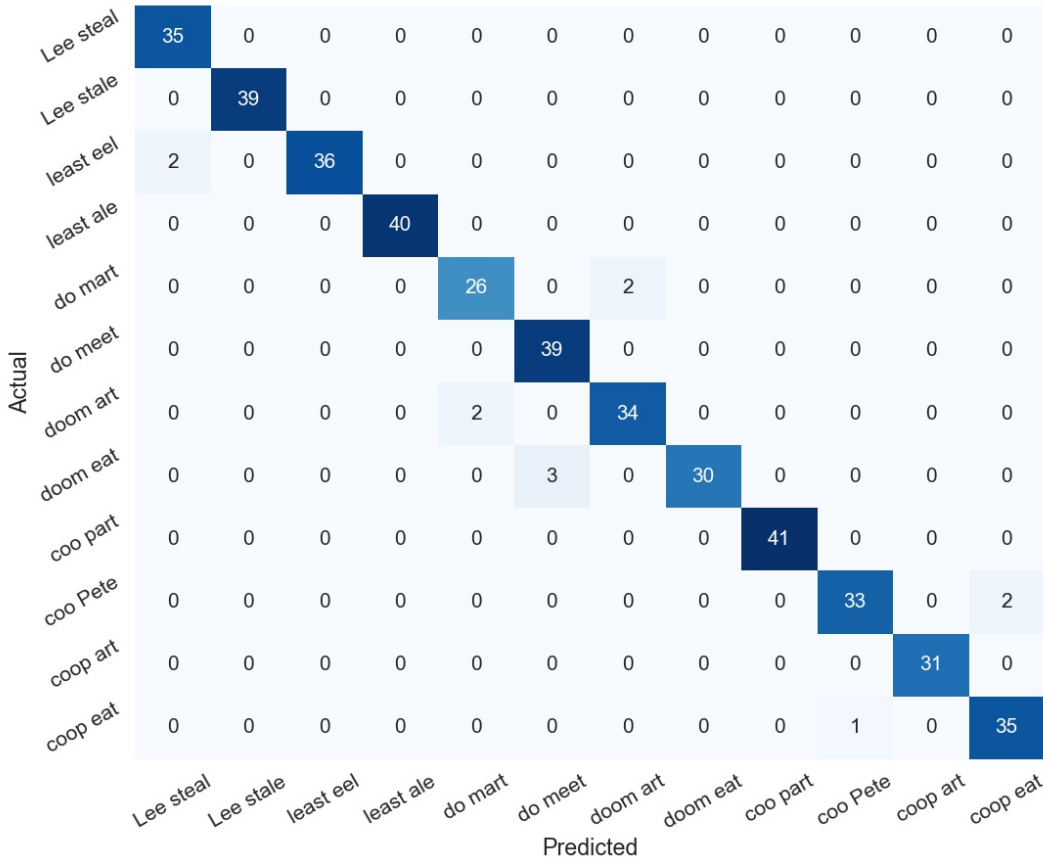525 **G. Duration analysis of NN-resyllabified and canonical onset**

526 **consonants**

527 Although resyllabified sequences may have become similar to their onset

528 counterparts in terms of spectral pattern, there is evidence that resyllabifed

529 codas retain their underlying coda status through duration (Gao & Xu,

28

530  2010; Lehiste, 1960). Specifically, the durations of the resyllabified

531  consonants are shorter compared to the canonical onsets. To test whether

532  duration differs between the two, the same acoustic intervals from the

533  previous section were used. Bayesian analysis with *linear* regression was

534  used to determine if duration of the acoustic interval was affected by

535  syllable affiliation (i.e. genuine onset vs. NN-resyllabified coda). Duration

536  was used as the dependent variable and item group and syllable affiliation

537  were used as the predictor. The likelihood function used the normal

538  Gaussian distribution. For the regression coefficient priors, we used weakly

539  informative Gaussian prior ($\beta \sim N(0, 5^2)$), and for the sigma prior we used

540  the half Cauchy distribution ($\sigma \sim Cauchy[0, 5^2]$).

541  **III.   Results**

542      **A. Resyllabification classifiers**

543  Fig. 8 shows the model performance of the word sequence classifiers. Since

544  we trained a model for each speaker separately, the result in Fig. 8 was

545  calculated by summing over each speaker's confusion matrix. As shown, the

546  classifiers achieved near ceiling accuracy on the test split for the slow

547  speaking rate, indicating that the models could distinguish the word

548  sequences very well.

29
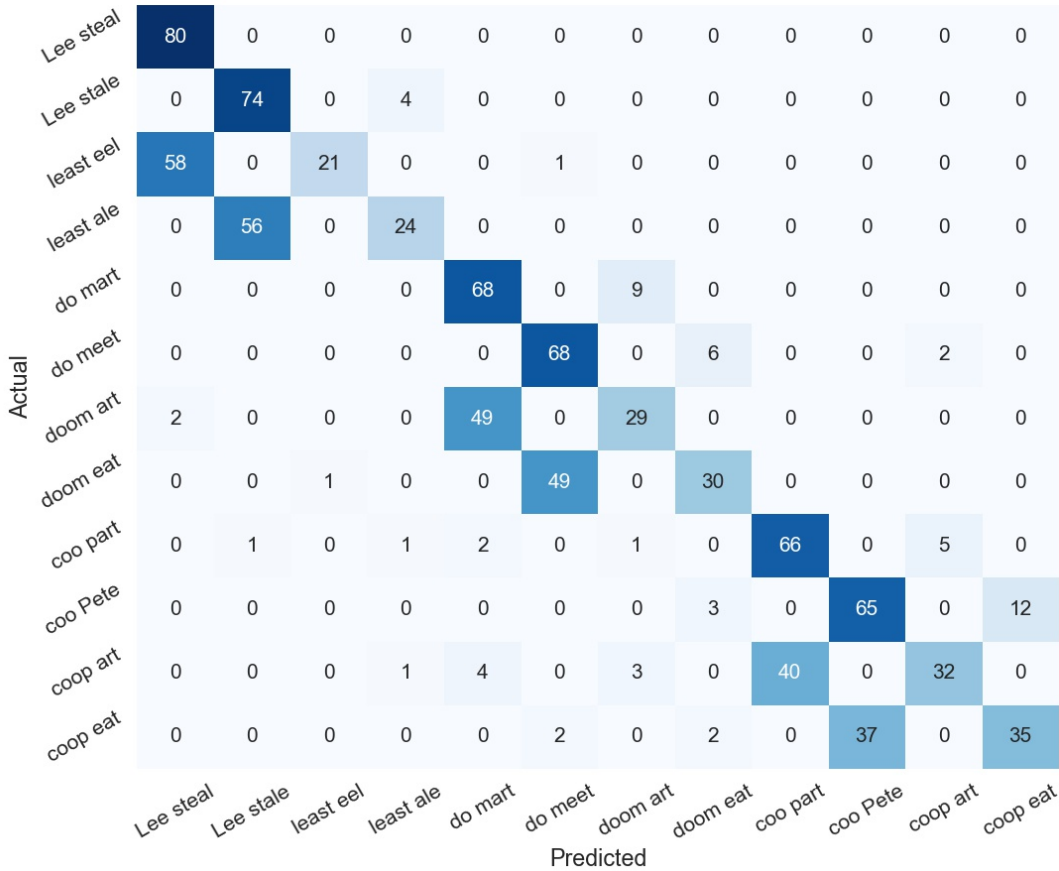
549

FIG. 8. Confusion matrix of model performance on the testing split of the slow speech rate. This is an element wise summation of all the speakers' confusion matrices. The colour intensity of tiles reflects numeric value.

Fig. 9 shows the model performance on the normal speaking rate by summing over the results from all the speakers. Table IV list the accuracy rate for the onset, coda and all sequences. As can be seen, most of the onset sequences were classified correctly. Thus, the classifiers trained on the slow speaking rate data also did well on the onset conditions spoken at a faster rate, such as "Lee steal" or "Lee stale". In the coda condition, the classifiers

30

559 misclassified a large portion of the sequences as their onset counterpart,

560 such as classifying "least eel" as "Lee steal". These misclassified sequences,

561 presumably due to resyllabification, are examined in detail later.

562



563 FIG. 9. Confusion matrix of model performance on the normal speech rate.

564 This is an element wise summation of all the speakers' confusion matrices.

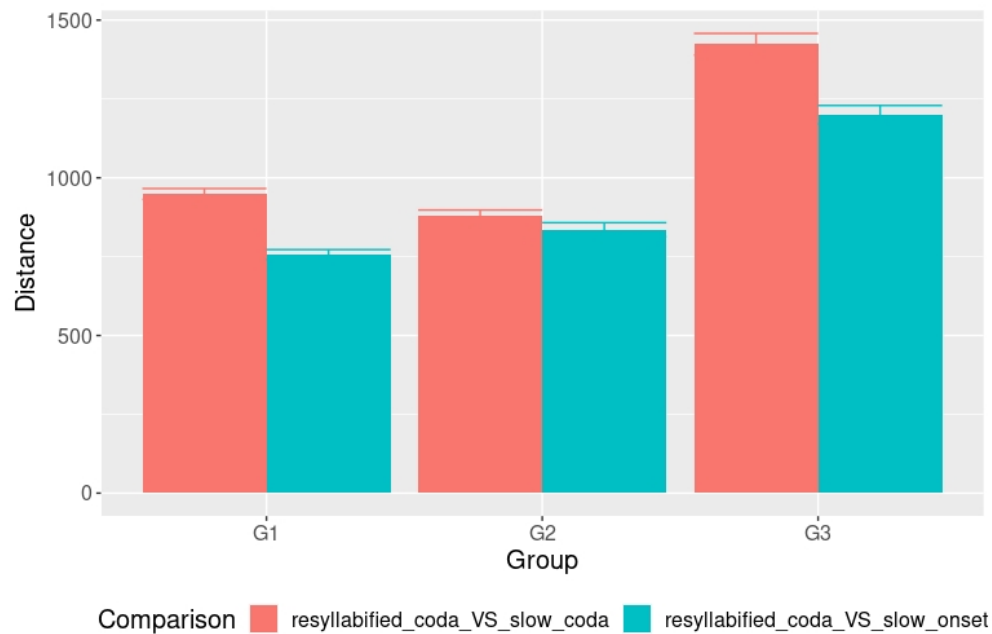565 The colour intensity of tiles reflects numeric value.


566 TABLE IV. Accuracy summary for the normal speech rate tokens.

31

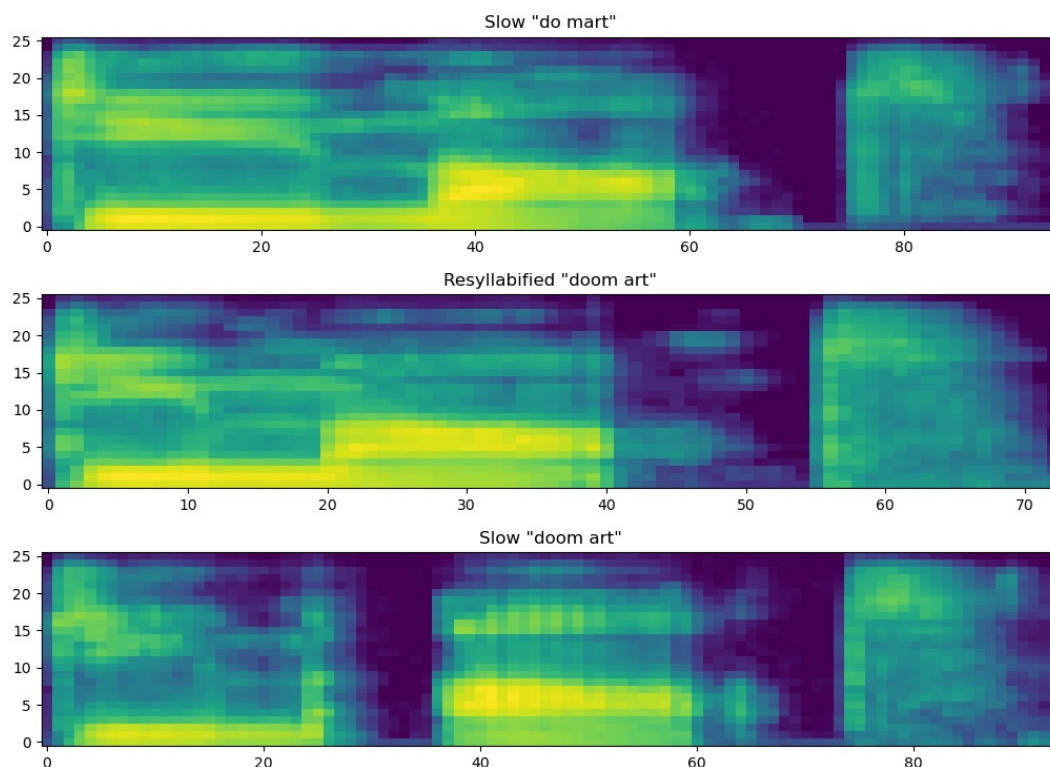| | |
|---|---|
| Coda | 0.36 |
| Onset | 0.90 |
| Overall | 0.63 |

567

**B. DTW analysis**

Fig. 10 shows a bar graph of the cosine distance between the NN-resyllabified tokens and the slow tokens. The NN-resyllabified sequences were only compared to slow sequences in the same group. The figure shows that, when minimising the effect of speech tempo, NN-resyllabified words such as "least eel" is more similar to its canonical onset counterpart "Lee steal" than to its non-resyllabified version. In other words, when comparing the NN-resyllabified condition with the slow onset condition, the cosine distance is smaller than when comparing with the slow true coda condition.

32

FIG. 10. DTW cosine distance between resyllabifed normal rate sequences and slow sequences. The error bars represent 95% of the confidence interval. G1 – "least eel", "least ale", "Lee stale", "Lee steel"; G2 – "doom art", "doom eat", "do mart", "do meet"; G3 – "coop art", "coop eat", "coo part", "coo Pete".
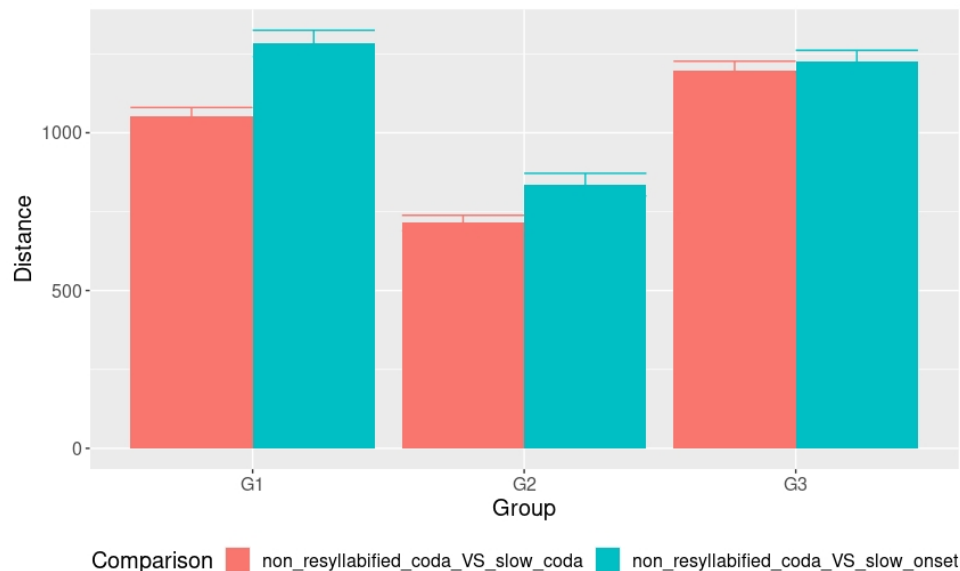
The result from the DTW analysis can be reflected by the spectrograms in Fig. 11. "doom art" in the middle of Fig. 11 was classified as "do mart" by the neural network in the previous section, therefore we treated it as a resyllabified token. The NN-resyllabified "doom art" appears to be more similar to the canonical onset version "do mart" in the top panel. The bottom panel shows "doom art" spoken in the slow condition, likely with a glottal stop at the beginning of the second syllable "art".

33

FIG. 11. Mel-spectrograms of three word sequences from one speaker.

Fig. 12 shows the DTW cosine distance between correctly classified normal rate coda tokens and the slow tokens. The opposite trend from Fig. 10 can be observed: the non NN-resyllabified sequences are more similar to their canonical coda form in the slow rate condition, which support the prediction that correctly classified coda tokens likely have not been resyllabified, unlike their misclassified counterparts.
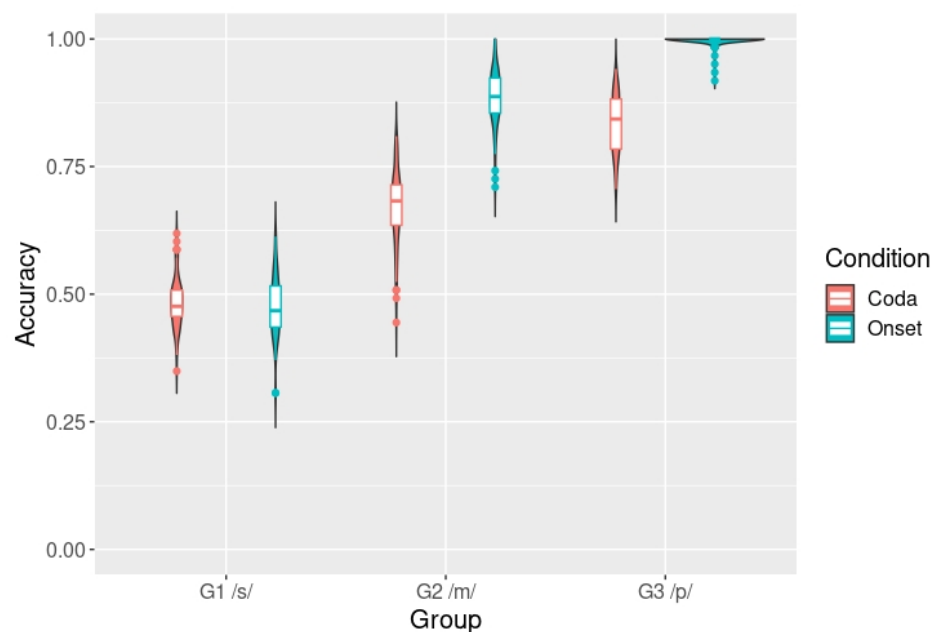
590

591

592

593

594

595

596

597

34

598

FIG. 12. DTW cosine distance between non NN-resyllabifed normal rate
sequences and slow sequences. The error bars represent 95% of the
confidence interval. G1 – "least eel", "least ale", "Lee stale", "Lee steel"; G2
– "doom art", "doom eat", "do mart", "do meet"; G3 – "coop art", "coop eat",
"coo part", "coo Pete".

## C. Intervocalic consonant alignment analysis

### 1. Results for slow speech rate

With the consonant intervals described in Section II E, we trained 80 neural
networks for each vowel minimal pair in Table I and obtained an accuracy
distribution from the test set. Fig. 13 shows the accuracy rate from the slow
speech rate condition. As the figure shows, for /s/ frication in the
intervocalic cluster (i.e. G1), the vowel classification accuracy is around
chance, indicating that little to no vowel information was picked up by the
binary classifier in the frication of /s/ for both the onset (e.g. "Lee stale")

35

613 and coda conditions (e.g. "least ale"). For G2, the intervocalic /m/ contains

614 more detectable vowel information as the onset of the second syllable, and

615 less so when it is the coda of the first syllable. Similar trends can be

616 observed for G3, although with overall higher accuracy – the binary

617 classifier performs better when /p/ is the onset of the second syllable.



618

619 FIG. 13. Vowel classification accuracy by group from the slow speech rate

620 condition. G1 – "least eel", "least ale", "Lee stale", "Lee steel"; G2 – "doom

621 art", "doom eat", "do mart", "do meet"; G3 – "coop art", "coop eat", "coo

622 part", "coo Pete".

623 To test hypothesis via model comparison, we use the Bayes Factor, which

624 can offer support for a model based on the observed data (Dienes, 2014,

625 Harm & Lakens, 2018). The posterior distributions of the model parameters

626 are not very informative as predictions need to be transformed with the

36

627    inverse logit function, and their details are included as supplementary

628    materials[4]. Therefore, the predicted distribution from 100 random samples

629    is shown in Fig. 14. As the figure shows, the model with an interaction term

630    shows the best predicative power. $BF_0$ was very close to zero (i.e. $BF_1$ is

631    larger than 10). Therefore, the data indicate that the alternative model, i.e.

632    onset and coda conditions are different, is highly more likely, because

633    model accuracy differs greatly. We also constructed a model with an

634    interaction effect between item group and syllable structure. $BF_{interaction}$ (the

635    BF indicating support for the interaction model over the full model) is larger

636    than 10, which provides strong support for the interaction model. To

637    conclude, the data shows strong evidence for the effect of syllable structure,

638    which differs greatly between groups. In other words, there is robust effect

639    of syllable structure for G2 and G3, but likely not for G1.

FIG. 14. Model predictions against 100 random samples for the slow rate, where y refers to the observed data and $y_{rep}$ refers to predictions. The columns correspond to item groups and the rows correspond to model type.

### 2. Results for normal speech rate

The accuracy distributions from the normal speech rate condition are shown in Fig. 15. Note that the coda condition only contained NN-resyllabified sequences. Fig. 15 shows that the amount of vowel information detected during the acoustic consonantal intervals (e.g. /s/ frication in "Lee stale") was very similar between the NN-resyllabified coda and onset sequences. The item group wise trends are similar to the slow rate condition in Fig. 13.

38

651 The aspiration from the plosive onset /p/ contains the most vowel related
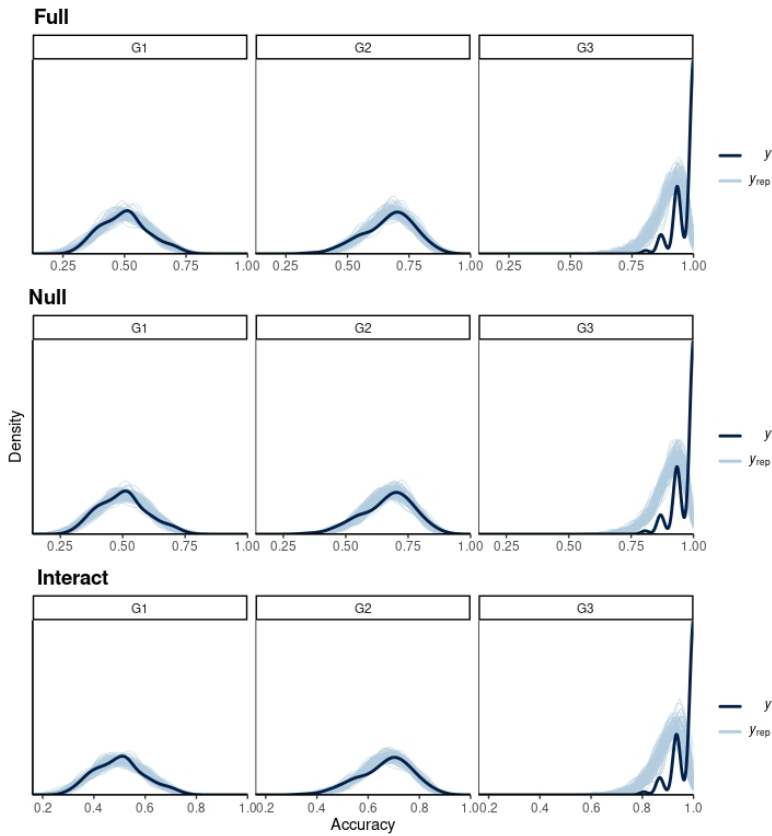
652 energy, and the nasal murmur from /m/ contained enough vowel information

653 for the classifier to perform above chance. For /s/ in G1, the accuracy

654 distributions are centered at chance level (i.e. 50%), indicating that little to

655 no vowel information was detected by the binary classifiers during the

656 frication intervals.



657

658 FIG. 15. Vowel classification accuracy by group from the normal speech

659 rate condition. The coda condition here refers to the NN-resyllabified coda

660 sequences in the normal speech rate condition. G1 – "least eel", "least ale",

661 "Lee stale", "Lee steel"; G2 – "doom art", "doom eat", "do mart", "do meet";

662 G3 – "coop art", "coop eat", "coo part", "coo Pete".

663 The predicted distributions from the Bayesian analysis results are shown in

664 Fig. 16. The posterior distributions of model parameters can be found in the

665 supplemented materials[5]. Visually, the predicted distributions do not differ

39

666　too much from one another. $BF_0$ was larger than 10, signifying that the data

667　provides more support for the null model. Fig. 15 indicates that model

668　accuracy might differ slightly between the NN-resyllabified coda and the

669　onset sequences for G1. In other words, there might be an interaction

670　between the effect of syllable structure and group. $BF_{interaction}$ (the BF

671　indicating support for the interaction model over the null model) is smaller

672　than 1/10, therefore, there is little to no evidence suggesting that accuracy

673　differs between onset and NN-resyllabified coda tokens for G1.

674



675　FIG. 16. Model predictions against 100 random samples for the normal rate,

676　where y refers to the observed data and $y_{rep}$ refers to predictions. The

677　columns correspond to item groups and the rows correspond to model type.

40

678    **D. Duration of intervocalic consonants**

679    The duration of the acoustic intervals for the canonical and NN-resyllabified

680    onsets are shown in Fig. 17. Congruent with previous findings (Gao & Xu,

681    2010; Lehiste, 1960), NN-resyllabified codas are shorter than the canonical

682    onsets. Predictions of the Bayesian analysis are shown in Fig. 18, and the

683    parameter posterior distributions are included as supplemented materials[6].

684    The effect of syllable structure was estimated to be around 0.01 ($\mu = 0.008$

685    [0.005, 0.012]). $BF_0$ is smaller than 1/10, which indicates that duration

686    differs between syllable structures.



687

688    FIG. 17. Duration of onset and NN-resyllabified consonants from the normal

689    speaking rate condition.

41

Fig. 18. Model predictions against 100 random samples for the duration results, where y refers to the observed data and $y_{rep}$ refers to predictions. The columns correspond to item groups and the rows correspond to model type.

## IV.    Discussion

Previous debates on the phenomenon of resyllabification have mainly relied on phonotactic analysis, listener judgment or phonetic properties such as voicing and aspiration. In this study we tested an alternative approach that examines articulatory coordination, and coarticulation, as reflected in the spectral patterns, using machine learning models with acoustic data. The findings have offered a new perspective on the nature of resyllabification.

42

## A. Overall findings

The results of computational analysis have largely confirmed the two predictions laid out in the introduction. The deep learning models trained on slow speech rate data misidentified coda sequences by classifying them as their onset counterparts, and DTW analysis showed that for all three consonants (i.e. /st/, /p/ and /m/), the sequences identified as resyllabified were more similar to their onset versions than the original coda versions. Moreover, the correctly classified sequences are more similar to their canonical coda version, which indicate that they likely have not undergone resyllabification. Therefore, the first prediction — codas in the NN-resyllabified sequences spectrally resemble canonical onsets more than their canonical coda version, was supported. The results from the binary classifiers confirm the second prediction by showing that there was a similar amount of vowel information detected in the NN-resyllabified onsets and canonical onsets, but not between the true codas and onsets from the slow condition. This suggests that the underlying articulation was alike between the NN-resyllabified and canonical onsets. Therefore, the results confirm previous findings of resyllabification in English (de Jong, 2001; Gao & Xu, 2010; Stetson, 1951). In connected speech, resyllabification can happen when a coda consonant is followed by a vowel initial syllable, and it applies to both singleton consonants and consonant clusters. The coda status of the NN-resyllabified consonants, however, seem to be partially retained through duration: Resyllabified codas are shorter compared to

43

725 canonical onsets. This is consistent with the findings of Lehiste (1960) and

726 more recently Gao and Xu (2010). Whether or not listeners can perceive the

727 durational cues, however, need to be tested in future studies. Furthermore,

728 future studies can investigate the effect of resyllabification and syllable

729 position on consonant duration by examining both NN-resyllabified and non

730 NN-resyllabified consonants.

731 It is also interesting to note the relation between resyllabification and

732 speech rate. When syllable duration is around 350 ms in the current study,

733 the rate of inferred resyllabification already reaches above 50%. At 2.86

734 syllables per second, this speech rate is rather slow, compared to the typical

735 normal articulation rate of 5-7 syllables per second in connected speech

736 (Eriksson, 2012; Tiffany, 1980). But this is consistent with the finding of de

737 Jong (2001) that resyllabification start to take place as speech rate

738 increases to around 350 ms per syllable, and resyllabification rate

739 approaches 100% at 150 ms per syllable. The implication is that the

740 tendency for resyllabification must be very strong so that it would be

741 difficult to avoidat normal speech rate.

742 The finding of resyllabification align with the syllable model shown in Fig. 1

743 based on which the predictions illustrated in Fig. 2 were derived. That is,

744 once a coda consonant is resyllabified as the onset of the next syllable, as

745 determined by the deep learning model and DTW analysis, its articulation is

746 overlapped with the vowel of the next syllable, as determined by the binary

44

747  classifiers. This is consistent with the recent finding that the movements

748  towards the vowel and onset C are synchronised at syllable onset (Liu et al.,

749  2022; Liu & Xu, 2021; Xu et al., 2019), which is denoted by the rime and

750  onset tiers in Fig. 1.

751  **B. Coarticulation resistance and dimension-specific sequential**

752  **target approximation (DSSTA)**

753  CV synchronisation does not mean that vowel information is always

754  detectable from the syllable onset or at the same time point, however,

755  partly due to *coarticulation resistance*, i.e. the ability of a segment to

756  restrain coarticulatory effects from adjacent segments (Bladon & Al-

757  Bamerni, 1976; Recasens, 1984). Recasens (1984) proposes that the degree

758  of coarticulation resistance is dependent on the amount of constraint that a

759  consonant or vowel places on the tongue body. Xu (2020) further proposes

760  that the phenomenon is a mechanism that resolves the articulatory conflicts

761  between consonants and vowels when they both involve the same

762  articulator while being co-produced to achieve C-V co-onset (Fig. 1).

763  According to this mechanism, namely, *dimension-specific sequential target*

764  *approximation mechanism*, different (e.g. vertical or horizontal) dimensions

765  of an articulator can be engaged in executing only a single target, which is

766  either consonantal or vocalic, during C-V coproduction. This mechanism

767  maximises the degree of C-V synchronisation while allowing individual

768  articulator dimensions to be engaged in only sequential target

769  approximation movements, i.e. without gestural blending (Saltzman &

45

770  Munhall, 1989) given its computational difficulty (Tilsen, 2019). The

771  following discussion will offer an account of the differences in the detected

772  vowel information in the present results that includes DSSTA as a critical

773  mechanism.

774  The amount of detectable vowel information in the consonant interval

775  follows the order of Group 1 (/s/) < Group 2 (/m/) < Group 3 (/p/). This order

776  may result from two different sources. The first, which is more obvious, is

777  the differences in their relative timing. The frication in Group 1 and nasal

778  murmur in Group 2 both correspond to the articulatory closure of the

779  consonants, whereas the aspiration in Group 3 corresponds to the

780  articulatory release, which occurs after the closure. This could partially

781  explain why more vowel information was detected in Group 3 than in the

782  other two groups. The second source is coarticulation resistance due to

783  DSSTA. The consonant /s/ in Groups 1 involves the tongue body to form a

784  groove needed to direct the airflow toward the front teeth (Borden, Harris

785  and Raphael, 2003). The involvement of the tongue body would generate

786  serious coarticulation resistance in /s/ in Group 1 because the horizontal

787  and vertical dimensions of the tongue body are likely both involved in

788  approaching the target of the sibilant (Recasens & Espinosa, 2009). In

789  contrast, the articulation of /m/ in Group 2 requires only lip closure without

790  constraints on the tongue. This would account for the greater amount of

791  detectable vowel information in Group 2 than in Group 1. The lack of tongue

792  involvement in labial consonants is true of /p/ in Group 3 as well. But there,

46

793 it is added on top of the fact that aspiration, where the binary classification

794 was performed, occurs after the stop closure, thus giving rise to the

795 maximal vowel information detected by the classifier. Note that had one of

796 the syllables in Group 1 contained a rounded vowel such as /u/, DSSTA

797 would predict that vowel information would be better detected, because lip

798 movements are not in direct conflict with the articulation of /s/. This

799 possibility can be tested in future research.

800 **C. Chance level performance of the binary classifier for G1**
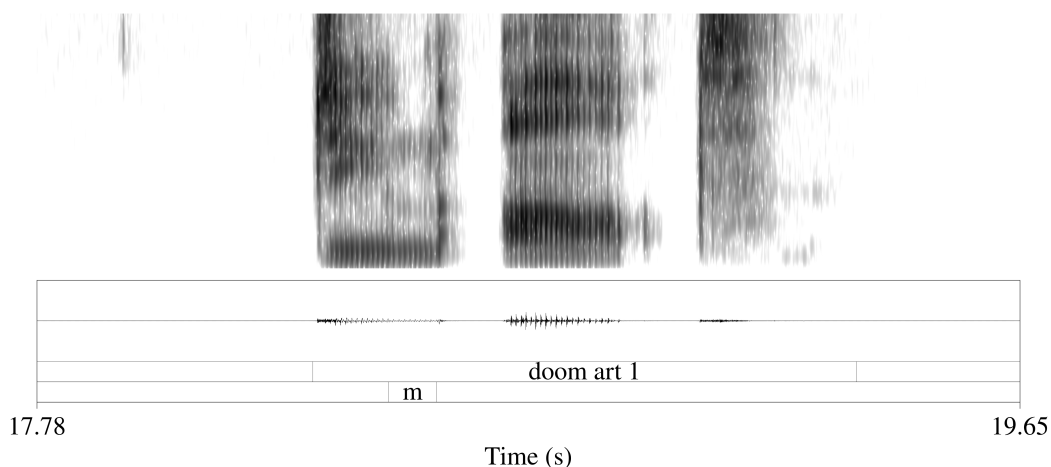
801 **sequences**

802 The lack of detectable vowel information in /st/ even in normal speech rate

803 may seem to contradict the recent finding that vowel articulation could be

804 detected at the same time as the onset of a consonant cluster (Liu & Xu,

805 2021). That study found that for a minimal triplet such as "slit" vs. "slot" vs.

806 "flot", the difference between "slit" and "slot" could be detected around the

807 same time as "slot" and "flot", *before* the frication onset. But we have noted

808 three major differences between Liu and Xu (2021) and the current study.

809 First, Liu and Xu (2021) only looked at clusters such as /sp/ and /sl/, but

810 not /st/ as in the current study. /p/ does not require any tongue movement,

811 thus is less coarticulation resistant than both /l/ and /t/. In terms of /l/ and

812 /t/, both being alveolars, Iskarous et al. (2013) found that /t/ is more

813 coarticulation resistant than /l/ in the vertical dimension for the jaw and the

814 tongue blade. This could be due to the requirement of a full closer for /t/ as

815 a plosive but not for the approximant /l/. /t/ being more coarticulation

47

816 resistant means that it may have delayed much of the vowel movements.

817 Second, much larger vowel contrasts were involved in Liu and Xu

818 (2021)—/slit/ vs. /slot/ than those in the present study—/steal/ vs. /stale/.

819 The greater the vowel contrast, the greater the magnitude of tongue

820 movement in the articulatory dimensions not essential for the consonant

821 articulation, and the more detectable the vowel information during the

822 frication interval. Third, the target words were produced with a carrier in

823 Liu and Xu (2021), which made the speech more fluent than the isolated

824 word sequences said in the present study. The average speech rate in Liu

825 and Xu (2021) was about 140 ms per syllable, compared to 350 ms per

826 syllable in this study. It is hard to tell, however, if any of these factors is

827 decisive, or all of them jointly contribute to blocking the vowel information

828 from being present in the /s/ frication.

829 **D. Above chance performance of the binary classifier for the slow**

830 **coda sequence in G2**

831 One of the most surprising results of this study is the finding that, as shown

832 in Fig. 13, for the slow speaking rate, there is information of the upcoming

833 vowel in the intervocalic consonants when they are in the coda position of

834 the first syllable (e.g. "doom art"; "coop art"), albeit less than when they are

835 in the onset position. The detection of vowel information in a non

836 resyllabified coda may seem particularly striking given the clear temporal

837 gap or glottalisation between the two syllables, as can be seen in Fig. 19

838 and Fig. 20. But the glottal component, as can be judged both auditorily and

48

839  spectrally, corresponds to a glottal stop or glotallisation (which is also a

840  form of glottal stop: Redi & Shattuck-Hufnagel, 2001; Garellek, 2013), that

841  serves as the onset of the syllable /art/. A glottal stop, just like other stops

842  such as /b, d, g/, would be fully coarticulated with the following vowel (Xu,

843  2020), as illustrated in Fig. 2. This means that the target approximation

844  of /a/ must have started some time well before the glottal closure (Liu et al.,

845  2022; Xu and Liu, 2007). This can indeed be seen in Fig. 19, i.e. the brief

846  yet clearly visible labial release after the nasal murmur of /m/, and the F2

847  transition from "doom" to "eat" during and right before the glottalised

848  interval in Fig. 20. The high vowel detection rate of around 80% for /p/ and

849  65% for /m/ means that the vowel target approximation may have started

850  during (though probably not before) the closure of the coda. Exactly when

851  during the closure, however, awaits future investigations.



852

853  FIG. 19. Spectrogram of "doom art" from a male speaker.

49

854

FIG. 20. Spectrogram of "doom eat" from a female speaker.

## E. Broader implications

The finding of a clear tendency toward resyllabification in this study provides further support for the synchronisation model of the syllable (Xu, 2020) beyond recent findings (Liu et al., 2022; Liu & Xu et al., 2021). According to the model, there is a strong demand for onset consonants to synchronise (i.e. fully overlap) with the vowel, and a high time pressure against the preservation of coda consonants. This is partially consistent with the maximum onset principle (Pulgram, 1970; Selkirk, 1982), but offers specific articulatory details that can be tested in the acoustic signals as done in the present study. Because the syllable is both essential and highly controversial for theoretical models in linguistics as well as psycholinguistics, the current results may have implications for many broader issues about speech production, but here we focus only on two major ones. The first is about the influential psycholinguistic model of

50

870  speech production (Levelt et al., 1999), which proposes a step-by-step

871  model of how speech production proceeds from lexical selection to

872  articulation. The results of the present study are relevant for the

873  phonological encoding to articulation stages in the model. The most

874  relevant result is probably the corroboration of previous findings that

875  resyllabification is contingent on local articulation rate: highly likely at

876  normal rate, but optional at slow rate (de Jong, 2001; Stetson, 1951). This

877  means that until local speech rate is known, the articulatory affiliation of

878  coda consonant is undetermined, which would suggest that either syllables

879  retrieved from memory (during phonological encoding) are incomplete in

880  terms of segment affiliation, or the retrieved syllables are reorganised by

881  resyllabification, and that this reorganisation would occur *after* the phonetic

882  encoding stage, just before articulation.

883  The finding of rate-dependency of resyllabification is further relevant to any

884  psycholinguistic model of speech production given the known extensive use

885  of speech timing by linguistic functions. Specifically, local articulation rate,

886  which is jointly determined by syllable duration and pause duration, is used

887  to encode multiple levels of boundary strength (Lehiste, 1972; Klatt, 1976;

888  Nakatani, O'Connor & Aston, 1981; Wagner, 2005; Wang & Xu, 2019). Thus

889  resyllabification is likely a regular variable of connected speech beyond

890  word-level phonetics. In fact, it is likely part of the process of producing

891  connected speech that involves many other phonetic reorganisations,

892  including deletion of intervocalic coda (as opposed to resyllabification in

51

893 some languages) (e.g. tone sandhi (Chen, 2000), intrusive /r/ (Gick, 1999),

894 and vowel hiatus breakers (Mudzingwa, 2013), etc). There is already

895 evidence that some of these reorganisations may be cognitively real, at least

896 in the case of tone sandhi (Zhang, Xia & Peng, 2015). These phonetic

897 reorganisation tactics could therefore be included in an enhanced

898 psycholinguistic model of speech production, and their cognitive reality

899 could be experimentally investigated.

900 The second broad issue is whether the present results can be interpreted in

901 terms of ambisyllabicity. The original proposal of ambisyllabicity was

902 motivated by the lack of phonetic means to clearly determine syllable

903 boundaries, so the affiliation of intervocalic segments had to rely on

904 phonotactic well-formedness, and for cases where *ill-formed* syllables would

905 occur if an intervocalic consonant can only have a single affiliation, e.g.,

906 *happy*, *attic*, *hobby*, the solution is ambisyllabicity, i.e., simultaneously

907 affiliation to both adjacent syllables (Kahn, 1976). Exactly how such double

908 association is realised phonetically, however, has remained unclear. Gick

909 (2003) has proposed that some intervocalic segments, e.g. /l/ and /w/,

910 actually consist of a C-gesture and a V-gesture, which are simultaneously

911 phased to the surrounding syllables, therefore ambisyllabified. The phonetic

912 evidence is in terms of different time delays in the *achievement of the*

913 *respective C and V gestural goals, which differs from the onset alignment*

914 *the current study has examined.* Although the present study is not designed

915 for examining ambisyllabicity, at least one phonetic cue is shown to have

52

916 the potential to indicate the original coda status of a consonant, namely, the

917 shorter duration of NN-resyllabified coda than the original onset consonant

918 (also c.f. Lehiste, 1960). However, if CV onset coarticulation is considered

919 as the sole indicator, the NN-resyllabified codas are unambiguously

920 overlapped with the following vowel according to the present data.

921 **F. Caveats**

922 Two of the resyllabification classifiers satisfied the early stopping criteria,

923 which meant that their training epochs were determined with the test split

924 rather than the pilot data. This could have slightly inflated the overall

925 accuracy reported for the slow condition in section III A. However, the use

926 of the classifier is to classify normal rate sequences, which is the focus of

927 the study and their accuracy has not been inflated as the normal rate data

928 were not used in any way during training.

929 The possibility of false negatives cannot be completely ruled out regarding

930 the chance level performance of the binary classifier for G1. Providing that

931 upcoming vowel related acoustic information exist during frication, two

932 scenarios could result in false negative detections:

933 1. Chance performance due to chance.

934 2. The neural networks are not powerful enough to detect the subtle

935 difference.

53

936 The first scenario refers to the opposite of what is described in Combrisson
937 and Jerbi (2015), namely, the model achieved chance performance by
938 chance. This could be due to the randomised nature of the data split and/or
939 model parameter initialisation (not hyperparameters). However, this
940 possibility is accounted for in the current study, by repeatedly training 80
941 classifiers on randomised train and test data and analysing the resultant
942 accuracy distributions. For the second scenario, despite tuning the
943 hyperparameters with data from pilot recordings, the neural network was
944 not tuned for each speaker and consonant type separately. In practice, it is
945 very difficult to construct a perfect network regardless of the type of data in
946 question. Therefore, there is a small possibility that the binary classifier
947 could not detect a difference between groups in G1 due to the lack of
948 robustness. Future studies could incorporate articulatory data, as it might
949 provide more detailed information than acoustic data in the current study
950 (Tilsen, 2020).

951 On the other hand, the possibility of false positives cannot be ruled out
952 either. Providing that the test dataset is large enough, machine learning
953 models cannot always achieve 100% accuracy. The same applies to the
954 word sequence classifiers in this study. This is evident in the results from
955 the slow speech rate in section III A. Although overall accuracy is high,
956 there were still coda sequences classified as their onset counterpart, as well
957 as cases where onset sequences were classified as their coda counterpart.
958 At the slow speech rate (2 syllables per second on average), is it unlikely

54

959   that resyllabification occurred, so these misclassifications are likely genuine

960   incorrect classifications (i.e. not due to syllabification). As for the normal

961   rate results, there should also exist genuine misidentifications like the slow

962   rate, which is likely why there are onset sequences classified as their coda

963   counterparts. This means that a small number of the NN-resyllabified

964   sequences might be genuine misidentification as well. However, the normal

965   rate results show that onset sequences reached an accuracy rate of 90%

966   and only 36% was achieved for the coda ones. Therefore, a large portion of

967   the NN-resyllabified tokens are likely due to syllabification structure and

968   not just simple false positives.

969   Also, the study did not conduct a parallel analysis of $V_2$ binary classification

970   for the correctly classified coda tokens. Unlike the DTW analysis, there are

971   too few correctly classified coda sequences in the normal rate for training

972   neural network classifiers, especially for G1 and G2. This issue is

973   exacerbated by the imbalance of speakers in the data, i.e. some speakers

974   had zero or a very small number of correctly classified tokens in certain

975   item groups. Future study can potentially avoid this issue by increasing the

976   number of repetitions in the normal rate condition.

977   Finally, as noted in section IV B, the lack of detectable vowel information in

978   Group 1 might have been avoided had one of the syllables in each pair

979   contained a rounded vowel. This is because, despite its involvement of the

980   tongue-body, the articulation of /s/ is not in direct conflict with the lip

55

981 movements of the co-produced vowel. This possibility can be investigated in

982 future research.

## V.    CONCLUSION

984 We used deep learning models with acoustic data to investigate the

985 phenomenon of resyllabification. The models trained on slow speech data

986 can be used to infer resyllabified sequences in normal speech rate data.

987 This was verified by DTW analysis, which revealed that, compared to slow

988 speech, NN-resyllabified sequences were more similar to the true onset

989 sequences than their original coda productions. The acoustic intervals of

990 intervocalic consonants were examined with bi-directional recurrent neural

991 network models. We found that similar amount of vowel information was

992 detected in the intervocalic consonants between the NN-resyllabified codas

993 and the genuine onsets, suggesting that the coarticulation structure of the

994 former resembles that of the latter. For slow speech rate, the results show

995 that the articulatory structures likely differed between the onset and coda

996 sequences. Surprisingly, however, vowel information can still be detected

997 from the closure and release of labial coda consonants, indicating that the

998 articulation of the vowel has started during the acoustic interval of a coda

999 consonant even when it is not resyllabified.

## APPENDIX

1001 The hyperparameter details for the multi-class classifier and the binary

1002 classifiers are shown in Table II and Table III, respectively.

56

1003    TABLE II. Hyperparameters for the multi-class classifiers.

| Hyperparameter | Value |
| --- | --- |
| Number of residual blocks | 3 |
| Number of GRU layers | 4 |
| Number of units in the GRU layers | 512 |
| Number of units in the linear layers | 512 |
| Dropout rate | 0.1 |
| Number of channels for the CNN layers | 32 |
| Batch size | 32 |
| Learning rate | 0.0001 |
| Optimiser | RMSprop |
| Epoch number | 120 |

1004

1005    TABLE III. Hyperparameters for the binary classifiers.

| Hyperparameter | Value |
| --- | --- |
| Number of units in the first LSTM | 60 |

57

| | |
|---|---|
| layer | |
| Number of units in the second LSTM layer | 30 |
| Dropout rate for the first LSTM layer | 0.1 |
| Dropout rate for the second LSTM layer | 0.2 |
| Number of units in the linear layer | 50 |
| Merge mode | Summation |
| Batch size | 16 |
| Optimiser | Adam |
| Learning rate | 0.001 |
| Epoch number | 70 |

1006

1007 [1]During data splitting, correlated samples due to augmentation were not

1008 included in the same dataset. e.g. the original "coo part" and its augmented

1009 version always ended up in the same split.

1010 [2]The full detail of models and data processing can be found at

1011 https://github.com/Clara-liu/deep_speech_resyllabification

58

1012   [3]The details of implementation of custom one-inflated-beta-distribution are

1013   available at

1014   https://github.com/Clara-liu/deep_speech_resyllabification/blob/main/

1015   one_inflated_beta.R

1016   [4]See supplementary materials at [URL] for details on the posterior

1017   distributions for the slow rate condition.

1018   [5]See supplementary materials at [URL] for details on the posterior

1019   distributions for the normal rate condition.

1020   [6]See supplementary materials at [URL] for details on the posterior

1021   distributions for the duration analysis.

1022   **REFERENCES (BIBLIOGRAPHIC)**

1023   Adda-Decker, M., de Mareüil, P. B., Adda, G., & Lamel, L. (2002).

1024        Investigating syllabic structure and its variation in speech. In

1025        *Pronunciation Modeling and Lexicon Adaptation for Spoken Language*

1026        *Technology* (p. 6).

1027   Amodei, D., Anubhai, R., Battenberg, E., Case, C., Casper, J., Catanzaro, B.,

1028        Chen, J., Chrzanowski, M., Coates, A., Diamos, G., Elsen, E., Engel, J.,

1029        Fan, L., Fougner, C., Han, T., Hannun, A., Jun, B., LeGresley, P., Lin,

1030        L., ... Zhu, Z. (2015). *Deep Speech 2: End-to-End Speech Recognition*

1031        *in English and Mandarin.* https://doi.org/10.48550/ARXIV.1512.02595

1032    Audacity Team. (2021). *Audacity*. https://audacityteam.org/

1033    Balakrishnan, N., & Nevzorov, V. (2003). A primer on statistical distribu-
1034          tions. Hoboken, NJ: John Wiley and Sons.

1035    Barlow, J. A., & Gierut, J. A. (1999). Optimality Theory in Phonological
1036          Acquisition. *Journal of Speech, Language, and Hearing Research*,
1037          *42*(6), 1482–1498. https://doi.org/10.1044/jslhr.4206.1482

1038    Bartelds, M., Richter, C., Liberman, M., & Wieling, M. (2020). A New
1039          Acoustic-Based Pronunciation Distance Measure. *Frontiers in*
1040          *Artificial Intelligence*, *3*, 39. https://doi.org/10.3389/frai.2020.00039

1041    Bermúdez-Otero, R. 2011. Cyclicity. In van Oostendorp, M., Ewen, C J.,
1042          Hume, E., & Rice, K. (Eds.), Blackwell companion to phonology.
1043          Chichester: Wiley-Blackwell, Vol. 4, pp. 2019–2048.

1044    Biel, A. L., & Friedrich, E. V. C. (2018). Why You Should Report Bayes
1045          Factors in Your Transcranial Brain Stimulation Studies. *Frontiers in*
1046          *Psychology*, *9*, 1125. https://doi.org/10.3389/fpsyg.2018.01125

1047    Birgit A. (2001). Regional Variation and Edges: Glottal Stop Epenthesis and
1048          Dissimilation in Standard and Southern Varieties of German.
1049          *Zeitschrift Für Sprachwissenschaft*, *20*(1), 3–41.
1050          https://doi.org/doi:10.1515/zfsw.2001.20.1.3

1051    Bladon, R. A. W., & Al-Bamerni, A. (1976). Coarticulation resistance in

1052        English /l/. *Journal of Phonetics*, *4*(2), 137–150.

1053        https://doi.org/10.1016/S0095-4470(19)31234-3

1054    Borden, G. J., Harris, K. S., and Raphael, L. J. (2003). *Speech Science*

1055        *Primer: Physiology, Acoustics, and Perception of Speech, 4th Edition*

1056        (Williams & Wilkins, Baltimore).

1057    Blevins, J. (2003). *Evolutionary phonology: The emergence of sound*

1058        *patterns*. Cambridge University Press.

1059    Boersma, P., & Weenink, D. (2022). *Praat: Doing phonetics by computer*

1060        (6.2.14) [Computer software]. http://www.praat.org/

1061    Browman, C. P., & Goldstein, L. (1992). Articulatory Phonology: An

1062        Overview. *Phonetica*, *49*(3–4), 155–180.

1063        https://doi.org/10.1159/000261913

1064    Chen, M. Y. (2000). *Tone Sandhi: Patterns across Chinese Dialects*

1065        (Cambridge University Press, Cambridge, UK).

1066    Clements, G. N., & Keyser, S. J. (1983). CV phonology. A generative theory

1067        of the syllable. *Linguistic Inquiry Monographs Cambridge*, *9*, 1–191.

1068    Combrisson, E., & Jerbi, K. (2015). Exceeding chance level by chance: The

1069        caveat of theoretical chance levels in brain signal classification and

1070        statistical assessment of decoding accuracy. *Journal of Neuroscience*

61

1071    *Methods*, *250*, 126–136.

1072    https://doi.org/10.1016/j.jneumeth.2015.01.010

1073    Content, A., Kearns, R. K., and Frauenfelder, U. H. (2001). Boundaries

1074    versus Onsets in Syllabic Segmentation. *Journal of Memory and*

1075    *Language* **45**, 177-199.

1076    de Jong, K. J. (2001). Rate-Induced Resyllabification Revisited. *Language*

1077    *and Speech*, *44*(2), 197–216.

1078    https://doi.org/10.1177/00238309010440020401

1079    de Jong, K. J., Lim, B., & Nagao, K. (2004). The Perception of Syllable

1080    Affiliation of Singleton Stops in Repetitive Speech. Language and

1081    Speech, 47(3), 241–266.

1082    Dienes, Z. (2014). Using Bayes to get the most out of non-significant results.

1083    *Frontiers in Psychology*, *5*. https://doi.org/10.3389/fpsyg.2014.00781

1084    Dienes, Z. (2016). How Bayes factors change scientific practice. *Journal of*

1085    *Mathematical Psychology*, *72*, 78–89.

1086    https://doi.org/10.1016/j.jmp.2015.10.003

1087    Douma, JC & Weedon, JT. (2019). Analysing continuous proportions in

1088    ecology and evolution: A practical introduction to beta and Dirichlet

1089    regression. *Methods Ecol Evol*. 10, 1412– 1430.

1090   Eriksson, A. (2012). Aural/acoustic vs. automatic methods in forensic
1091       phonetic case work. in *Forensic Speaker Recognition* (Springer), pp.
1092       41-69.

1093   Gao, H., & Xu, Y. (2010). Ambisyllabicity in English: How real is it?  *In*
1094       *Proceeding of the 9th Phonetic Conference of China*, Tianjin.

1095   Gao, M. (2009). Gestural coordination among vowel, consonant and tone
1096       gestures in Mandarin Chinese. *Chinese Journal of Phonetics*, *2*, 43–50.

1097   Garellek, M. (2012). Glottal stops before word-initial vowels in American
1098       English: distribution and acoustic characteristics. *UCLA Working*
1099       *Papers in Phonetics*, 110, 1-23.

1100   Garellek, M. (2013). Production and perception of glottal stops (Doctoral
1101       dissertation, UCLA). Retrieved from
1102       https://escholarship.org/uc/item/7zk830cm

1103   Gaskell, M. G., Spinelli, E., & Meunier, F. (2002). Perception of
1104       resyllabification in French. *Memory & Cognition*, 30(5), 798–810.
1105       https://doi.org/10.3758/BF03196435

1106   Gelman, A., Carlin, J.B., Stern, H.S., Dunson, D.B., Vehtari, A., & Rubin,
1107       D.B. (2013). Bayesian Data Analysis (3rd ed.). Chapman and Hall/CRC.
1108       https://doi.org/10.1201/b16018

63

1109    Gick, B. (1999). A gesture-based account of intrusive consonants in English.

1110        *Phonology*, *16*(1), 29–54. https://doi.org/10.1017/S0952675799003693

1111    Gick, B. (2003). Articulatory correlates of ambisyllabicity in English glides

1112        and liquids. Phonetic interpretation: Papers in laboratory phonology,

1113        6, 222-236.

1114    Goldstein, L., Byrd, D., & Saltzman, E. (2006). The role of vocal tract

1115        gestural action units in understanding the evolution of phonology. In

1116        *Action to Language via the Mirror Neuron System* (pp. 215–249).

1117        https://doi.org/10.1017/CBO9780511541599.008

1118    Goslin, J., & Frauenfelder, U. H. (2001). A Comparison of Theoretical and

1119        Human Syllabification. *Language and Speech*, *44* (4), 409–436.

1120        https://doi.org/10.1177/00238309010440040101

1121    Gronau, Q.F., Wagenmakers, EJ. Limitations of Bayesian Leave-One-Out

1122        Cross-Validation for Model Selection. *Comput Brain Behav,* 2, 1–11

1123        (2019). https://doi.org/10.1007/s42113-018-0011-7

1124    Harms, C., & Lakens, D. (2018). Making 'Null Effects' Informative:

1125        Statistical Techniques and Inferential Frameworks. *Journal of Clinical*

1126        *and Translational Research*, 24.

1127    He, K., Zhang, X., Ren, S., & Sun, J. (2015). *Deep Residual Learning for*

1128        *Image Recognition* (arXiv:1512.03385). arXiv.

1129        http://arxiv.org/abs/1512.03385

1130 Ioffe, S., & Szegedy, C. (2015). *Batch Normalization: Accelerating Deep*

1131    *Network Training by Reducing Internal Covariate Shift*

1132    (arXiv:1502.03167). arXiv. http://arxiv.org/abs/1502.03167

1133 Iskarous, K., Mooshammer, C., Hoole, P., Recasens, D., Shadle, C. H.,

1134    Saltzman, E., & Whalen, D. H. (2013). The coarticulation/invariance

1135    scale: Mutual information as a measure of coarticulation resistance,

1136    motor synergy, and articulatory invariance. *The Journal of the*

1137    *Acoustical Society of America*, *134*(2), 1271–1282.

1138    https://doi.org/10.1121/1.4812855

1139 Jacewicz, E., Fox, R. A., O'Neill, C., & Salmons, J. (2009). Articulation rate

1140    across dialect, age, and gender. *Language variation and change*,

1141    *21*(2), 233–256. https://doi.org/10.1017/S0954394509990093

1142 Jeffreys, H. (1961). *The theory of probability (3rd ed.*). Oxford University

1143    Press.

1144 Kahn, D. (1976). Syllable-based generalizations in English phonology.

1145    ((Doctoral dissertation, MIT).

1146 Klatt, D. H. (1976). Linguistic uses of segmental duration in English:

1147    Acoustic and perceptual evidence. *J. Acoust. Soc. Am.* 59, 1208-1221.

1148 Ko, T., Peddinti, V., Povey, D., & Khudanpur, S. (2015). *Audio augmentation*

1149    *for speech recognition*. Interspeech 2015, 3586–3589.

1150    https://doi.org/10.21437/Interspeech.2015-711

1151 Kogan, J. A., & Margoliash, D. (1998). Automated recognition of bird song

1152      elements from continuous recordings using dynamic time warping and

1153      hidden Markov models: A comparative study. *The Journal of the*

1154      *Acoustical Society of America*, *103*(4), 2185–2196.

1155      https://doi.org/10.1121/1.421364

1156 Lakens, D., McLatchie, N., Isager, P. M., Scheel, A. M., & Dienes, Z. (2020).

1157      Improving Inferences About Null Effects With Bayes Factors and

1158      Equivalence Tests. *The Journals of Gerontology: Series B*, *75*(1), 45–

1159      57. https://doi.org/10.1093/geronb/gby065

1160 Lee, M., & Wagenmakers, E. (2014). Bayesian Cognitive Modeling: A

1161      Practical Course. Cambridge: Cambridge University Press.

1162      doi:10.1017/CBO9781139087759

1163 Lehiste, I. (1960). An acoustic-phonetic study of internal open juncture.

1164      *Phonetica Supplement.*

1165 Lehiste, I. (1972). The timing of utterances and linguistic boundaries. *J.*

1166      *Acoust. Soc. Am.* 51, 2018-2024.

1167 Lerato, L., & Niesler, T. (2019). Feature trajectory dynamic time warping

1168      for clustering of speech segments. *EURASIP Journal on Audio,*

1169      *Speech, and Music Processing*, *2019* (1), 6.

1170      https://doi.org/10.1186/s13636-019-0149-9

1171  Levelt, W. J. M., Roelofs, A., & Meyer, A. S. (1999). A theory of lexical

1172       access in speech production. *Behavioral and Brain Sciences*, *22*(01).

1173       https://doi.org/10.1017/S0140525X99001776

1174  Liu, Z., & Xu, Y. (2021). *Segmental Alignment of English Syllables with*

1175       *Singleton and Cluster Onsets*. Interspeech 2021, 3969–3973.

1176       https://doi.org/10.21437/Interspeech.2021-187

1177  Liu, Z., Xu, Y., & Hsieh, F. (2022). Coarticulation as synchronised CV co-

1178       onset – Parallel evidence from articulation and acoustics. *Journal of*

1179       *Phonetics*, *90*, 101116. https://doi.org/10.1016/j.wocn.2021.101116

1180  Luo, D., Zou, Y., & Huang, D. (2018). *Investigation on Joint Representation*

1181       *Learning for Robust Feature Extraction in Speech Emotion*

1182       *Recognition*. Interspeech *2018*, 152–156.

1183       https://doi.org/10.21437/Interspeech.2018-1832

1184  MacNeilage, P. F. (1998). The frame/content theory of evolution of speech

1185       production. *Behavial and Brain Sciences* 21, 499–546.

1186  Marin, S., & Pouplier, M. (2014). Articulatory synergies in the temporal

1187       organization of liquid clusters in Romanian. *Journal of Phonetics*, *42*,

1188       24–36. https://doi.org/10.1016/j.wocn.2013.11.001

1189  Mirzaei, M. S., Meshgi, K., & Kawahara, T. (2018). Exploiting automatic

1190       speech recognition errors to enhance partial and synchronized

67

1191        caption for facilitating second language listening. *Computer Speech &*

1192        *Language*, *49*, 17–36. https://doi.org/10.1016/j.csl.2017.11.001

1193  Mok, P. P. K. (2012). Effects of consonant cluster syllabification on vowel-to-

1194        vowel coarticulation in English. *Speech Communication*, *54*(8), 946–

1195        956.

1196  Morey, R. D., Romeijn, J. W., & Rouder, J. N. (2016). The philosophy of

1197        Bayes factors and the quantification of statistical evidence. *Journal of*

1198        *Mathematical Psychology*, 72, 6-18.

1199  Mudzingwa, C. (2013). Hiatus resolution strategies in Karanga (Shona).

1200        *Southern African Linguistics and Applied Language Studies*, *31*(1), 1–

1201        24. https://doi.org/10.2989/16073614.2013.793953

1202  Mullooly, R. (2003). *An Electromagnetic Articulography study of*

1203        *resyllabification of rhotic consonants in English*. International

1204        Conference of Phonetic Sciences, Barcelona.

1205  Nakatani, L. H., O'connor, K. D., and Aston, C. H. (1981).Prosodic aspects of

1206        American English speech rhythm. *Phonetica* 38, 84-106.

1207  Nam, H. (2007). *Articulatory modeling of consonant release gesture*.

1208        International Conference of Phonetic Sciences, Saarbrücken.

1209  Nam, H., Goldstein, L., & Saltzman, E. (2009). Self-organization of syllable

1210        structure: A coupled oscillator model. In F. Pellegrino, E. Marsico, I.

1211       Chitoran, & C. Coupé (Eds.), *Approaches to Phonological Complexity*

1212       (pp. 297–328). Walter de Gruyter.

1213       https://doi.org/10.1515/9783110223958.297

1214  Nam, H., Mitra, V., Tiede, M., Hasegawa-Johnson, M., Espy-Wilson, C.,

1215       Saltzman, E., & Goldstein, L. (2012). A procedure for estimating

1216       gestural scores from speech acoustics. *The Journal of the Acoustical*

1217       *Society of America*, *132*(6), 3980–3989.

1218       https://doi.org/10.1121/1.4763545

1219  Ni Chiosáin, M. N., Welby, P., & Espesser, R. (2012). Is the syllabification of

1220       Irish a typological exception? An experimental study. *Speech*

1221       *Communication*, *54*(1), 68–91.

1222       https://doi.org/10.1016/j.specom.2011.07.002

1223  Ospina, R., Ferrari, S L.P. (2012). A general class of zero-or-one inflated

1224       beta regression models. *Computational Statistics & Data Analysis,*

1225       56(6), 1609-1623.

1226  Ojala, M., & Garriga, G. C. (2009). Permutation Tests for Studying Classifier

1227       Performance. *2009 Ninth IEEE International Conference on Data*

1228       *Mining*, 908–913. https://doi.org/10.1109/ICDM.2009.108

1229  Park, D. S., Chan, W., Zhang, Y., Chiu, C.-C., Zoph, B., Cubuk, E. D., & Le,

1230       Q. V. (2019). *SpecAugment: A Simple Data Augmentation Method for*

*Automatic Speech Recognition*. Interspeech 2019, 2613–2617.

https://doi.org/10.21437/Interspeech.2019-2680

Pastätter, M., & Pouplier, M. (2014). *The articulatory modeling of German coronal consonants using TADA*. International seminar on speech production, Cologne, Germany.

Perkell, J., & Chiang, C. M. (1986). *Preliminary support for a 'hybrid' model of anticipatory coarticulation*. International congress of acoustics, Toronto.

Pervaiz, A., Hussain, F., Israr, H., Tahir, M. A., Raja, F. R., Baloch, N. K., Ishmanov, F., & Zikria, Y. B. (2020). Incorporating Noise Robustness in Speech Command Recognition by Noise Augmentation of Training Data. *Sensors*, *20*(8), 2326. https://doi.org/10.3390/s20082326

Prom-on, S., Xu, Y., & Thipakorn, B. (2009). Modeling tone and intonation in Mandarin and English as a process of target approximation. *The Journal of the Acoustical Society of America*, *125*(1), 405–424. https://doi.org/10.1121/1.3037222

Pulgram, E. (1970). *Syllable, word, nexus, cursus* (Mouton, The Hague).

Recasens, D. (1984). Vowel-to-vowel coarticulation in Catalan VCV sequences. *J. Acoust. Soc. Am.* 76, 1624-1635.

Recasens, D., & Espinosa, A. (2005). Articulatory, positional and coarticulatory characteristics for clear /l/ and dark /l/: Evidence from two Catalan dialects. *Journal of the International Phonetic Association*, *35*(1), 1–25. https://doi.org/10.1017/S0025100305001878

Recasens, D., & Espinosa, A. (2009). An articulatory investigation of lingual coarticulatory resistance and aggressiveness for consonants and vowels in Catalan. *The Journal of the Acoustical Society of America*, *125*(4), 2288–2298. https://doi.org/10.1121/1.3089222

Redi, L., and Shattuck-Hufnagel, S. (2001). Variation in the realization of glottalization in normal speakers. *J. Phonetics,* 29, 407-429.

Sakoe, H., & Chiba, S. (1978). Dynamic programming algorithm optimization for spoken word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, *26*(1), 43–49. https://doi.org/10.1109/TASSP.1978.1163055

Saltzman, E., & Byrd, D. (2000). Task-dynamics of gestural timing: Phase windows and multifrequency rhythms. *Human Movement Science*, *19*(4), 499–526. https://doi.org/10.1016/S0167-9457(00)00030-0

Saltzman, E. L., & Munhall, K. G. (1989). A Dynamical Approach to Gestural Patterning in Speech Production. *Ecological Psychology*, *1*(4), 333–382. https://doi.org/10.1207/s15326969eco0104_2

1270 Schiller, N. O., Mever, A. S., & Levelt, W.J. (1997). The syllabic structure of
1271       spoken words: Evidence from the syllabification of intervocalic
1272       consonants. *Language and Speech*, *40*(2), 103–140.

1273 Schönbrodt, F. D. & Wagenmakers, E.-J. Bayes factor design analysis:
1274       planning for compelling evidence. *Psychon. Bull. Rev.* 25, 128–142
1275       (2018).

1276 Selkirk, E. O. (1982). "The syllable," in *The structure of phonological*
1277       *representations, Part II*, edited by H. v. d. Hulst, and N. Smith (Foris
1278       Publications, Dordrecht, The Netherlands), pp. 337–383.

1279 Semeniuta, S., Severyn, A., & Barth, E. (2016). *Recurrent Dropout without*
1280       *Memory Loss* (arXiv:1603.05118). arXiv.
1281       http://arxiv.org/abs/1603.05118

1282 Sharma, J., Granmo, O.-C., & Goodwin, M. (2020). *Environment Sound*
1283       *Classification Using Multiple Feature Channels and Attention Based*
1284       *Deep Convolutional Neural Network*. Interspeech 2020, 1186–1190.
1285       https://doi.org/10.21437/Interspeech.2020-1303

1286 Shattuck-Hufnagel, S. (2011). The role of the syllable in speech production
1287       in American English: A fresh consideration of the evidence. In
1288       *Handbook of the Syllable* (pp. 197–224). Brill.

1289 Shaw, J. A., Gafos, A. I., Hoole, P., & Zeroual, C. (2011). Dynamic invariance
1290       in the phonetic expression of syllable structure: A case study of

1291        Moroccan Arabic consonant clusters. *Phonology*, *28*(3), 455–490.

1292        https://doi.org/10.1017/S0952675711000224

1293  Smith, J. L. (2001). Lexical Category and Phonological Contrast. *Workshop*

1294        *on the Lexicon*, 61–72.

1295  Soltau, H., Liao, H., & Sak, H. (2016). *Neural Speech Recognizer: Acoustic-*

1296        *to-Word LSTM Model for Large Vocabulary Speech Recognition*

1297        (arXiv:1610.09975). arXiv. http://arxiv.org/abs/1610.09975

1298  Steriade, D. (1999). Alternatives to syllable-based accounts of consonantal

1299        phonotactics. *Item Order in Language and Speech*, 205–245.

1300  Stetson, R. H. (1951). *Motor Phonetics:A study of Speech Movements in*

1301        *Action*. North Holland.

1302  Stone, J. V. (2013). *Bayes' rule: A tutorial introduction to Bayesian analysis*.

1303        Sebtel press.

1304  Strycharczuk, P., & Kohlberger, M. (2016). Resyllabification Reconsidered:

1305        On the Durational Properties of Word-Final /s/ in Spanish. *Laboratory*

1306        *Phonology*, *7*, 1–24. https://doi.org/10.5334/labphon.5

1307  Strycharczuk, P., & Scobbie, J. M. (2017). Fronting of Southern British

1308        English high-back vowels in articulation and acoustics. *The Journal of*

1309        *the Acoustical Society of America*, *142*(1), 322–331.

1310        https://doi.org/10.1121/1.4991010

1311 Tiffany, W. R. (1980). "The effects of syllable structure on diadochokinetic

1312      and reading rates," *J. Speech Hear. Res.* **23**, 894-908.

1313 Tilsen, S. (2017). Exertive modulation of speech and articulatory phasing.

1314      *Journal of Phonetics*, *64*, 34–50.

1315      https://doi.org/10.1016/j.wocn.2017.03.001

1316 Tilsen, S. (2019). Motoric Mechanisms for the Emergence of Non-local

1317      Phonological Patterns. *Frontiers in Psychology*, *10*, 2143.

1318      https://doi.org/10.3389/fpsyg.2019.02143

1319 Tilsen, S. (2020). Detecting anticipatory information in speech with signal

1320      chopping. *Journal of Phonetics*, *82*, 100996.

1321      https://doi.org/10.1016/j.wocn.2020.100996

1322 Tilsen, S. (2022). *An informal logic of feedback-based temporal control*.

1323      https://doi.org/10.13140/RG.2.2.14017.28003/1

1324 Tilsen, S., Kim, S. E., & Wang, C. (2021). Localizing category-related

1325      information in speech with multi-scale analyses. *PloS one*, *16*(10),

1326      e0258178. https://doi.org/10.1371/journal.pone.0258178

1327 Tuller, B., & Kelso, J. A. . S. (1990). Phase transitions in speech production

1328      and their perceptual consequences. In M. Jeannerod (Ed.), *Attention*

1329      *and performance* (Vol. 13). Erlbaum.

1330 Tuller, B., & Kelso, J. A. . S. (1991). The production and perception of

1331     syllable structure. *Journal of Speech and Hearing Research*, *34*, 501–

1332     508.

1333 Uffmann, C. (2007). Intrusive [r] and optimal epenthetic consonants.

1334     *Language Sciences*, *29*(2–3), 451–476.

1335 Wagenmakers, E.-J., Morey, R. D., & Lee, M. D. (2016). Bayesian benefits

1336     for the pragmatic researcher. *Current Directions in Psychological*

1337     *Science*, 25(3), 169–176.

1338 Wagner, M. (2005). Prosody and Recursion. (Massachusetts Institute of

1339     Technology).

1340 Wang, C., Xu, Y., and Zhang, J. (2019). Mandarin and English use different

1341     temporal means to mark major prosodic boundaries. in *The 19th*

1342     *International Congress of Phonetic Sciences* (Melbourne, Australia).

1343 Wu, S.-L., Shire, M. L., Greenberg, S., & Morgan, N. (1997). Integrating

1344     syllable boundary information into speech recognition. *1997 IEEE*

1345     *International Conference on Acoustics, Speech, and Signal*

1346     *Processing*, *2*, 987–990. https://doi.org/10.1109/ICASSP.1997.596105

1347 Xu, A., Birkholz, P., & Xu, Y. (2019). *Coarticulation as synchronized*

1348     *dimension-specific sequential target approximation: An articulatory*

1349     *synthesis simulation*. 205–109.

1350 Xu, Y. (1986). Acoustic-phonetic characteristics of junctures in Mandarin

1351       Chinese. *Journal of Chinese Linguistics*, *4*, 353–360.

1352 Xu, Y. (2020). *Syllable is a synchronization mechanism that makes human*

1353       *speech possible* [Preprint]. PsyArXiv.

1354       https://doi.org/10.31234/osf.io/9v4hr

1355 Xu, Y., & Liu, F. (2006). Tonal alignment, syllable structure and

1356       coarticulation: Toward an integrated model. *Italian Journal of*

1357       *Linguistics*, *18*, 125–159.

1358 Xu, Y., & Liu, F. (2007). Determining the temporal interval of segments with

1359       the help of F0 contours. *Journal of Phonetics*, *35*(3), 398–420.

1360       https://doi.org/10.1016/j.wocn.2006.06.002

1361 Zhang, X., Sun, J., & Luo, Z. (2014). One-against-All Weighted Dynamic

1362       Time Warping for Language-Independent and Speaker-Dependent

1363       Speech Recognition in Adverse Conditions. *PLoS ONE*, *9*(2), e85458.

1364       https://doi.org/10.1371/journal.pone.0085458