# Deriving the onset and offset times of planning units from acoustic and articulatory measurements

**Joe Rodd,**[1,2,a] **Hans Rutger Bosker,**[1,3] **Louis ten Bosch,**[2,1,3] **and Mirjam Ernestus**[2,1]

[1] *Max Planck Institute for Psycholinguistics, P.O. Box 310, 6500 AH, Nijmegen, The Netherlands*

[2] *Radboud University, Centre for Language Studies, P.O. Box 9103, 6500 HD, Nijmegen, The Netherlands*

[3] *Radboud University, Donders Institute for Brain, Cognition and Behaviour, P.O. Box 9104, 6500 HE, Nijmegen, The Netherlands*

*joe.rodd@mpi.nl,*

*hansrutger.bosker@mpi.nl,*

*l.tenbosch@let.ru.nl,*

*m.ernestus@let.ru.nl*

1

1     **Abstract:**    Many psycholinguistic models of speech sequence plan-

2     ning make claims about the onset and offset times of planning units,

3     such as words, syllables, and phonemes. These predictions typically go

4     untested, however, since psycholinguists have assumed that the tem-

5     poral dynamics of the speech signal is a poor index of the temporal

6     dynamics of the underlying speech planning process. This article ar-

7     gues that this problem is tractable, and presents and validates two

8     simple metrics that derive planning unit onset and offset times from

9     the acoustic signal and articulatographic data.

[a] Author to whom correspondence should be addressed.

## 1. Background

Typically, the inverse mapping between the acoustic signal and the articulator configuration is characterized as highly non-linear and one-to-many, in that many speech sounds can be produced by multiple configurations of the vocal tract (e.g. Lindblom, 1983). This assumed intractability complicated the evaluation of psycholinguistic models of speech planning, specifically claims about the implementation of abstract linguistic planning units by speech motor programs.

While it is the case that speakers can make use of alternative vocal tract configurations to achieve speech sounds when articulatory freedom is constrained (Lindblom, 1979), or to reduce the required movement from the previous configuration (e.g. Boyce and Espy-Wilson, 1997), the opacity of the correspondences between acoustics, articulation, and the dynamics of higher planning processes may be overestimated (Hogden *et al.*, 1996). This paper posits that the problem is tractable, and proposes methods to characterize the dynamics of higher planning processes from the acoustic signal or from tracked articulator movements. Thus, the testing of previously untestable predictions of psycholinguistic models is facilitated.

### 1.1 Acoustic change largely reflects articulatory change

Despite assumptions to the contrary, in practice, the inverse mapping from the acoustic signal to articulatory configurations can be defined in a appropriate way to predict articulatory configurations from the acoustic signal, within a certain tolerance for deviations in the articulatory domain. For speech sounds that intrinsically consist of multiple acoustic events (such as diphthongs, plosives), the mapping results in an estimated trajectory in articula-

31 tory space. For a subset of stable speech sounds, 'codebooks' of articulatory configurations

32 associated with acoustic outcomes can be compiled (e.g. Hogden *et al.*, 1996). Moreover,

33 machine learning approaches that can make use of contextual information and sufficiently

34 large corpora of training data have proven successful in predicting articulatory configuration

35 from the acoustic signal with no constraints on speech materials (e.g. Illa and Ghosh, 2018;

36 Richmond, 2006; Uria *et al.*, 2011).

37     Relatedly, it holds that when the vocal tract is in a stable configuration, the acoustic

38 output is also stable, and that when the acoustic output is changing, the vocal tract con-

39 figuration must also be changing. This observation has been exploited in blind speech seg-

40 mentation, where frame-by-frame changes in the acoustic spectrum are tracked, and peaks

41 in spectral change are detected. These peaks correspond to perceptually relevant phone

42 boundaries (e.g. Dusan and Rabiner, 2006; Hoang and Wang, 2015; ten Bosch and Cranen,

43 2007). These approaches are intended to automate the preparation of corpora to test speech

44 recognition systems, and assume that segments are concatenated without overlap, making

45 these algorithms unsuited for the retrieval of onset and offset times of overlapping planning

46 units predicted by psycholinguistic models. They can, however, serve as inspiration for the

47 development of new techniques to retrieve planning unit dynamics.

48     Note that although changes in the acoustic signal must reflect changes in the articula-

49 tory configuration, it does not follow that when the vocal tract configuration is changing, the

50 acoustic signal always changes with it, since for many speech sounds, the precise positioning

51 of non-critical articulators is unimportant (such as tongue position during the realization of

52 /m/).

53 *1.2 The mapping between planning units and acoustics and articulation*

54 A class of psycholinguistic speech production models (which we will term phoneme-based

55 models) characterize the units that mediate between formulation (lexical access and phono-

56 logical encoding) and execution (speech motor programming and articulation itself) as

57 phonemes, or sequences of phonemes, such as syllables, demi-syllables, or whole words (e.g.

58 Dell and O'Seaghdha, 1992; Levelt, 1989; Levelt *et al.*, 1999; Tourville and Guenther, 2011).

59 Phoneme-based models also conceptualize the execution process as an obedient servant of

60 formulation (e.g. Dell and O'Seaghdha, 1992; Levelt, 1989; Levelt *et al.*, 1999; Tourville

61 and Guenther, 2011), which entails that the observable movements of the articulators and

62 the resulting speech acoustics are inherently a consequence of planning units in formulation

63 becoming active and subsequently being deactivated. That the dynamics of the activation

64 of planning units directly influences the articulatory configuration and thereby the acoustic

65 output seems plausible in the light of findings that competing representations in the formu-

66 lation phase exert some influence on fine detail in articulation (e.g. Goldrick and Blumstein,

67 2006).

68     The DIVA model (Tourville and Guenther, 2011) operationalizes the planning units

69 by defining them in terms of upper and lower bounds for articulator positions, and upper

70 and lower bounds of the expected auditory outcome in terms of fundamental frequency and

71 formants. Planning units typically overlap in time, and all simultaneously active planning

⁷² units exert influence on both the articulatory configuration and speech acoustics directly via

⁷³ the feedforward route. They also influence articulation and acoustics indirectly by shaping

⁷⁴ the expected acoustic and somatosensory outcomes, which in turn lead to corrective feedback.

⁷⁵ The temporal overlap of adjacent planning units (at the output stage of phoneme-

⁷⁶ based psycholinguistic speech planning models) results in local coarticulation in the overt

⁷⁷ speech. Equivalently, low level pre-activation (priming) of upcoming planning units and

⁷⁸ incomplete deactivation of preceding planning units result in longer-range coarticulation in

⁷⁹ the overt speech.

⁸⁰

⁸¹ The retrieval of planning units from articulatory measurements has previously been

⁸² attempted by Steiner and Richmond (2009), who developed an analysis-by-resynthesis ap-

⁸³ proach that reconstructs a gestural score from electromagnetic articulography (EMA) data in

⁸⁴ terms of vocalic and consonantal gestures for the VocalTractLab (VTL) synthesizer (Birkholz

⁸⁵ *et al.*, 2007). This representation differs somewhat from that inherent to phoneme-based

⁸⁶ models, in that vowel and consonants are treated as fundamentally distinct units of repre-

⁸⁷ sentation on distinct tiers of the gestural score, while phoneme-based models instead predict

⁸⁸ a chain of potentially overlapping planning units of the same class, on the same tier.

⁸⁹ Vaz *et al.* (2016) described an algorithm to retrieve underlying structure from multi-

⁹⁰ variate time series data, and tested it on vocal tract constriction distances measured from

⁹¹ real-time MRI vocal tract data. The algorithm was able to construct an inventory of ges-

⁹² tures from the data, and an activation time series for each of these gestures, which are

collectively analogous to a gestural score in the articulatory phonology (AP) framework. AP diverges from phoneme-based production models in that the planning units it supposes are not phonemes or sequences of phonemes, but rather articulatory gestures defining articulatory events, such as opening of the glottal aperture, or the creation of a labial closure (Browman and Goldstein, 1992), which cannot easily be translated into phonemes.

The direct retrieval of the timings of planning units from the acoustic signal has been attempted by Nam *et al.* (2012), again with an analysis-by-synthesis approach, and similarly rooted in the articulatory phonology (AP) framework. Their procedure involves constructing a task dynamic gestural score (encoding the speech to be produced in terms of degrees of constriction at different positions in the vocal tract) from an orthographic transcription of the speech. Then, the TADA model (Nam *et al.*, 2004; Saltzman and Munhall, 1989) is used to predict time-varying vocal tract dimensions from the gestural score, which is then synthesized to produce a speech signal. Next, dynamic time warping (DTW) is applied between the synthesized and natural speech signals. This involves stretching and compressing the synthesized speech signal in the temporal dimension, to improve the temporal alignment with the natural speech signal. The result of the DTW is a warping scale, which can then be applied to the gestural score, yielding a warped gestural score from which activation and deactivation times of individual gestures can be established.

Aside from requiring potentially difficult to acquire articulatory measurements (EMA in the case of Nam *et al.* (2012), real time MRI in the case of Vaz *et al.* (2016)), these procedures that construct multivariate gestural scores cannot readily be applied to phoneme-

114 based models of speech production, since the gestures are not consistent with or easily

115 mapped to the planning units hypothesized by phoneme-based models of lexical access and

116 multi-word processes of speech production (e.g. Bohland *et al.*, 2010; Dell and O'Seaghdha,

117 1992; Levelt, 1989; Levelt *et al.*, 1999). An additional concern is that the process leaves

118 the researcher relatively unconstrained in the construction of the gestural score for a given

119 utterance, either directly or through their parameterization of the linguistic model.

120 **2. Study aims**

121 This study aims to provide a means to estimate the onset and offset times of phoneme-based

122 planning units (such as words, syllables or phonemes) from recorded speech materials. The

123 tight temporal locking between formulation and execution processes in speech production

124 (e.g. Goldrick and Blumstein, 2006) suggests that reconstructing the activation dynamics

125 of planning units from measurements of articulator movement is feasible. That the inverse

126 mapping between acoustics and articulation is transparent enough to construct codebooks

127 describing the mapping implies that reconstructing the activation dynamics of planning units

128 from the acoustic signal should also be feasible for a constrained repertoire of speech sounds.

129 We propose two approaches to retrieve planning unit onset and offset times from

130 speech materials; from the acoustic signal, and from EMA data. We compare the outcomes

131 of the two techniques, to establish that recovering planning unit onset and offset times from

132 the acoustic signal is broadly equivalent to recovering planning unit timing from articulato-

133 graphic data.

The first metric uses fleshpoint position data gathered by electromagnetic articu-lography, and begins by deriving upper and lower bounds for each fleshpoint position for each segment from corpus data. Subsequently, a multi-dimensional, time-varying target for a multi-segmental speech sequence is constructed, the temporal parameters of which are adjusted to achieve a good fit to the observed data.

The second is a metric that exploits the acoustic signal directly with no need to record articulator motion, but constrains the speech sounds that can be evaluated. This metric depends on the claim that acoustic instability mirrors articulatory instability, which in turn reflects simultaneous activation of multiple planning units.

Neither metric is predicated on any specific theoretical treatment of speech produc-tion, aside from the assumption that planning units are phonemes or sequences of phonemes, and the parameterization of both metrics is wholly data-driven. For the experimental psy-cholinguist, a metric that can be collected from the acoustic signal alone is clearly preferable, since that reduces the burden of data collection on both researcher and participant, and makes recording of electrophysiological or other measures during speech production possible because no articulatographic data needs to be collected.

The two metrics were tested on acoustic and articulatory data for the same vowel-consonant sequences, taken from the electromagnetic articulography subset of the mngu0 corpus (Richmond *et al.*, 2011), where monophthongs transitioned into continuant conso-nants. The choice of this limited subset was driven by the need to use segments that were acoustically stable during realization, for the acoustic metric. Comparing the performance of

155 the metrics against a 'gold standard' baseline annotation of the onsets and offsets of speech

156 planning units is clearly impossible, given that any hand annotation of speech planning unit

157 onsets and offsets would inherently be largely arbitrary and noisy.

158 **3. Speech materials**

159 The EMA subset of the mngu0 corpus (Richmond *et al.*, 2011) was used, which consists

160 of TIMIT sentences read by a single male speaker of British English. EMA sensors were

161 placed on the lower and upper lips, at the tongue tip, blade and dorsum and on the lower

162 incisors (to track jaw motion). A further sensor was placed on the upper incisors to serve as a

163 reference for the others. For technical details relating to the data collection and preparation

164 see Richmond *et al.* (2011).

165 *3.1 Post processing and annotation*

166 From the 1263 sentences of the mngu0 corpus, vowel - consonant (VC) sequences of inter-

167 est were identified, where a monophthong transitioned into a continuant consonant. The

168 sequences of interest were all one of the following: /am/, /aʃ/, /av/, /ɪʃ/, /ɪv/, /im/, /iv/,

169 /ʌm/, /is/, /ʌs/, /ɒn/. n/. Note that in the context of phoneme-based speech planning mod-

170 els, where no distinction is made between planning units for different classes of phonemes,

171 there is no reason to suppose that sequences of a different composition (CVs, or CCs, for

172 instance) would behave any differently from the VCs tested here. This means that the pre-

173 dictions of phoneme-based speech planning models can effectively be tested by this reduced

174 set of sequences. This yielded 775 sequences of interest, which were identified based on the

175 forced aligned transcriptions available in the corpus. Analysis intervals from the temporal

176 center of the forced aligned vowel to the temporal center of the forced-aligned consonant

177 were defined (see Figure 1(a)). The analysis interval served as a landmark to identify the

178 planning unit transitions found; so the precision of the start and end points of the interval

179 was not critical, as long as the transition between the planning units was included.

180 In the EMA data, lateral movement was discarded, yielding articulator positions on

181 the mid-sagittal plane only. To facilitate annotation, the remaining two dimensional data

182 was rotated independently for each sensor by means of principal component analysis, so

183 that PC1 captured the most informative direction of movement for that sensor, which in all

184 cases was the open-close dimension. Since PC2 is orthogonal to PC1, it captured forward-

185 backward movement of each sensor. Then, manual annotation was undertaken (by the first

186 author) to identify articulatory stable periods of each segment for use in the preparation of

187 the targets used in the articulatory metric. In the manual annotation procedure, movement

188 tracks in PC1-PC2 dimensions were displayed on a graphical interface, in which the periods

189 of stability associated with the vowel and continuant consonant could be highlighted. The

190 articulatory configuration was considered stable if there was little to no change (assessed

191 visually) in several sensors. Since the targets were defined in terms of 95% highest density

192 intervals (see section 4.1), some noise in this annotation procedure was acceptable.

193 **4. Planning unit timing from articulatory measurements**

194 The articulatory metric approaches the identification of planning unit onset and offset times

195 from EMA data by essentially inverting the motor control process: reconstructing a mul-

196 tidimensional articulatory target that could have lead to the recorded movements during
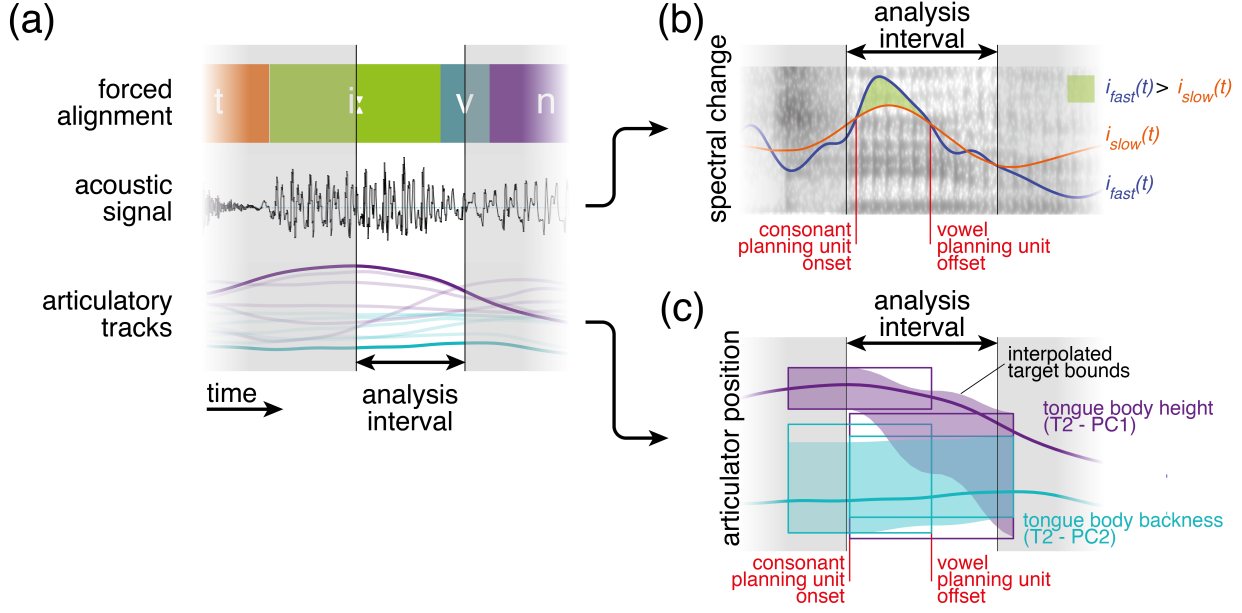
Fig. 1.    (color online) An example analysis. (a) An analysis interval is defined that stretches

from the temporal center of the forced aligned vowel to the center of the forced aligned consonant.

(b) The acoustic metric. Lines show $I_{slow}(t)$ and $I_{fast}(t)$, the Gaussian smoothed, interpolated

spectral distance functions used to identify the acoustically evident planning unit overlap during

the analysis interval. $I_{slow}(t)$ has a 90 ms kernel, $I_{fast}(t)$ has a 30 ms kernel. Shading identifies

periods of atypically fast acoustic change (where $I_{fast}(t) > I_{slow}(t)$), from which the onset of the

consonant planning unit and the offset of the vowel planning unit are derived. (c) The articulatory

metric. The heavy lines indicate the recorded movement of the tongue body sensor, in the open-

close dimension (PC1) and the forward-backward dimension (PC2). The outlined boxes indicate

the segmental targets, the shading indicates the interpolated sequence level target.

197    a vowel-consonant sequence. This was done separately for each vowel-consonant transition

198    token, using a parameter optimization routine which adjusted the onset and offset times of

199    the segment targets to construct a target that fitted the recorded movements well.

200    *4.1 Establishing segment targets*

201    First, separate segmental targets are established for the vowels and for the consonants,

202    defined in terms of upper and lower bounds for the positions for each fleshpoint (lower

203    jaw, upper and lower lips, tongue tip, blade and dorsum) on the two dimensions (principal

204    components) of the mid-sagittal plane. These maxima and minima are derived from the

205    distribution of sensor positions during the hand-annotated stable periods of those segments

206    in the corpus, irrespective of context, by extracting the 95% highest density interval(s).

207    When the positioning of a fleshpoint is of crucial importance to the identity of the segment,

208    the positioning of that fleshpoint varies little between realizations, and the target is therefore

209    narrow (e.g. the positioning of the tongue tip in /s/). When the positioning of a fleshpoint

210    is only marginally relevant for the identity of the segment, the target is broad (e.g. the

211    positioning of the tongue back in /v/), since there is lots of variability in the source data.

212    *4.2 Combining segmental targets to form a sequence target*

213    The sequence targets were constructed by temporally-overlapping the vowel and consonant

214    targets. Figure 1(c) depicts an example of the construction of the targets, for the sequence

215    /iːv/, showing the target bounds for each segment as boxes (purple for PC1, blue for PC2),

216    for the tongue body sensor. The segmental targets are fixed at the outer edges, such that

217    the vowel target begins at the hand-annotated onset of vowel stability, and the consonant

²¹⁸ target ends at the hand-annotated offset of consonant stability. The other two temporal

²¹⁹ parameters, the offset of the vowel target and the onset of the consonant target are free

²²⁰ parameters that can be optimized.

²²¹ The upper bound of the sequence target is calculated as an exponential moving aver-

²²² age (with a window of 20 ms) of the upper bounds of the segmental targets over time. This

²²³ means that for time points when only the vowel target is engaged, the upper bound is equal

²²⁴ to the upper bound of the vowel target. When both segmental targets are engaged, however,

²²⁵ the upper bound switches smoothly from following the upper bound of the vowel target to

²²⁶ reflecting the average upper bound of both targets. Once the vowel target is disengaged, the

²²⁷ upper bound again smoothly shifts to reflect the upper bound of the consonant target. The

²²⁸ lower bound of the target is derived in the same way.

²²⁹ *4.3 Parameter optimization*

²³⁰ For each analysis interval, an independent parameter optimization routine is conducted. Two

²³¹ parameters, the onset time of the consonant target and the offset time of the vowel target,

²³² are optimized with the BOBYQA algorithm (Powell, 2009; Ypma *et al.*, 2018).

²³³ To evaluate how well a sequence target defined by a pair of consonant target onset

²³⁴ and vowel target offset times fitted the observed movements, the proportion of time points

²³⁵ where the recorded sensor positions are outside the bounds of the multidimensional tar-

²³⁶ get is counted. This proportion is used as a score to be minimized during the parameter

²³⁷ optimization process.

238  For each realization, 200 starting points for these parameters are tried, sampled from

239  normal distributions ($SD = 25$ ms) centered around the annotated end of vowel stability

240  (this is the center-point of the starting distributions for the consonant onset parameter) and

241  the annotated beginning of consonant stability (this is the center-point of the starting distri-

242  butions for the vowel offset parameter). A search space constraint ensures that the algorithm

243  only considers solutions where the overlap between the segment targets is greater than 0.

244  Having multiple starting points allows us to assess how consistently the algorithm selects

245  the best performing parameter sets, and offers more protection from premature convergence

246  to local minima. To select a single vowel offset time and a single consonant onset time from

247  the distributions that resulted from the 200 initializations, a two-dimensional distribution

248  is estimated from the resulting parameters, where the dimensions are the vowel offset time

249  parameter and consonant onset time parameter. The distribution is weighted by one minus

250  the score achieved in each attempt, so as to weight the best performing solutions most heav-

251  ily, and the peak is identified. The coordinates of this peak define the planning unit onset

252  and offset times.

253  **5. Planning unit timing from the acoustic signal**

254  The acoustic metric quantifies the rate of change in the acoustic signal (the spectral change).

255  Local peaks in this signal identify periods where the speech acoustics, and therefore the

256  underlying vocal tract configuration, are changing. At the transition between two planning

257  units, this change is due to the interaction of the two overlapping planning units, and the

258  duration of the instability is equated with the duration of the overlap. We term this overlap

259  'acoustically evident planning unit overlap'. To be able to establish the onset and offset

260  of instability, a method is required to transform a continuous signal into a categorical one:

261  to distinguish acoustic stability from instability. This is done by overlaying two different

262  smoothings of this signal; a 'fast' smooth that captures local changes in the signal, and

263  a 'slow' smooth that captures longer trends. We identify periods when the 'fast' smooth

264  exceeds the 'slow' smooth as unstable, and other periods as stable. The onset of the second

265  planning unit is equated with the start of such a period of instability. The offset of the first

266  planning unit is equated with the end of that same period of instability. This is illustrated

267  in Figure 1(b).

268  *5.1 Step 1: quantifying acoustic change*

269  To identify the period of overlap, the MFCC vectors (mel frequency cepstral coefficient; 25

270  ms analysis frame length, samples every 10 ms) for the analysis intervals (with a margin of

271  40 ms before and after) are extracted using the HTK front end (Young *et al.*, 2006). MFCC

272  vectors may be seen as a numeric representation of the spectral content of the speech signal

273  during a short (25 ms) window, and are one of the best spectro-temporal representations of

274  speech acoustics. From each frame to the next, the Euclidean distance in MFCC space was

275  calculated as follows, where $j$ is the index of the MFCC coefficient and $t$ is the index of the

276  frame:

$$D_{spec} = \sqrt{\sum_{j=0}^{12} (\text{MFCC}_{j_t} - \text{MFCC}_{j_{t+1}})^2} \tag{1}$$

277    This gives $D_{spec}(t)$, a spectral distance function quantifying the degree of spectral

278    change evident in the acoustic signal, sampled every 10 ms.

279    *5.2 Step 2: identifying periods of fast acoustic change*

280    This spectral distance function is smoothed twice, once with a 30 ms wide Gaussian kernel,

281    yielding $D_{fast}(t)$, which captures relatively fast changes in the spectral distance function;

282    and once with a 90 ms wide Gaussian kernel, yielding $D_{slow}(t)$, which captures longer term

283    trends in the function.

284    Spline interpolation (every 0.1 ms) is then applied to these functions in order to

285    improve temporal resolution, yielding $I_{fast}(t)$ and $I_{slow}(t)$ . The two interpolated functions

286    are overlaid, and parts of the signal in each analysis interval where $I_{fast}(t)$ is larger than

287    $I_{slow}(t)$ are identified as candidate overlaps (in Figure 1(b) shown as green shading). Where

288    $I_{fast}(t)$ exceeds $I_{slow}(t)$, atypically fast acoustic change is occurring: acoustically evident

289    planning unit overlap. It is possible that there are multiple periods where $I_{fast}(t)$ exceeds

290    $I_{slow}(t)$, however, typically one period is longer and the associated peak is larger. Therefore,

291    a heuristic is engaged to select precisely one period per analysis window: the duration of each

292    of these periods is calculated. Periods that cross the boundaries of the analysis interval (into

293    the margins) are discarded. When an analysis interval still contains multiple periods, all but

294    the longest candidate are discarded. This yields precisely one period of acoustically evident

295    planning unit overlap per analysis interval. The onset of the remaining period of overlap

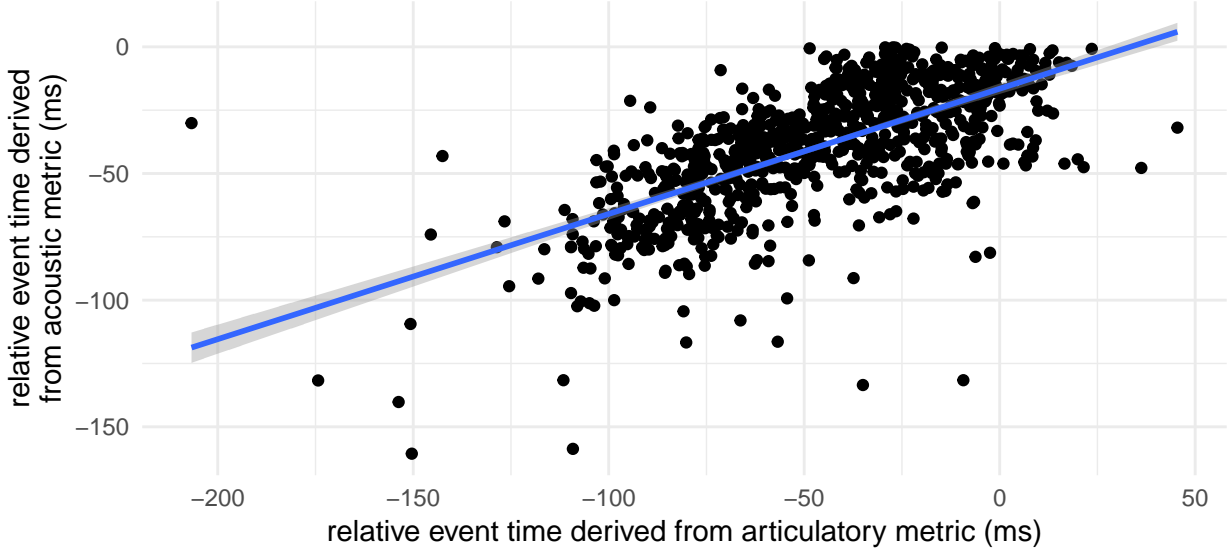296    (where $I_{fast}(t)$ becomes larger than $I_{slow}(t)$) yields the onset of the consonant planning unit.

Fig. 2. (color online) The correlation between planning unit onset and offset times, derived from the articulatory (x-axis) and acoustic metrics (y-axis). All event times are relative to the forced-aligned offset of the consonant segment, meaning that times less than 0 are to be expected.

297 The offset of the overlap (where $I_{fast}(t)$ becomes smaller than $I_{slow}(t)$) yields the offset of

298 the vowel planning unit.

299 This procedure was refined by testing various kernel widths and interpolations via a

300 grid search, in which the the parameters that resulted in the highest spectral change peak

301 were selected.

302 R scripts implementing the two metrics and the data preprocessing method are avail-

303 able from https://git.io/fh8EM.

304 **6. Results and discussion**

305 *6.1 Validity of the metrics*

Figure 2 shows the onsets and offsets of planning units (event times) as predicted by the articulatory (x-axis) and acoustic metrics (y-axis). All event times are relative to the forced-aligned offset of the consonant segment, meaning that times less than 0 are to be expected. An $r^2$ of 0.447 was calculated between the event times derived by the two metrics. This moderately high correlation between the predictions of the two metrics indicates that they both capture the same underlying dynamic process of planning unit activation.

The intercept of -10.64 indicates that the acoustic metric systematically predicts earlier event times than the articulatory metric does. This is approximately half the width of the 25 ms analysis window employed in the acoustic metric, which suggests that this anticipation may be an artifact of the spectral analysis inherent to the acoustic metric.

*6.2 Reliability*

The metrics were evaluated by comparing the planning unit onset and offset times predicted by each metric. Because the two metrics are so divergent in the modality of the data used and the approach used to derive event times from the data, we interpreted the finding that the two metrics predicted comparable event times as evidence that they are both indexing the onset and offset times of planning units. This is of course weaker evidence in support of the validity of a metric than comparison against data capturing the ground truth, but the ground truth is clearly unobtainable for psychological processes such as the activation dynamics of planning units. Comparison against the results obtained by Nam *et al.* (2012) is also problematic given the AP theoretical framing inherent to their procedure.

*6.3 Applicability and ecological validity*

The articulatory metric is in principle equally suited to examining transitions between any pair of segments where there is at least a short period of articulatory stability in each segment, including stops. Of course, given the metric-comparison approach we took to evaluate the performance of the two metrics, the articulatory metric was only tested on materials also suitable for the acoustic metric.

The acoustic metric is inherently limited to identifying planning unit onset and offset times at transitions between a subset of segment types involving at least a short period of articulatory stability and incomplete obstruction of the airflow: monophthong vowels, nasals and continuant fricatives. Nevertheless, for the experimental psycholinguist, the convenience of the acoustic-only recording may well outweigh the disadvantage of constrained material selection.

Both metrics share the inherent assumption that the onsets of all the movements or gestures involved in the production of a phoneme are synchronized. This assumption is inherent to the class of phoneme-based models, which form the mainstream in psycholinguistic models of higher speech planning. Adhering to it was necessary to achieve this paper's goal of making it possible to test and refine phoneme-based models by relating activation dynamics to the speech signal. Models based on a multivariate gestural score may achieve better fits to the data given that they are not constrained by this synchronicity assumption.

The metrics were developed and tested using the mngu0 corpus (Richmond *et al.*, 2011), which contains a large quantity of English data from a single speaker, rather than

348 smaller quantities of data from multiple speakers available in other corpora (e.g. the Wis-

349 consin x-ray microbeam database, Westbury *et al.*, 1990). The mngu0 corpus was selected

350 because we sought to have a large number of realizations of each segment to reliably compute

351 the static segment targets for the articulatory metric. It remains to be seen how the articu-

352 latory metric would perform given a smaller dataset from which to derive target boundaries.

353 A requirement for a large speaker-specific dataset would be disadvantageous in the context

354 of experimental psycholinguistics, where it is typically desirable to test multiple speakers

355 on a small set of materials, though recent success in using a generalized background model

356 and a speaker-specific adaptive model in acoustic-to-articulatory inversion (Illa and Ghosh,

357 2018) offers hope that a comparable approach could work for this metric too.

## 7. Conclusion

359 This paper presented two techniques to identify planning unit onsets and offsets from artic-

360 ulographic and acoustic data in the context of phoneme-based models of speech production.

361 The first metric requires articulographic recording, but imposes less constraint on speech

362 material selection. The second metric exploits the acoustic signal directly, with no need to

363 record articulator motion, but constrains the speech sounds that can be evaluated. This

364 metric depends on the claim that acoustic instability mirrors articulatory instability, which

365 in turn reflects simultaneous activation of multiple planning units. The two metrics are ag-

366 nostic to the duration of planning units (syllables, demi-syllables, phonemes, entire words),

367 and make minimal assumptions about precisely what is encoded by the planning unit, other

368 than that upper and lower bounds for articulatory positions are encoded. A moderately

high correlation between the event times predicted by the two metrics indicates that they

capture the same underlying dynamic process of planning unit activation. This correlation

means in turn that temporal predictions arising from phoneme-based psycholinguistic mod-

els of speech planning can be tested using the acoustic signal without the need to collect

articulographic data.

## Acknowledgments

## References and links

Birkholz, P., Steiner, I., and Breuer, S. (**2007**). "Control concepts for articulatory speech
synthesis," in *Sixth ISCA Tutorial and Research Workshop on Speech Synthesis*, edited by
P. Wagner, J. Abresch, S. Breuer, and W. Hess, Bonn.

Bohland, J. W., Bullock, D., and Guenther, F. H. (**2010**). "Neural representations and mech-
anisms for the performance of simple speech sequences," Journal of cognitive neuroscience
**22**(7), 1504–1529.

Boyce, S., and Espy-Wilson, C. Y. (**1997**). "Coarticulatory stability in American English
/r/," The Journal of the Acoustical Society of America **101**(6), 3741–3753, https://asa.

scitation.org/doi/abs/10.1121/1.418333, doi: 10.1121/1.418333.

Browman, C. P., and Goldstein, L. (**1992**). "Articulatory phonology: An overview," Phonetica **49**(3-4), 155–180.

Dell, G. S., and O'Seaghdha, P. G. (**1992**). "Stages of lexical access in language production," Cognition **42**(1–3), 287–314, http://www.sciencedirect.com/science/article/pii/001002779290046K, doi: 10.1016/0010-0277(92)90046-K.

Dusan, S., and Rabiner, L. (**2006**). "On the relation between maximum spectral transition positions and phone boundaries," in *Ninth International Conference on Spoken Language Processing*.

Goldrick, M., and Blumstein, S. E. (**2006**). "Cascading activation from phonological planning to articulatory processes: Evidence from tongue twisters," Language and Cognitive Processes **21**(6), 649–683, doi: 10.1080/01690960500181332.

Hoang, D.-T., and Wang, H.-C. (**2015**). "Blind phone segmentation based on spectral change detection using Legendre polynomial approximation," The Journal of the Acoustical Society of America **137**(2), 797–805.

Hogden, J., Lofqvist, A., Gracco, V., Zlokarnik, I., Rubin, P., and Saltzman, E. (**1996**). "Accurate recovery of articulator positions from acoustics: New conclusions based on human data," The Journal of the Acoustical Society of America **100**(3), 1819–1834.

Illa, A., and Ghosh, P. K. (**2018**). "Low resource acoustic-to-articulatory inversion using bi-directional long short term memory," Proc. Interspeech 2018 3122–3126.

Levelt, W. J. M. (**1989**). *Speaking: From intention to articulation* (MIT press, Cambridge, MA).

Levelt, W. J. M., Roelofs, A., and Meyer, A. S. (**1999**). "A theory of lexical access in speech production," Behavioral and brain sciences **22**(01), 1–38.

Lindblom, B. (**1979**). "Formant frequencies of some fixed-mandible vowel and a model of speech motor programming by predictive simulation," Journal of Phonetics **7**, 147–162.

Lindblom, B. (**1983**). "Economy of Speech Gestures," in *The Production of Speech*, edited by P. F. MacNeilage (Springer New York), pp. 217–245, http://link.springer.com/chapter/10.1007/978-1-4613-8202-7_10, doi: 10.1007/978-1-4613-8202-7_10.

Nam, H., Goldstein, L., Saltzman, E., and Byrd, D. (**2004**). "TADA: An enhanced, portable Task Dynamics model in MATLAB," The Journal of the Acoustical Society of America **115**(5), 2430–2430.

Nam, H., Mitra, V., Tiede, M., Hasegawa-Johnson, M., Espy-Wilson, C., Saltzman, E., and Goldstein, L. (**2012**). "A procedure for estimating gestural scores from speech acoustics," The Journal of the Acoustical Society of America **132**(6), 3980–3989.

Powell, M. J. (**2009**). "The BOBYQA algorithm for bound constrained optimization without derivatives," Cambridge NA Report NA2009/06, University of Cambridge, Cambridge 26–46.

Richmond, K. (**2006**). "A trajectory mixture density network for the acoustic-articulatory inversion mapping," in *Ninth International Conference on Spoken Language Processing*.

Richmond, K., Hoole, P., and King, S. (**2011**). "Announcing the electromagnetic articulography (day 1) subset of the mngu0 articulatory corpus," in *Twelfth Annual Conference of the International Speech Communication Association*.

Saltzman, E. L., and Munhall, K. G. (**1989**). "A dynamical approach to gestural patterning in speech production," Ecological psychology **1**(4), 333–382.

Steiner, I., and Richmond, K. (**2009**). "Towards unsupervised articulatory resynthesis of German utterances using EMA data," in *Tenth Annual Conference of the International Speech Communication Association*.

ten Bosch, L. F. M., and Cranen, B. (**2007**). "A computational model for unsupervised word discovery," in *Proc. Interspeech 2007*, pp. 1481–1884.

Tourville, J. A., and Guenther, F. H. (**2011**). "The DIVA model: A neural theory of speech acquisition and production," Language and Cognitive Processes **26**(7), 952–981, http://dx.doi.org/10.1080/01690960903498424, doi: 10.1080/01690960903498424.

Uria, B., Renals, S., and Richmond, K. (**2011**). "A deep neural network for acoustic-articulatory speech inversion," in *NIPS 2011 Workshop on Deep Learning and Unsupervised Feature Learning*, Citeseer.

Vaz, C., Toutios, A., and Narayanan, S. S. (**2016**). "Convex Hull Convolutive Non-Negative Matrix Factorization for Uncovering Temporal Patterns in Multivariate Time-Series Data.," in *INTERSPEECH*, pp. 963–967.

Westbury, J., Milenkovic, P., Weismer, G., and Kent, R. (**1990**). "X?ray microbeam speech production database," The Journal of the Acoustical Society of America **88**(S1), S56–S56,

449  https://asa.scitation.org/doi/abs/10.1121/1.2029064, doi: 10.1121/1.2029064.

450  Young, S., Evermann, G., Gales, M., Hain, T., Kershaw, D., Liu, X., Moore, G., Odell, J.,

451  Ollason, D., Povey, D. *et al.* (**2006**). "The HTK book (for HTK version 3.4)," Cambridge

452  University Engineering Department **2**(2).

453  Ypma, J., Borchers, H. W., and Eddelbuettel, D. (**2018**). "Package 'nloptr'" R package.