# Word Segmentation on Discovered Phone Units with Dynamic Programming and Self-Supervised Scoring

Herman Kamper

*Abstract*—Recent work on unsupervised speech segmentation has used self-supervised models with phone and word segmentation modules that are trained jointly. This paper instead revisits an older approach to word segmentation: bottom-up phone-like unit discovery is performed first, and symbolic word segmentation is then performed on top of the discovered units (without influencing the lower level). To do this, I propose a new unit discovery model, a new symbolic word segmentation model, and then chain the two models to segment speech. Both models use dynamic programming to minimize segment costs from a self-supervised network with an additional duration penalty that encourages longer units. Concretely, for acoustic unit discovery, duration-penalized dynamic programming (DPDP) is used with a contrastive predictive coding model as the scoring network. For word segmentation, DPDP is applied with an autoencoding recurrent neural as the scoring network. The two models are chained in order to segment speech. This approach gives comparable word segmentation results to state-of-the-art joint self-supervised segmentation models on an English benchmark. On French, Mandarin, German and Wolof data, it outperforms previous systems on the ZeroSpeech benchmarks. Analysis shows that the chained DPDP system segments shorter filler words well, but longer words might require some external top-down signal.

*Index Terms*—Unsupervised word segmentation, phone segmentation, acoustic unit discovery, zero-resource speech.

## I. INTRODUCTION

THE Zero-Resource Speech Team at the 2012 JHU CLSLP workshop had a simple goal: they wanted to develop an unsupervised approach that takes unlabelled speech and provides a full segmentation of the audio into word-like units [1]. This could enable speech technology in very low-resource settings, and could be used in cognitive models that attempt to mimic how infants acquire language. It seemed doable at the time because there were several models that could perform full-coverage word segmentation from transcribed symbolic input (phoneme or phone sequences) [2]–[4]. The team's idea was to use a clustering model to map speech to phone-like acoustic units—a Gaussian mixture model (GMM) was used—and to then apply symbolic segmentation on top of the discovered units. Unfortunately, they found that this chained approach—where bottom-up acoustic unit discovery is performed independently of word segmentation—gave very poor results, with word token $F_1$ scores of around 5%. Compared to idealized symbolic input, the discovered units were too noisy and variable for a symbolic segmentation approach to pick up any meaningful signal [1].

E&E Engineering, Stellenbosch University, South Africa.

In the following decade, several groups moved to joint speech segmentation models where bottom-up information from discovered acoustic units can influence top-down word segmentation and vice versa. The hierarchical Bayesian HMM [5] and the later nonparametric Bayesian double-articulation analysers [6]–[8] are some examples. Other groups tried to model higher-level units like syllables or words directly, circumventing the need for explicit phone modelling [9]–[12]. E.g., the Bayesian embedded segmental GMM (BES-GMM) [10] performs probabilistic clustering and segmentation on fixed-dimensional representations of whole word-like segments, as does the related embedded segmental $K$-means (ES-KMeans) model [11]. Much more recently, self-supervised neural networks using contrastive predictive coding (CPC) [13] have been modified to include separate acoustic unit and word segmentation modules that are trained jointly [14], [15]. These joint models give state-of-the-art word segmentation results (details in §II).

Given that direct whole-word modelling [10], [11] and joint phone and word modelling [14], [15] have proven to be such fruitful directions for speech segmentation, maybe the original premise of [1] was flawed? This paper argues that their idea— where bottom-up acoustic unit discovery and word segmentation are performed separately—should be revisited. Support for this comes from the big recent improvements in unsupervised unit discovery itself, specifically through self-supervised models that are coupled with a clustering step [16], [17]. These improved acoustic units (learned in a purely bottom-up fashion) could provide the signal necessary for symbolic word segmentation.

In this paper I specifically describe a duration-penalized dynamic programming (DPDP) procedure that combines a segment scoring function with a duration penalty. Duration modelling is not a new idea in unsupervised word segmentation [6]–[8], [10], [18], [19], but coupling it with self-supervised neural scoring functions is. The proposed DPDP approach can be used for either acoustic unit discovery or symbolic word segmentation, given an appropriate scoring network. I describe separate DPDP models for these two tasks, and then chain the two models in order to do speech segmentation. For acoustic unit discovery, I apply DPDP with a CPC clustering model, an approach based on [20], [21]. For symbolic word segmentation, DPDP is applied with an autoencoding recurrent neural network (AE-RNN), an approach inspired by [22]. Chaining these two new DPDP models gives similar performance to the joint self-supervised segmental [14] and aligned [15] CPC models on English data. The DPDP system also achieves some of the best-reported word segmentation scores on the French,

Mandarin, German and Wolof ZeroSpeech 2017 and 2020 benchmarks [17], [23], where direct whole-word models have performed particularly well in the past [11].

This work can be seen as an extension of the conference paper [21], where some of the phone segmentation models in §IV were first proposed and evaluated on English. Here I improve on these models by using better self-supervised scoring networks. The main contributions of the paper are therefore as follows. (1) I propose a new model for acoustic unit discovery. (2) I propose a new approach for symbolic word segmentation. (3) These two DPDP-based models are then chained to segment speech. The chained system is applied to English and non-English data and compared to previous self-supervised and other state-of-the-art approaches. (4) I analyse the combination of different acoustic unit discovery and symbolic segmentation models (both DPDP and non-DPDP methods) for segmenting speech. Taken together, the paper shows—for the first time—that using symbolic word segmentation on top of bottom-up discovered units can give comparable or better results compared to state-of-the-art joint self-supervised and direct whole-unit speech segmentation approaches.

The paper is structured as follows. In §III, the general DPDP framework is introduced. The new acoustic unit discovery and symbolic word segmentation approaches are then respectively presented in §IV and §V. These two sections have the same structure: I introduce the model, review related work, and give an intermediate evaluation on the model's specific task. In §VII the two DPDP models are then chained, compared to previous speech segmentation approaches, and analysed. As stated, I describe related work as it becomes relevant in the respective sections. But it is worth starting by reviewing work that is relevant to the paper as a whole.

## II. RELATED WORK

I briefly outline two specific approaches that I compare to throughout. Both can be seen as extensions of contrastive predictive coding (CPC) and both can be used for either phone (§IV-C) or word segmentation from speech (§VII-A).

Standard CPC speech models learn continuous features by trying to classify future observations (at different time-steps) from among a set of negative examples [13]. The idea is that a model would need to learn meaningful phonetic contrasts while being invariant to nuisance factors such as speaker. Bhati et al. [14] extend this by using a second segment-level CPC layer. The segmental CPC (SCPC) consists of a frame-level CPC module, a differentiable boundary detector operating on the learned features, and a segment-level CPC module operating on aggregated features from the lower layer. The model is trained end-to-end using the combination of two losses: a next-frame classification loss for the lower-level CPC and a next-segment classification loss for the higher-level CPC. For phone segmentation, peak detection is used to find points of high dissimilarity between the learned features. For word segmentation, latent segment representations are compared.

Instead of treating classification independently for each future time-step as in standard CPC, the aligned CPC (ACPC) model of Chorowski et al. [24] outputs a sequence of predictions

that are then aligned to future time-steps. Since the model encourages piece-wise constant latent features, the idea is that changes in these features would correspond to phone boundaries. To extend this model to word segmentation, the learned features are used for differentiable boundary detection and passed on to a higher-level ACPC module—very similar to the SCPC. This multi-level ACPC (mACPC) is trained end-to-end on the combination of the lower- and higher-level ACPC losses [15]. The mACPC gives some of the best unsupervised word segmentation results on the Buckeye benchmark (§VII-A).

Since both the SCPC and mACPC are trained end-to-end, phone discovery at the lower level can influence word discovery at the higher level and vice versa. This is in contrast to the approach I propose, where phone discovery is performed without any influence from a word segmentation module—acoustic unit discovery is purely bottom-up. My argument is not necessarily that the bottom-up approach is superior, but we will see that it is competitive with these CPC-based models.

## III. DURATION-PENALIZED DYNAMIC PROGRAMMING (DPDP)

In this section I describe duration-penalized dynamic programming (DPDP) in its general form. In the next sections I describe specific instances of DPDP-based models that can be used for acoustic unit discovery or symbolic word segmentation.

We have an input sequence $X = x_{1:T} = (x_1, \ldots, x_T)$ that we want to segment. As illustrated in Figure 1, this can either be a sequence of symbols, e.g. each $x_t \in \{1, \ldots, K\}$, or a sequence of speech features, e.g. $x_t \in \mathbb{R}^D$. Now imagine we have some scoring network $w(\cdot)$ that takes a subsequence and gives a segment cost. If $w(x_{a:b})$ is low, this indicates that the segment $x_{a:b}$ is modelled well by the scoring network. Our goal is then to find the segmentation $S$ that gives the lowest overall cost when we sum up all the individual segment costs:

$$\arg\min_S \sum_{(a,b) \in S} w(x_{a:b}) \qquad (1)$$

We can solve this problem efficiently using dynamic programming. This corresponds to finding the shortest path through a directed acyclic graph (DAG) with edge weights given by the scoring network. A segmentation $S$ can be specified as a sequence of (start, end) tuples. The red-solid path in Figure 1 corresponds
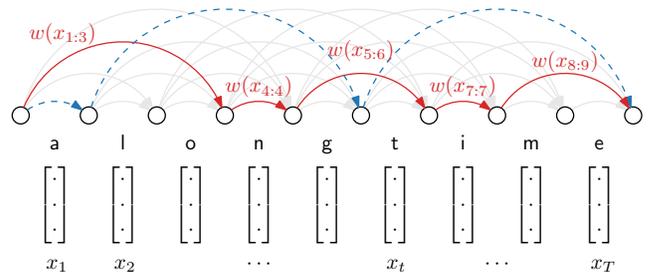


Fig. 1. The directed acyclic graph for DPDP with a maximum segment length of four. The segment cost $w(\cdot)$ should indicate how well that segment is modelled by a scoring network. E.g., on the dashed-blue path, $w(x_{2:5}) = w(\text{"long"})$ should be low if "long" is modelled well. The input can either be symbolic (shown as characters) or continuous (vectors), depending on the scoring network. Word segmentation is performed by finding the shortest path through the graph.

to the segmentation $S = ((1,3),(4,4),(5,6),(7,7),(8,9))$, giving the result "alo n gt i me". The blue-dashed path is a different segmentation with a different cost. The problem is to find the cheapest overall segmentation or, equivalently, the shortest path through the DAG.

To find this path using dynamic programming, we can define

$$\alpha_t \triangleq \min_{S_t} \sum_{(a,b) \in S_t} w(x_{a:b})$$

as the cost for the optimal segmentation up to step $t$, where $S_t$ denotes a segmentation up to this intermediate point. These forward variables can be calculated recursively:

$$\alpha_t = \min_{j=0}^{t-1} \{\alpha_j + w(x_{j+1:t})\} \qquad (2)$$

We start with $\alpha_0 = 0$ and calculate $\alpha_t$ for $t = 1, 2, \ldots, T$. We keep track of the optimal choice ($\arg\min$) for each $\alpha_t$ in (2), and the overall optimal segmentation is then obtained by starting from the final step $t = T$ and moving backwards, repeatedly choosing the optimal boundary.

I have not explained the structure of the scoring network and will only do so in the next sections. But it is worth briefly touching on one potential issue. Depending on the properties of the scoring network, we could end up with a trivial solution. One example is where the network always gives a very low cost when the input is a single symbol or single speech frame; this results in over-segmentation. To deal with this, a penalty term is added to penalize shorter segments:

$$\begin{aligned} w(x_{a:b}) &= w_{\text{seg}}(x_{a:b}) + \lambda\, w_{\text{dur}}(\text{dur}(x_{a:b})) \\ &= w_{\text{seg}}(x_{a:b}) + \lambda\, w_{\text{dur}}(b - a + 1) \end{aligned} \qquad (3)$$

Here I explicitly distinguish the segment cost $w_{\text{seg}}(\cdot)$ from the duration penalty $w_{\text{dur}}(\cdot)$ and include a duration weight $\lambda$ that controls the relative importance of the two terms.

I refer to this procedure as duration-penalized dynamic programming (DPDP). Note that nothing in this description is new: it is really just a specific formulation of dynamic programming (see e.g. [18], [19] for related duration-based models). But the formulation here gives a common way to talk about and mathematically formulate the following two new models that are contributions of this work: one is used for acoustic unit discovery (§IV) and the other for symbolic word segmentation (§V). In the final experiments (§VII), I chain the two models to do word segmentation from speech.

## IV. DPDP FOR ACOUSTIC UNIT DISCOVERY

Given a collection of unlabelled speech utterances, the goal of acoustic unit discovery is to learn a finite set of phone-like units representing the speech sounds that make up the language. A model that can do this can also be used for unsupervised phone segmentation, predicting phone boundaries. Several acoustic unit discovery and unsupervised phone segmentation models have been proposed [25]–[30]. Recently, large gains have been achieved by combining self-supervised neural networks with a clustering component. Self-supervised speech models can produce continuous features that accurately capture phonetic contrasts while being invariant

to nuisance factors (e.g. speaker) [13], [31]–[33]. To obtain discrete units, a self-supervised model can be coupled with a vector quantization (VQ) module, either as part of the model itself or by introducing a clustering step after training [32]–[37]. Despite improvements in phone discrimination tasks [17], these approaches still encode speech at a much higher bitrate than true phone sequences [38]. The main reason is that the assignment of features to VQ codebook vectors is done independently for each speech frame at a fixed rate. In reality, adjacent frames are likely to belong to the same speech unit, and these units occur at a variable rate. The SCPC and ACPC (§II) try to address this problem through intermediate boundary detection. Below I propose a different approach.

### A. DPDP on vector-quantized self-supervised speech features

The approach here constrains a VQ model so that contiguous feature vectors are assigned to the same code, resulting in a variable-rate encoding of the input speech. I first explain the approach by itself before formulating it as an instance of DPDP.

Given input speech, continuous feature vectors are first extracted using an encoder from a self-supervised speech model. Let us denote these features as $x_{1:T}$, where each $x_t \in \mathbb{R}^D$ is a learned $D$-dimensional feature vector. We then solve a constrained optimization problem through which this sequence is divided into segments. Figure 2 illustrates this, showing two possible segmentations. All the features within a segment are assigned to the same code from a VQ codebook $\{e_k\}_{k=1}^K$, with each $e_k \in \mathbb{R}^D$. The overall cost for a particular segmentation is the sum of the squared distances of the features and the representative code of each segment. In Figure 2, this cost corresponds to adding up the squared lengths of the arrowed lines. Our objective is to find the segmentation $S$ that minimizes the overall summed distances:

$$\arg\min_S \sum_{(a,b) \in S} \sum_{x_t \in x_{a:b}} ||x_t - \hat{x}_t||^2 \qquad (4)$$

where $\hat{x}_t$ is the codebook vector to which the segment $x_{a:b}$ is assigned, and the segmentation $S$ is specified as in §III. One problem with this overall cost is that the best segmentation will always place each $x_t$ in its own segment, assigning it to the
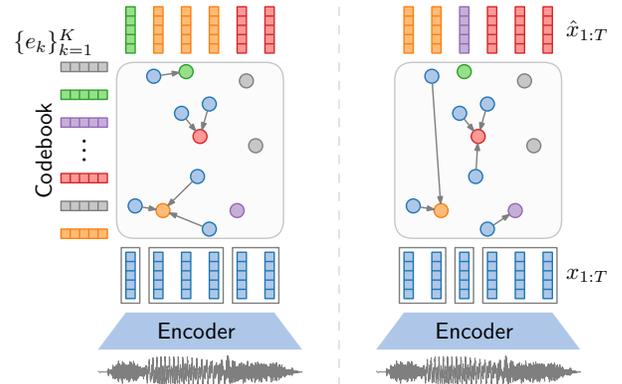


Fig. 2. Two segmentations of an utterance. The features in each segment are assigned to the same code. The left will have a smaller sum of squared distances between the features and their assigned codes than the right.

code closest to it. Additional constraints are therefore required: I specifically introduce a duration penalty to encourage longer but fewer segments. To incorporate this penalty, let us now frame this as an instance of DPDP (§III).

To get to the overall cost in (4) (without the duration penalty), we need the following segment cost:

$$w_{\mathrm{seg}}(x_{a:b}) = \min_{k=1}^{K} \sum_{x_t \in x_{a:b}} ||x_t - e_k||^2$$

We then add a duration penalty $w_{\mathrm{dur}}(l)$. There are several options. One is to use a probabilistic prior over acoustic unit duration (more on this in §V-A). Here I use a simpler linear penalty, based on [39]: $w_{\mathrm{dur}}(l) = -l + 1$. Following the DPDP formulation, the combined weight is then as in (3), and the best segmentation in (1) is obtained by recursively calculating (2).

### B. Related work

The above approach was first introduced in [21], where it was specifically coupled with VQ-VAE and VQ-CPC scoring networks with a VQ layer that is learned during self-supervised training. In the intermediate evaluation below, I show that better performance can be achieved by applying DPDP on a large CPC model that uses $K$-means clustering to obtain VQ representations after self-supervised training.

A significant shortcoming of [21] was that it failed to connect the approach to much older work done at BBN in the 1980s. Building on an earlier heuristic approach [40], Roucos and Dunham [20] proposed an approach for low-bitrate speech coding: linear prediction speech frames are jointly segmented and quantized using a codebook through a dynamic programming procedure. Because they used an alignment-based segment cost, they did not have to include an explicit duration penalty as I do here. The only other difference is that their approach operates directly on speech frames, while the one here operates on learned self-supervised features.

The DPDP approach in this section is also a generalization of the more recent work of Chorowski et al. [39]. Instead of using a duration penalty, they enforce a prespecified number of segments into which a sequence needs to be divided. It can actually be shown that, if we set the duration penalty $w_{\mathrm{dur}}(l) = -l + 1$ as above and have a unique $\lambda$ for each utterance, their formulation is the dual of the one presented here. This only holds for this particular duration penalty and not in general. Moreover, their implementation used a greedy approximation; it was shown in [21] that better phone segmentation results are achieved when doing full dynamic programming.

### C. Intermediate evaluation: Unsupervised phone segmentation

I present a brief intermediate evaluation here. The goal is not to achieve state-of-the-art performance, but rather to verify that the DPDP method can give reasonable phone segmentation results (before using it as input for symbolic word segmentation). Complete experimental details are given in §VI, but in short, here I follow the experimental setup of [41] on the English Buckeye speech corpus. In addition to phone boundary precision, recall and $F_1$, I also report over-segmentation (OS): how many fewer/more boundaries are proposed compared to

TABLE I
INTERMEDIATE PHONE BOUNDARY SEGMENTATION RESULTS (%) ON BUCKEYE TEST DATA FOR STATE-OF-THE-ART MODELS AND DPDP SEGMENTATION.

| Model | Prec. | Rec. | $F_1$ | OS | $R$-val. |
|---|---|---|---|---|---|
| *Unsupervised:* | | | | | |
| GRU next-frame prediction [28] | 69.3 | 65.1 | 67.2 | −6.1 | 72.1 |
| GRU gate activation [29] | 69.6 | 72.6 | 71.0 | −4.1 | 74.8 |
| Self-sup. contrastive [41] | 75.8 | 76.9 | 76.3 | **−1.4** | 79.7 |
| SCPC [14] | **76.5** | 78.7 | **77.6** | - | **80.7** |
| ACPC [15] | 74.4 | 76.3 | 75.3 | - | 78.7 |
| Merged CPC+$K$-means (no DP) | 36.9 | **97.2** | 53.5 | 164.5 | −40.5 |
| DPDP VQ-CPC | 69.6 | 73.4 | 71.5 | 5.4 | 75.1 |
| DPDP VQ-VAE | 70.8 | 85.6 | 77.5 | 20.9 | 74.8 |
| DPDP CPC+$K$-means | 73.2 | 77.7 | 75.4 | 6.2 | 78.3 |
| *Supervised:* | | | | | |
| LSTM [44] | 87.8 | 83.3 | 85.5 | 5.4 | 87.2 |
| LSTM structured loss [45] | 85.4 | 89.1 | 87.2 | −4.1 | 88.8 |

the ground truth. The ideal is 0%. $F_1$ is not always sensitive enough to the trade-off between recall and over-segmentation, which motivates the $R$-value [42]: it gives a perfect score (100%) when a method has perfect recall and perfect OS.

Table I gives the results for five state-of-the-art approaches [14], [15], [28], [29], [41]. DPDP is used with three self-supervised models. The VQ-VAE and VQ-CPC scoring networks are from [32]. Both networks' encoders take log-mel spectrograms as input, downsample it by two, and discretize the feature vectors using a VQ layer with 512 codes. The networks (including the VQ layers) are trained on the 15-hour English training set from ZeroSpeech 2019 [16]. The third DPDP system uses the CPC-big model from [36], trained on LibriLight unlab-6k [43]. CPC-big takes raw speech input. After training the network, $K$-means is applied to features extracted from the second layer in its four-layer LSTM context network; features are extracted from the LibriSpeech train-clean-100 set. The result is a VQ codebook with 50 codes. Since the structure of the codebooks differ in the three DPDP systems, the duration weight in (3) is set individually: $\lambda = 3$ for the DPDP VQ-VAE, $\lambda = 400$ for DPDP VQ-CPC, and $\lambda = 2$ for DPDP CPC+$K$-means. These weights were tuned on Buckeye development data (see §VI).

Table I shows that the DPDP approach outperforms [28], [29], but performs slightly worse on most metrics compared to the recent state-of-the-art unsupervised phone segmentation approaches [14], [15], [41]. In terms of precision, $F_1$, OS and $R$-value, the DPDP systems also all outperform a system where no duration penalty is applied (no DP), i.e. repeated cluster indices are simply merged. This no-DP approach results in severe over-segmentation. Of the three DPDP systems, the CPC+$K$-means model (here used as a DPDP scoring network for the first time) achieves the best precision and $R$-value, while the DPDP VQ-VAE gives better recall and $F_1$ scores. Because of these intermediate results, I do not report results with the VQ-CPC in §VII.

## V. DPDP FOR UNSUPERVISED WORD SEGMENTATION FROM SYMBOLIC INPUT

The goal in unsupervised word segmentation from symbolic input is to break up an input sequence (normally of phonemes or phones) into subsequences representing words. My aim here is specifically to develop a model that can operate on the noisy symbolic sequences from an acoustic unit discovery model (such as the ones in §IV). The new model that I introduce is a simplified version of [22], as explained below in §V-B.

### A. DPDP autoencoder recurrent neural network (AE-RNN)

This approach is again an instance of DPDP. Let us say we have an autoencoding recurrent neural network (AE-RNN): an encoder-decoder model that takes a sequence of symbols as input, summarizes these into a single embedding vector using an encoder RNN, conditions a decoder RNN on the embedding, and finally tries to produce an output sequence that matches the input. Such a model could be seen as compressing the input sequence into a fixed-dimensional embedding, with the decoder then acting as a decompression algorithm. After training, the AE-RNN could give an accurate reconstruction of some inputs, while for some other inputs the reconstruction could be bad.

We can use this property to perform segmentation of an input symbol sequence $x_{1:T}$. Our goal is to divide the sequence into segments, each of which is accurately modelled by the AE-RNN. Figure 3 illustrates this for one possible segmentation. Here each $x_t \in \{k\}_{k=1}^K$ takes on a discrete value. Each of the four segments in Figure 3 can be scored based on its negative log likelihood according to the AE-RNN:

$$w_{\text{seg}}(x_{a:b}) = - \sum_{x_t \in x_{a:b}} \log P(x_t | x_{a:b}; \boldsymbol{\theta})$$

where $P(x_t | x_{a:b}; \boldsymbol{\theta})$ is the $t^{\text{th}}$ output of the decoder (obtained after a softmax layer) when the encoder is presented with segment $x_{a:b}$. $\boldsymbol{\theta}$ is the parameters of the AE-RNN. An overall segmentation cost is obtained by summing the individual segment costs. A different segmentation in Figure 3 will lead to a different overall cost, and we want to find the minimum.

There is a trivial solution: an AE-RNN that places a high output probability only on the first input symbol would give a result where every symbol is in its own segment. I therefore

again introduce a duration penalty $w_{\text{dur}}(l)$ and formulate this approach as an instance of DPDP (§III).

There are different options for $w_{\text{dur}}(l)$. In the main experiments in §VII, I use a linear penalty for the DPDP AE-RNN. But for the intermediate evaluation here (§V-C), I use a fixed probability mass function from a truncated gamma density (truncated at 50 symbols). This DPDP can correspond to a valid probabilistic model, but then we also need to model the number of segments in an utterance: I use a geometric distribution and set $\lambda = 1$ in (3). With these settings, this DPDP corresponds to a hidden semi-Markov model (HSMM) [46] with a fixed duration distribution for all states and the same emission distribution (parametrized by the AE-RNN).[1]

Up to now I implicitly assumed that we have a trained AE-RNN, but in reality we will need to fit its parameters $\boldsymbol{\theta}$. One way would be to start with a random segmentation, train the AE-RNN on these random segments, find the optimal segmentation under this AE-RNN, update the AE-RNN, re-segment, and so on. Another approach is that of [22]: instead of only considering the single best segmentation, they probabilistically sample likely segmentations. This could help better explore the segmentation space. Instead of either of these approaches, I found that a simple approach gave robust results: simply train the AE-RNN on the complete full-length utterances in the dataset. This AE-RNN is trained once and then fixed as the DPDP scoring network. This worked better than iterative refinement in developmental experiments.

### B. Related work

Word segmentation of symbolic sequences has a long history, particularly in the cognitive science community where these models are used to investigate how infants learn to segment words and discover the lexicon of their native language [2]. Most models are applied on transcribed phonemic input where every word is represented by the same sequence of phonemes.[2]

Some of the first approaches relied on the assumption that the transitions between symbols within a word are more predictable than across words [48]. Based on this, a word boundary can be predicted when the transition probability between two symbols dips below a threshold [49]. Another approach is to explicitly model word units. The hierarchical Dirichlet process model [2] uses a bigram language model over inferred word tokens with
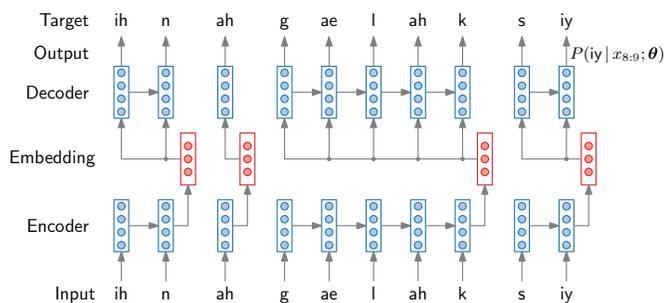


Fig. 3. An example where a symbolic sequence is divided into four segments, each processed through an autoencoding RNN (AE-RNN). The DPDP AE-RNN considers different segmentations and selects the one with the lowest combined negative log likelihood and duration cost.

[1]HSMMs [46] are probabilistic models where the duration that you stay in a hidden state is explicitly modelled. If we use a duration penalty that corresponds to a valid probability distribution, set $\lambda = 1$, model transitions between states (e.g. uniform), and model the number of segments in an utterance, then the DPDP AE-RNN *is* an HSMM with a fixed duration distribution and emission distributions parametrized by the AE-RNN (a parametrization not seen before, as far as I know). This is the variant of the DPDP AE-RNN used in the intermediate evaluation here (§V-C). However, the more general formulation of the DPDP AE-RNN given in this section is not bound by the constraint of a valid probabilistic duration prior, and indeed the model used in the speech segmentation experiments (§VII) uses a simple linear penalty (this gave better performance on development data). Also note that the ES-KMeans [11] and BES-GMM [10] that I compare to in §VII can also be formulated as HSMMs. HSMMs are therefore widespread, but the details of how they are parametrized are crucial and can lead to very different results.

[2]Some studies also consider variation in word pronunciation, taking phonetic input [47]. Although not evaluated in the intermediate evaluation in this section, the DPDP AE-RNN can also deal with such variation, as implicitly illustrated in the main experiments in §VII where I apply it on noisy discovered units.

| | Word boundary | | | Token |
| Model | Prec. | Rec. | $F_1$ | $F_1$ |
|---|---|---|---|---|
| Every phoneme as a word (no DP) | 27 | 100 | 43 | 3 |
| Transition probability [49] | 59 | 71 | 64 | 47 |
| Hierarchical Dirichlet process [2] | **90** | 74 | **87** | 74 |
| Adaptor grammar [3] | - | - | - | **88** |
| RNN memory segmenter [22] | 81 | **85** | 83 | 72 |
| DPDP AE-RNN | 78 | **85** | 81 | 69 |

priors to encourage predictable word sequences and a small vocabulary. A related nonparametric Bayesian approach is the adaptor grammar [3], which assumes that a corpus is generated from a set of re-write rules that can be learned probabilistically; below I specifically use the "colloc" variant, which includes bigram-like rules for co-occurring word units.

More recently, researchers have turned to neural sequence models, e.g. [50]. The RNN memory segmenter of Elsner and Shain [22] is particularly relevant here, since the DPDP AE-RNN can be seen as a simplification of this model. As explained above, I train the AE-RNN once on the full-length utterances in the dataset and then find the single best segmentation using (1) and (2). Instead of taking the single best segmentation, [22] iteratively refines the AE-RNN by probabilistically sampling boundaries. To help with the search problem, they train a separate RNN to propose boundaries, and use its smoothed output to sample boundaries. Below I show that their approach performs marginally better on phonemic data; but the DPDP AE-RNN was easier to tune in the chained approach where it is coupled with an acoustic unit discovery model (§VII).

### C. Intermediate evaluation: Unsupervised word segmentation of phoneme sequences

In this intermediate evaluation I verify that the DPDP AE-RNN gives reasonable word segmentation results when applied to phonemic sequences.[3] Again the goal is not state-of-the-art results (although I compare to established systems). I specifically consider results on the Brent corpus [48], a standard benchmark for symbolic word segmentation.

I tune the hyperparameters of the DPDP AE-RNN on the first 1000 utterances of the corpus. This includes the parameters for the gamma-based duration and geometric sentence length priors (§V-A). The AE-RNN consists of a three-layer GRU encoder and a single-layer GRU decoder, all with 200-dimensional hidden vectors; the latent embedding has 25 dimensions. The model is trained and evaluated on the full 9790-utterance corpus, as is the practice in other work [22]; results are very similar when removing the first 1000 utterances.

Table II shows word segmentation results. Apart from word boundary scores, where each boundary decision is evaluated

[3]I release a separate repository for the DPDP AE-RNN to reproduce the experiments in this section: `http://github.com/kamperh/dpdp_aernn`.

separately, the table also shows the word token $F_1$ score, which requires both boundaries of a word to be correctly predicted without any intermediate boundary proposals. This metric therefore also implicitly penalizes over-segmentation. We see that the DPDP AE-RNN gives similar or slightly worse performance compared to [22], on which it is based. Although the adaptor grammar [3] gives the best token $F_1$ score here, we will see in §VII that this method performs much worse when applied on discovered acoustic units—so better performance on phonemic sequences does not necessarily translate to better downstream speech segmentation results.

## VI. EXPERIMENTAL SETUP

Our main experiments involve the task of unsupervised word segmentation from speech. Since this is an unsupervised task, a model is typically trained and evaluated on the same data [51]. This means that developmental experiments need to be performed carefully, and I follow the practice of [23], [52]: Hyperparameters are chosen by training and evaluating on a development dataset. All hyperparameters are then fixed. For testing, the model is then trained and evaluated on a different set (potentially from another language) without any changes to any hyperparameters from the developmental experiments.

Development is performed on a 3.5-hour development set from Buckeye [53]. For evaluation on English, testing is then done on the 5-hour test set from Buckeye. There is no speaker overlap between these two sets. I then also apply the DPDP system to four completely unseen languages using the French, Mandarin, German and Wolof datasets from the ZeroSpeech 2017 challenge [23], respectively containing 24, 2.5, 25 and 10 hours of active speech. This data was also used in ZeroSpeech 2020. I report word boundary precision, recall, $F_1$, OS and $R$-value, as well as word token $F_1$, all with a 20 ms tolerance.

The goal is to see whether DPDP-based symbolic word segmentation (§V) can be applied on top of DPDP-discovered acoustic units (§IV) in order to segment speech. For acoustic unit discovery, I use the DPDP models already described in §IV-C. All the scoring networks (e.g. VQ-VAE and CPC+$K$-means) are pretrained on English data and the hyperparameters are set exactly as explained in that subsection.

For word segmentation, I use the DPDP AE-RNN. Here it is set up slightly differently from the one in the intermediate evaluation in §V-C: the one here has a 10-dimensional symbol embedding layer, a single 500-dimensional GRU encoder layer, a 50-dimensional latent embedding, and a single 500-dimensional GRU decoder. Instead of using a probabilistic duration penalty and sequence length prior, I use the simple duration penalty $w_{dur}(l) = -l+1$ with a weight of $\lambda = 3$. This is also the duration penalty used in the DPDP unit discovery models. This simpler DPDP AE-RNN gave better development performance when applied on discovered units. The AE-RNN is trained with Adam optimization [54] using a learning rate of $1 \cdot 10^{-3}$ for 1500 steps on the full tokenized utterances. The combined time for training and doing forward DPDP inference in the AE-RNN on the Buckeye development data is roughly 15 minutes on a machine with an NVIDIA GeForce RTX 3070 GPU and a single 2.5 GHz CPU. This is apart from

the DPDP forward inference time for acoustic unit discovery, which is roughly 12 minutes on its own. I release code at `http://github.com/kamperh/vqwordseg`.

## VII. Experiments: DPDP for Unsupervised Word Segmentation from Speech

The goal is to see how bottom-up acoustic unit discovery followed by symbolic word segmentation compare to other approaches. My specific new proposal is to chain the DPDP models from §IV and §V. I first compare this chained DPDP system to state-of-the-art systems on English (§VII-A). Then I consider different combinations of symbolic segmentation with acoustic unit discovery models (§VII-B). Finally I compare to other existing systems on non-English data (§VII-C).

### A. Comparing to joint self-supervised and direct whole-unit approaches on English

Table III shows word segmentation results on the English Buckeye test data, comparing the chained DPDP system (row 8) to existing approaches.

Like the DPDP system, row 1 is a system where bottom-up acoustic unit discovery is followed with symbolic word segmentation: a 50-component GMM is trained on unlabelled speech, the speech is encoded according to the most probable component for each MFCC frame, repeated components are merged, and the resulting tokenization is segmented with an adaptor grammar [3]. This is representative of the approaches followed in [1] (see §I). The DPDP system (row 8) outperforms this method on all metrics. This therefore represents a large improvement in bottom-up discovery followed by segmentation—the methodology shared by the two systems.

Instead of explicitly learning phone-like acoustic units, the systems in rows 2 to 4 all try to model higher-level units directly without an explicit lower-level acoustic unit layer. They either treat syllables [9] or words [10], [11] as the basic modelling unit (see §I). While these approaches outperform the adaptor grammar on GMM units (row 1), the DPDP system (row 8) outperforms all these direct whole-unit models across all metrics (except for [9] in row 2 giving a better OS score).

Finally I compare to two recent state-of-the-art CPC-based models (rows 5 and 6), described in §II. In contrast to the row 2 to 4 systems, these models have separate modules for acoustic unit learning and word segmentation which are learned jointly. This allows top-down information from word segmentation to influence bottom-up acoustic unit learning [14], [15]. The DPDP system performs better on word boundary recall, $F_1$ and OS, but the CPC-based systems achieve better precision and $R$-values (word token $F_1$ was not reported in these papers).

As a sanity check, I also give the result that would be obtained if no duration penalty was used on top of the self-supervised CPC+$K$-means model (row 7), i.e. repeated cluster indices are simply merged and then treated as words. As was the case in the intermediate phone segmentation experiments (Table I), this no-DP approach over-segments heavily, resulting in high boundary recall but poor precision.

In summary, the DPDP system outperforms bottom-up and direct whole-unit modelling, and performs similarly to joint self-supervised CPC-based models on an English benchmark.

### B. Symbolic word segmentation on discovered acoustic units

It is clear that symbolic segmentation on top of discovered units has come a long way since [1]. But the evaluation above only considers two particular combinations of acoustic unit discovery and symbolic word segmentation models (rows 1 and 8). These could be combined in different configurations.

Table IV shows token $F_1$ scores on Buckeye development data where different symbolic word segmentation models are applied on top of different acoustic unit discovery approaches. The best overall $F_1$ of 24.1% is achieved by applying DPDP AE-RNN word segmentation on DPDP CPC+$K$-means acoustic units (this is row 8 in Table III). Looking only at the word segmentation models (i.e. comparing columns) we see that the DPDP AE-RNN performs better than the transition probability [49] and adaptor grammar [3] models, irrespective of which noisy acoustic unit tokenization it is applied on.

This raises the question: is it necessary to apply DPDP for acoustic unit discovery, or is it sufficient to just apply the DPDP AE-RNN directly on merged cluster indices from some VQ model? Stated differently: is the combination of the two DPDP models important? We see that the combination is indeed important since it gives the best overall performance (24.1%). But we also see that applying the DPDP AE-RNN

TABLE III
UNSUPERVISED WORD SEGMENTATION RESULTS (%) ON BUCKEYE TEST DATA.

| Model | Word boundary | | | | | Token |
| --- | --- | --- | --- | --- | --- | --- |
| | Prec. | Rec. | $F_1$ | OS | $R$-val. | $F_1$ |
| 1: Adaptor gr. on GMM [1] | 15.9 | 57.7 | 25.0 | 261.5 | −139.9 | 4.4 |
| 2: $K$-means on syll. [9] | 27.7 | 28.9 | 28.3 | **4.5** | 37.7 | 19.3 |
| 3: ES-KMeans [11] | 30.3 | 16.6 | 21.4 | −45.1 | 39.1 | 19.2 |
| 4: BES-GMM [10] | 31.5 | 12.4 | 17.8 | −60.5 | 37.2 | 18.6 |
| 5: SCPC [14] | 34.8 | 31.0 | 32.8 | −10.8 | 44.5 | - |
| 6: mACPC [15] | **42.1** | 30.3 | 35.1 | −26.2 | **47.4** | - |
| 7: Merged CPC+$K$-means | 9.1 | **94.5** | 16.6 | 936.6 | −701.4 | 1.5 |
| 8: DPDP AE-RNN on DPDP CPC+$K$-means | 35.3 | 37.7 | **36.4** | 6.7 | 44.3 | **25.0** |

TABLE IV
WORD TOKEN $F_1$ SCORES (%) ON BUCKEYE DEVELOPMENT DATA USING DIFFERENT COMBINATIONS OF ACOUSTIC UNIT SEGMENTATION METHODS (ROWS) AND DIFFERENT WORD SEGMENTATION APPROACHES (COLUMNS).

| | | Symbolic word segmentation | | |
| --- | --- | --- | --- | --- |
| | | Trans. prob. [49] | Adaptor gram. [3] | DPDP AE-RNN |
| Acoustic unit discovery | MFCC+GMM | 5.8 | 4.7 | 17.2 |
| | Merged CPC+$K$-means (no DP) | 6.5 | 5.5 | 22.5 |
| | DPDP VQ-VAE | 9.4 | 5.3 | 16.4 |
| | DPDP CPC+$K$-means | 9.7 | 8.7 | **24.1** |

directly on merged cluster indices from the CPC+$K$-means model—without using DPDP—comes close (22.5%).

Comparing the DPDP VQ-VAE and DPDP CPC+$K$-means models, we also see the importance of training the self-supervised scoring network on substantial amounts of data: the VQ-VAE is trained on roughly 15 hours of data while the CPC+$K$-means model is trained on 6k hours (§IV-C). The latter leads to substantially better word segmentation scores.[4]

### C. Comparing to existing approaches on non-English data

Although the systems above were developed and tested on different corpora, they were still developed and tested on the same language (English). It is therefore unclear how the DPDP system will perform on an unseen zero-resource language when it is applied with the same hyperparameters. For the final quantitative evaluation, I apply the chained DPDP system without alteration to French, Mandarin, German and Wolof data and compare it to Track 2 submissions to ZeroSpeech 2017 and 2020 [17], [23]. Because the same data was used in the two challenges, this is one of the most comprehensive word segmentation benchmarks. As scoring networks on the non-English data here, I still use the CPC-big model trained on English (§IV-C).

Results are shown in Table V. I use the ZeroSpeech evaluation suite[5] for calculating word boundary and token scores here. The challenge baseline system [55] is an unsupervised term discovery system that optimizes for precision and therefore achieves very poor recall. The challenge topline system is an adaptor grammar applied to phonemic transcriptions of the training data. I only compare to a relevant subset of submitted systems.[6] With the exception of token $F_1$ on German, we see that the DPDP system gives the best performance on the word boundary $F_1$ and word token $F_1$ scores in all other cases. On Mandarin, the token $F_1$ score is improved by more than 14% absolute over the previous best result.[7]

Despite the DPDP system's improvements, it is clear on both languages that there is still a long way to go to get to the topline performance; on French, German and Wolof, the best scores are still far from the word token $F_1$ score achieved by the idealized topline system that segments phonemic input. In work done concurrently with this current paper (but only appearing after this paper's preprint), an extension of an earlier nonparametric Bayesian model to speech was proposed [58]; despite also achieving improvements over previous ZeroSpeech submissions, their approach is also still far from the topline.

---

[4]As a further comparison, note that the SCPC [14] in Table III is also trained on the same 15-hour English dataset from ZeroSpeech 2019 as the VQ-VAE in Table IV. On the Buckeye development data, this SCPC achieves a word boundary $F_1$ and $R$-value of respectively 33.0% and 45.6%. In comparison, the DPDP AE-RNN on DPDP VQ-VAE achieves 23.2% and 38.0%.

[5]https://github.com/zerospeech/zerospeech2020

[6]All submission results are available at https://zerospeech.com.

[7]I should note that the DPDP system implicitly makes use of a large amount of English data for acoustic unit discovery, while the other systems only use the training set of the respective languages. Through cross-lingual transfer, this could provide a benefit to the DPDP system, but it could also be detrimental since it never sees any within-language data for acoustic feature learning.

TABLE V
WORD SEGMENTATION RESULTS (%) FOR A SELECTION OF SYSTEMS FROM ZEROSPEECH 2017 AND 2020.

| Model | Word boundary | | | Token |
| --- | --- | --- | --- | --- |
| | Prec. | Rec. | $F_1$ | $F_1$ |
| *French:* | | | | |
| Baseline: Sparse term discovery [55] | 32.5 | 0.6 | 1.2 | 0.0 |
| ES-KMeans [11] | 37.0 | 52.2 | 43.3 | 6.3 |
| Probabilistic DTW [56] | 31.6 | **86.4** | 46.3 | 5.1 |
| Self-expressing autoencoder [57] | 34.0 | 83.9 | 48.4 | 8.3 |
| DPDP AE-RNN on DPDP CPC+$K$-means | **49.8** | 57.9 | **53.5** | **12.2** |
| Topline adaptor grammar on phonemes | 83.1 | 89.3 | 86.1 | 57.0 |
| *Mandarin:* | | | | |
| Baseline: Sparse term discovery [55] | 54.3 | 1.3 | 2.5 | 0.2 |
| ES-KMeans [11] | 42.6 | 75.6 | 54.5 | 8.1 |
| Probabilistic DTW [56] | 34.2 | 87.4 | 49.2 | 4.4 |
| Self-expressing autoencoder [57] | 36.5 | **91.9** | 52.2 | 12.1 |
| DPDP AE-RNN on DPDP CPC+$K$-means | **66.2** | 70.7 | **68.3** | **26.3** |
| Topline: Adaptor grammar on phonemes | 66.2 | 100 | 79.7 | 34.9 |
| *German:* | | | | |
| Baseline: Sparse term discovery [55] | 36.3 | 1.5 | 2.8 | 0.4 |
| ES-KMeans [11] | 42.9 | 66.9 | 52.3 | **14.5** |
| Probabilistic DTW [56] | 24.6 | **85.2** | 38.2 | 2.9 |
| Self-expressing autoencoder [57] | 27.1 | 82.8 | 40.9 | 7.5 |
| DPDP AE-RNN on DPDP CPC+$K$-means | **50.5** | 61.7 | **55.6** | 9.0 |
| Topline: Adaptor grammar on phonemes | 70.0 | 98.3 | 81.8 | 50.5 |
| *Wolof:* | | | | |
| Baseline: Sparse term discovery [55] | 49.9 | 1.4 | 2.7 | 0.2 |
| ES-KMeans [11] | 50.8 | 55.0 | 52.8 | 10.9 |
| Probabilistic DTW [56] | 35.2 | 48.0 | 40.6 | 4.2 |
| Self-expressing autoencoder [57] | 39.9 | **84.7** | 54.2 | 14.8 |
| DPDP AE-RNN on DPDP CPC+$K$-means | **63.1** | 56.5 | **59.6** | **15.0** |
| Topline: Adaptor grammar on phonemes | 81.3 | 93.2 | 86.9 | 60.2 |

### D. Qualitative analysis

The similar performance of the DPDP system compared to state-of-the-art CPC-based models on English (§VII-A), and its superior performance to other systems on the non-English ZeroSpeech benchmarks (§VII-C), show that the idea of bottom-up acoustic unit discovery followed by symbolic word segmentation should not be discarded. However, although word token $F_1$ scores have improved from around 5% (as in [1]) to around 25% (Tables III and V), absolute scores are still low. What are the characteristics of the words that are correctly segmented? And which words do the DPDP system still struggle to segment?

Figure 5 shows the token recall for the 15 word types that are most often segmented correctly (top) and the words that are incorrectly segmented most often (bottom).[8] Several of the correctly segmented words are shorter filler words: "mm" and "um" have the highest recall. In contrast, the word types with the lowest recall are longer words. Several of these end

---

[8]To get the recall for a particular word type, the number of correctly segmented word tokens (where both boundaries are correct without an intermediate prediction) is divided by the total number of tokens of this type.
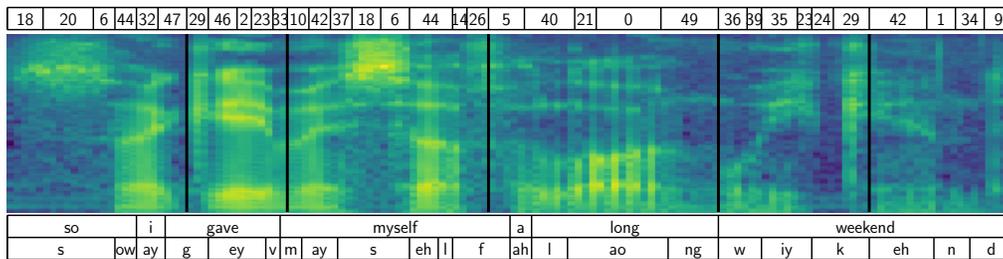
Fig. 4. An example segmentation from the chained DPDP system. The codes inferred from the DPDP CPC+$K$-means model are shown at the top. The vertical black lines on the spectrogram indicate where the DPDP AE-RNN places word boundaries. Ground truth word and phone boundaries are shown at the bottom.
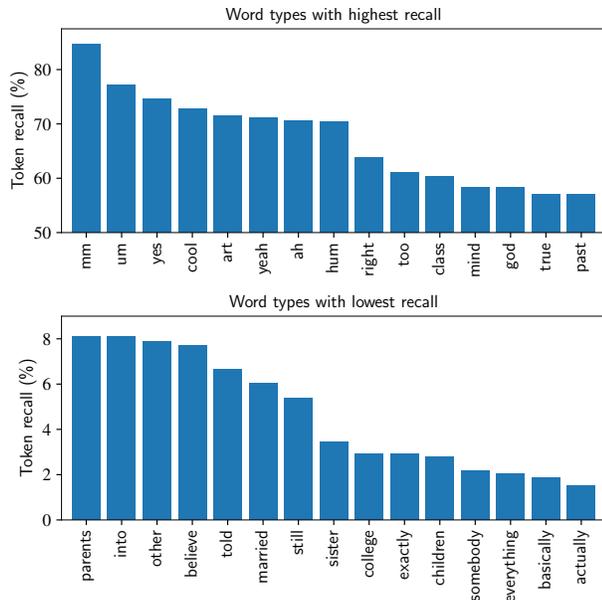


Fig. 5. Word token recall from the DPDP system for individual word types. The best (top) and poorest recalled words (bottom) are shown. Note the differences in scales on $y$-axes of the two plots.

on [l iy], and a manual review of segmentations reveals that a boundary is often placed just before these phones. In the segmentation example in Figure 4 we see a similar type of erroneous segmentation where "weekend" is segmented into two segments: [w iy k] and [eh n d].

Another indication that the DPDP system has a bias towards shorter words is if you look at the duration of correctly segmented word tokens: 231 ms with a standard deviation of 115 ms. Compare this to the same statistics for the ES-KMeans: 286 ms ± 124 ms. There is clearly large variation in these durations, but it is evident that the DPDP system's segmentation leans towards shorter words.

## VIII. CONCLUSION, DISCUSSION AND FUTURE WORK

I have described a model for unsupervised phone segmentation and a model for unsupervised symbolic word segmentation. I framed both as instances of a duration-penalized dynamic programming (DPDP) procedure, where a self-supervised neural scoring function is combined with a penalty term that encourages longer segments. I chained the two models to do word segmentation from speech and compared this to existing

direct whole-word and joint self-supervised approaches on standard benchmarks. The results showed that purely bottom-up phone discovery followed by symbolic segmentation— the methodology exemplified in the chained DPDP system— performs competitively to state-of-the-art models. This includes a comparison to two recent joint models that extend contrastive predictive coding (CPC) for segmentation [14], [15]. My argument is not that the chained methodology would in general be superior to joint modelling, but rather that it still has a place for further investigation for the task of speech segmentation. (It is even possible that the DPDP approach could be extended and improved through joint training in the future.)

Despite the encouraging results, absolute word segmentation scores remain low: on French, Mandarin, German and Wolof benchmarks, the best approaches are still far from a topline system where symbolic word segmentation is performed on transcribed phoneme sequences. On English, it is also interesting how similar the scores from the DPDP system are compared to the joint CPC-based approaches—despite the big differences in underlying methodology. Are we reaching the limits of what can be learned purely from unlabelled acoustic data? If so, how can we bridge the gap that remains between these systems and the idealized topline systems?

I believe that other sources of top-down information are needed. The qualitative analyses in this paper revealed a bias towards segmenting shorter (often filler) words. One approach would be to incorporate top-down information from a sparse term discovery system tailored towards discovering recurring but longer words [55], [56], [59], which could be used to bootstrap segmentation. Another approach would be to incorporate information from another modality; we know that infants have access to cross-situational cues from different modalities that can aid word learning [60]. In analogy, top-down information from visually grounded speech (VGS) models [61]–[63] could be used to help with segmenting particular words. In the analysis of [64], words that are accurately localized with an unsupervised VGS system include several words that are longer than those in the error analysis here (§VII-D). So the visual grounding information could provide a complementary segmentation signal (see [65] for very recent work on this). However, incorporating top-down signals for segmentation comes with its own challenges: [66] shows that speech segmentation is difficult even in a setting where text transcriptions or translations are used as an additional modality.

Apart from these future directions, there are two particular shortcomings of the DPDP system that need to be addressed. First, both DPDP approaches use a very coarse duration model, where the same duration penalty is applied irrespective of the acoustic unit or word segment under consideration. To address this, probabilistic models that explicitly model per-unit durations could be integrated or adapted to the DPDP approaches proposed here [46], [67]. Secondly, in contrast to some older speech segmentation models [10], [68], the DPDP autoencoding recurrent neural network (AE-RNN) does not infer an explicit lexicon—it can predict word boundaries, but word categories are not learned (categories are also not learned in [14], [15]). I did initial experiments where the latent DPDP AE-RNN representations are clustered with $K$-means, and also tried to then re-segment the acoustic units while taking this clustering into account, but this hurt overall segmentation performance. Building up a lexicon within the DPDP AE-RNN will therefore require further investigation.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] A. Jansen, E. Dupoux, S. J. Goldwater, M. Johnson, S. Khudanpur, K. Church, N. Feldman, H. Hermansky, F. Metze, R. Rose *et al.*, "A summary of the 2012 JHU CLSP workshop on zero resource speech technologies and models of early language acquisition," in *Proc. ICASSP*, 2013.

[2] S. J. Goldwater, T. L. Griffiths, and M. Johnson, "A Bayesian framework for word segmentation: Exploring the effects of context," *Cognition*, vol. 112, no. 1, pp. 21–54, 2009.

[3] M. Johnson and S. J. Goldwater, "Improving nonparameteric Bayesian inference: Experiments on unsupervised word segmentation with adaptor grammars," in *Proc. NAACL*, 2009.

[4] M. Elsner, S. J. Goldwater, and J. Eisenstein, "Bootstrapping a unified model of lexical and phonetic acquisition," in *Proc. ACL*, 2012.

[5] C.-y. Lee, T. O'Donnell, and J. R. Glass, "Unsupervised lexicon discovery from acoustic input," *Trans. ACL*, vol. 3, pp. 389–403, 2015.

[6] T. Taniguchi, S. Nagasaka, and R. Nakashima, "Nonparametric Bayesian double articulation analyzer for direct language acquisition from continuous speech signals," *IEEE Trans. Cogn. Developmental Syst.*, vol. 8, no. 3, pp. 171–185, 2016.

[7] R. Nakashima, R. Ozaki, and T. Taniguchi, "Unsupervised phoneme and word discovery from multiple speakers using double articulation analyzer and neural network with parametric bias," *Front. Robot. AI*, vol. 6, 2019.

[8] Y. Okuda, R. Ozaki, and T. Taniguchi, "Double articulation analyzer with prosody for unsupervised word and phoneme discovery," *arXiv preprint arXiv:2103.08199*, 2021.

[9] O. J. Räsänen, G. Doyle, and M. C. Frank, "Unsupervised word discovery from speech using automatic segmentation into syllable-like units," in *Proc. Interspeech*, 2015.

[10] H. Kamper, A. Jansen, and S. J. Goldwater, "A segmental framework for fully-unsupervised large-vocabulary speech recognition," *Comput. Speech Lang.*, vol. 46, pp. 154–174, 2017.

[11] H. Kamper, K. Livescu, and S. J. Goldwater, "An embedded segmental K-means model for unsupervised segmentation and clustering of speech," in *Proc. ASRU*, 2017.

[12] Y.-H. Wang, H.-y. Lee, and L.-s. Lee, "Segmental audio word2vec: Representing utterances as sequences of vectors with applications in spoken term detection," in *Proc. ICASSP*, 2018.

[13] A. van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.

[14] S. Bhati, J. Villalba, P. Żelasko, L. Moro-Velazquez, and N. Dehak, "Segmental contrastive predictive coding for unsupervised word segmentation," in *Proc. Interspeech*, 2021.

[15] S. Cuervo, M. Grabias, J. Chorowski, G. Ciesielski, A. Łańcucki, P. Rychlikowski, and R. Marxer, "Contrastive prediction strategies for unsupervised segmentation and categorization of phonemes and words," *arXiv preprint arXiv:2110.15909*, 2021.

[16] E. Dunbar, R. Algayres, J. Karadayi, M. Bernard, J. Benjumea, X.-N. Cao, L. Miskic, C. Dugrain, L. Ondel, A. W. Black *et al.*, "The Zero Resource Speech Challenge 2019: TTS without T," in *Proc. Interspeech*, 2019.

[17] E. Dunbar, J. Karadayi, M. Bernard, X.-N. Cao, R. Algayres, L. Ondel, L. Besacier, S. Sakti, and E. Dupoux, "The Zero Resource Speech Challenge 2020: Discovering discrete subword and word units," in *Proc. Interspeech*, 2020.

[18] D. Mochihashi, T. Yamada, and N. Ueda, "Bayesian unsupervised word segmentation with nested Pitman-Yor language modeling," in *Proc. ACL*, 2009.

[19] K. Uchiumi, H. Tsukahara, and D. Mochihashi, "Inducing word and part-of-speech with Pitman-Yor hidden semi-Markov models," in *Proc. ACL*, 2015.

[20] S. Roucos and M. O. Dunham, "A comparison of two methods for very-low-rate speech coding," in *Proc. MILCOM*, 1985.

[21] H. Kamper and B. van Niekerk, "Towards unsupervised phone and word segmentation using self-supervised vector-quantized neural networks," in *Proc. Interspeech*, 2021.

[22] M. Elsner and C. Shain, "Speech segmentation with a neural encoder model of working memory," in *Proc. EMNLP*, 2017.

[23] E. Dunbar, X. N. Cao, J. Benjumea, J. Karadayi, M. Bernard, L. Besacier, X. Anguera, and E. Dupoux, "The Zero Resource Speech Challenge 2017," in *Proc. ASRU*, 2017.

[24] J. Chorowski, G. Ciesielski, J. Dzikowski, A. Łancucki, R. Marxer, M. Opala, P. Pusz, P. Rychlikowski, and M. Stypułkowski, "Aligned contrastive predictive coding," in *Proc. Interspeech*, 2021.

[25] B. Varadarajan, S. Khudanpur, and E. Dupoux, "Unsupervised learning of acoustic sub-word units," in *Proc. ACL*, 2008.

[26] H. Gish, M.-H. Siu, A. Chan, and B. Belfield, "Unsupervised training of an HMM-based speech recognizer for topic classification," in *Proc. Interspeech*, 2009.

[27] C.-y. Lee and J. R. Glass, "A nonparametric Bayesian approach to acoustic model discovery," in *Proc. ACL*, 2012.

[28] P. Michel, O. Räsänen, R. Thiolliere, and E. Dupoux, "Blind phoneme segmentation with temporal prediction errors," *arXiv preprint arXiv:1608.00508*, 2016.

[29] Y.-H. Wang, C.-T. Chung, and H.-y. Lee, "Gate activation signal analysis for gated recurrent neural networks and its correlation with phoneme boundaries," in *Proc. Interspeech*, 2017.

[30] L. Ondel, H. K. Vydana, L. Burget, and J. Černocký, "Bayesian subspace hidden Markov model for acoustic unit discovery," *arXiv preprint arXiv:1904.03876*, 2019.

[31] Y.-A. Chung, W.-N. Hsu, H. Tang, and J. R. Glass, "An unsupervised autoregressive model for speech representation learning," in *Proc. Interspeech*, 2019.

[32] B. van Niekerk, L. Nortje, and H. Kamper, "Vector-quantized neural networks for acoustic unit discovery in the ZeroSpeech 2020 challenge," in *Proc. Interspeech*, 2020.

[33] B. van Niekerk, L. Nortje, M. Baas, and H. Kamper, "Analyzing speaker information in self-supervised models to improve zero-resource speech processing," in *Proc. Interspeech*, 2021.

[34] J. Chorowski, R. J. Weiss, S. Bengio, and A. van den Oord, "Unsupervised speech representation learning using WaveNet autoencoders," *IEEE Trans. Audio, Speech, Language Process.*, vol. 27, no. 12, pp. 2041–2053, 2019.

[35] A. Baevski, S. Schneider, and M. Auli, "vq-wav2vec: Self-supervised learning of discrete speech representations," in *Proc. ICLR*, 2020.

[36] T. A. Nguyen, M. de Seyssel, P. Rozé, M. Rivière, E. Kharitonov, A. Baevski, E. Dunbar, and E. Dupoux, "The Zero Resource Speech Benchmark 2021: Metrics and baselines for unsupervised spoken language modeling," in *NeurIPS SAS Workshop*, 2020.

[37] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE Trans. Audio, Speech, Language Process.*, 2021.

[38] C. Coupé, Y. M. Oh, D. Dediu, and F. Pellegrino, "Different languages, similar encoding efficiency: Comparable information rates across the human communicative niche," *Sci. Adv.*, vol. 5, no. 9, 2019.

[39] J. Chorowski, N. Chen, R. Marxer, H. Dolfing, A. Łańcucki, G. Sanchez, T. Alumäe, and A. Laurent, "Unsupervised neural segmentation and clustering for unit discovery in sequential data," in *NeurIPS PGR Workshop*, 2019.

[40] S. Roucos, R. Schwartz, and J. Makhoul, "Segment quantization for very-low-rate speech coding," in *Proc. ICASSP*, 1982.

[41] F. Kreuk, J. Keshet, and Y. Adi, "Self-supervised contrastive learning for unsupervised phoneme segmentation," in *Proc. Interspeech*, 2020.

[42] O. J. Räsänen, U. K. Laine, and T. Altosaar, "An improved speech segmentation quality measure: The R-value," in *Proc. Interspeech*, 2009.

[43] J. Kahn, M. Rivière, W. Zheng, E. Kharitonov, Q. Xu, P.-E. Mazaré, J. Karadayi, V. Liptchinsky, R. Collobert, C. Fuegen *et al.*, "Libri-light: A benchmark for ASR with limited or no supervision," in *Proc. ICASSP*, 2020.

[44] J. Franke, M. Mueller, F. Hamlaoui, S. Stueker, and A. Waibel, "Phoneme boundary detection using deep bidirectional LSTMs," in *in Proc. Speech Commun. ITG Symposium*, 2016.

[45] F. Kreuk, Y. Sheena, J. Keshet, and Y. Adi, "Phoneme boundary detection using learnable segmental features," in *Proc. ICASSP*, 2020.

[46] K. P. Murphy, "Hidden semi-Markov models (HSMMs)," 2002. [Online]. Available: http://www.cs.ubc.ca/~murphyk/mypapers.html

[47] M. Elsner, S. J. Goldwater, N. Feldman, and F. Wood, "A joint learning model of word segmentation, lexical acquisition and phonetic variability," in *Proc. EMNLP*, 2013.

[48] M. R. Brent, "An efficient, probabilistically sound algorithm for segmentation and word discovery," *Mach. Learn.*, vol. 34, no. 1-3, pp. 71–105, 1999.

[49] A. Saksida, A. Langus, and M. Nespor, "Co-occurrence statistics as a language-dependent cue for speech segmentation," *Devel. Sci.*, vol. 20, no. 3, p. e12390, 2017.

[50] K. Kawakami, C. Dyer, and P. Blunsom, "Learning to discover, ground and use words with segmental neural language models," in *Proc. ACL*, 2019.

[51] M. Versteegh, R. Thiollière, T. Schatz, X. N. Cao, X. Anguera, A. Jansen, and E. Dupoux, "The Zero Resource Speech Challenge 2015," in *Proc. Interspeech*, 2015.

[52] H. Kamper, A. Jansen, and S. J. Goldwater, "Unsupervised word segmentation and lexicon discovery using acoustic word embeddings," *IEEE Trans. Audio, Speech, Language Process.*, vol. 24, no. 4, pp. 669–679, 2016.

[53] M. A. Pitt, K. Johnson, E. Hume, S. Kiesling, and W. Raymond, "The Buckeye corpus of conversational speech: Labeling conventions and a test of transcriber reliability," *Speech Commun.*, vol. 45, no. 1, pp. 89–95, 2005.

[54] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2015.

[55] A. Jansen and B. Van Durme, "Efficient spoken term discovery using randomized algorithms," in *Proc. ASRU*, 2011.

[56] O. Räsänen and M. A. C. Blandón, "Unsupervised discovery of recurring speech patterns using probabilistic adaptive metrics," in *Proc. Interspeech*, 2020.

[57] S. Bhati, J. Villalba, P. Zelasko, and N. Dehak, "Self-expressing autoencoders for unsupervised spoken term discovery," in *Proc. Interspeech*, 2020.

[58] R. Algayres, T. Ricoul, J. Karadayi, H. Laurençon, S. Zaiem, A. Mohame, B. Sagot, and E. Dupoux, "DP-Parse: Finding word boundaries from raw speech with an instance lexicon," *Trans. ACL*, vol. 10, pp. 1051–1065, 2022.

[59] A. S. Park and J. R. Glass, "Unsupervised pattern discovery in speech," *IEEE Trans. Audio, Speech, Language Process.*, vol. 16, no. 1, pp. 186–197, 2008.

[60] O. Räsänen and H. Rasilo, "A joint model of word segmentation and meaning acquisition through cross-situational learning," *Psychol. Rev.*, vol. 122, no. 4, pp. 792–829, 2015.

[61] D. Harwath, A. Torralba, and J. R. Glass, "Unsupervised learning of spoken language with visual context," in *Proc. NIPS*, 2016.

[62] L. Gelderloos and G. Chrupała, "From phonemes to images: Levels of representation in a recurrent neural model of visually-grounded language learning," in *Proc. COLING*, 2016.

[63] O. Scharenborg *et al.*, "Linguistic unit discovery from multi-modal inputs in unwritten languages: Summary of the "Speaking Rosetta" JSALT 2017 Workshop," in *Proc. ICASSP*, 2018.

[64] K. Olaleye, D. Oneață, and H. Kamper, "Keyword localisation in untranscribed speech using visually grounded speech models," *IEEE J. Sel. Topics Signal Process.*, vol. 16, no. 6, pp. 1454–1466, 2022.

[65] P. Peng and D. Harwath, "Word discovery in visually grounded, self-supervised speech models," in *Proc. Interspeech*, 2022.

[66] R. Sanabria, H. Tang, and S. J. Goldwater, "On the difficulty of segmenting words with attention," in *EMNLP Insights Workshop*, 2021.

[67] M. J. Johnson and A. S. Willsky, "Bayesian nonparametric hidden semi-markov models," *J. Mach. Learn. Res.*, vol. 14, pp. 673–701, 2013.

[68] L.-s. Lee, J. R. Glass, H.-y. Lee, and C.-a. Chan, "Spoken content retrieval—beyond cascading speech recognition with text retrieval," *IEEE Trans. Audio, Speech, Language Process.*, vol. 23, no. 9, pp. 1389–1420, 2015.