

# PHONEME SEGMENTATION USING SELF-SUPERVISED SPEECH MODELS

Luke Strgar, David Harwath

University of Texas at Austin  
Department of Computer Science

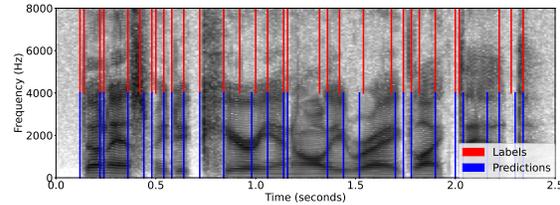
## ABSTRACT

We apply transfer learning to the task of phoneme segmentation and demonstrate the utility of representations learned in self-supervised pre-training for the task. Our model extends transformer-style encoders with strategically placed convolutions that manipulate features learned in pre-training. Using the TIMIT and Buckeye corpora we train and test the model in the supervised and unsupervised settings. The latter case is accomplished by furnishing a noisy label-set with the predictions of a separate model, it having been trained in an unsupervised fashion. Results indicate our model eclipses previous state-of-the-art performance in both settings and on both datasets. Finally, following observations during published code review and attempts to reproduce past segmentation results, we find a need to disambiguate the definition and implementation of widely-used evaluation metrics. We resolve this ambiguity by delineating two distinct evaluation schemes and describing their nuances. We provide a publicly available implementation of our work on Github<sup>1</sup>.

**Index Terms**— phonetic boundary detection, speech segmentation, self-supervised pre-training, transfer learning

## 1. INTRODUCTION

Phoneme boundary detection involves labeling the temporal boundaries between discrete phonemic units in a speech signal. Previously, phoneme segmentation has been studied and benchmarked in the supervised [1, 2, 3, 4] and unsupervised settings [5, 6]. In the former case, models are allowed to leverage a ground truth reference segmentation - a vector of phoneme onset, offset times - during training. In the latter case, the model only sees the input speech signal and is thus tasked with producing a segmentation by relying on the statistics of the underlying data alone. A third setting, known as forced-alignment or text-dependent phoneme segmentation, extends the supervised case by adding a temporally ordered list of phonetic identities to the model input. Conditioning on categorical phonetic identity means that model performance in the forced alignment setting typically supersedes text-independent supervised phoneme segmentation, which super-



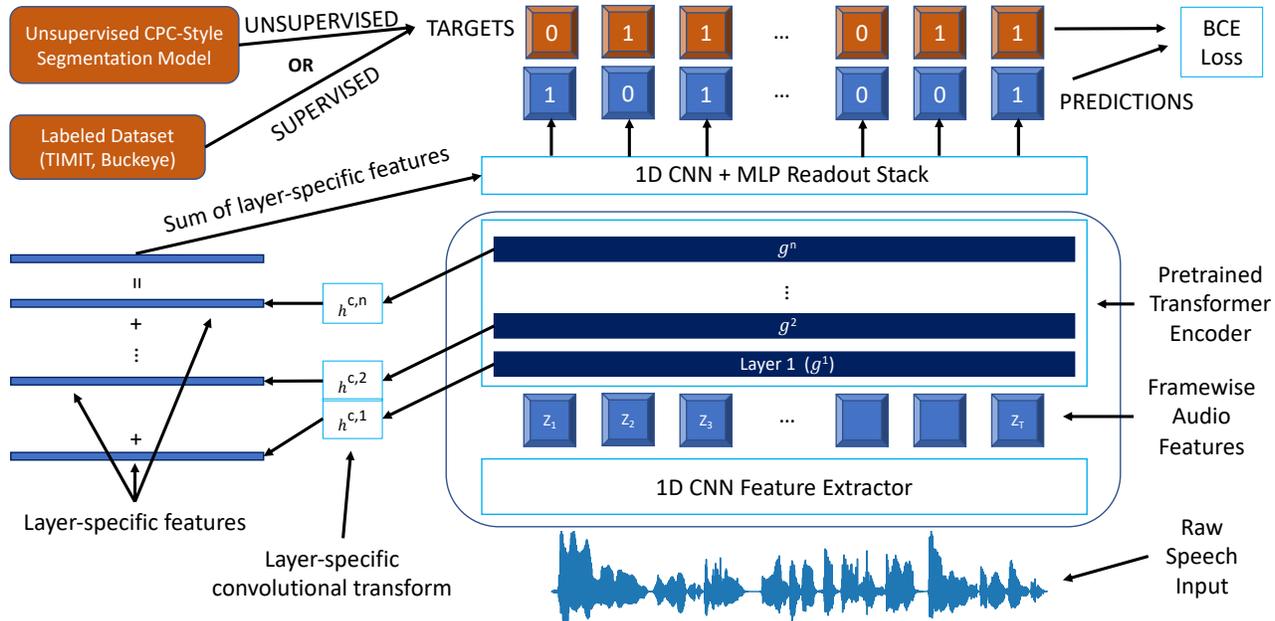
**Fig. 1.** Example spectrogram with ground truth and supervised model predicted boundaries.

sedes unsupervised predictions. In this paper, we focus on and report results for the unsupervised and text-independent supervised cases.

Self-supervised learning is a subclass of unsupervised learning in which training targets are derived from the input data itself. Recently, the speech processing field has benefited from the discovery and refinement of self-supervised strategies. Such heuristic strategies are often employed in a so-called model pre-training phase, and latter pre-trained models are fine-tuned or transfer learning is applied on specific downstream tasks. Numerous speech processing tasks have achieved new state-of-the-art (SotA) performances via application of fine-tuning and transfer learning to the information rich representations learned using self-supervised objectives. These include automatic speech recognition (ASR) [7, 8, 9], emotion recognition [10, 11], and speaker verification [11, 12], among others.

Inspired by the broad successes of self-supervised pre-training in speech processing, in this paper we explore its utility for phoneme segmentation. Specifically, we utilize pre-trained model checkpoints for two well-known and widely used self-supervised speech models, wav2vec2.0 [7] and HuBERT [8], and apply different strategies to refine these models' frame-wise representations for phoneme segmentation. In one case, we freeze the model's weights and extend its architecture with strategically placed, trainable, convolutional probe layers to manipulate and synthesize hierarchical features to output a binary predictor for each frame corresponding to the presence of a boundary. In a separate case, we append a simple projection layer to the pre-trained model's encoder and train the projection layer as well as all model

<sup>1</sup><https://github.com/lstrgar/self-supervised-phone-segmentation>



**Fig. 2.** Readout model architecture schematic. A pre-trained model extracts hierarchical features from the raw waveform. Features are processed by a series of convolutional networks and probability scores are computed. Finally, binary cross entropy loss is evaluated using model predictions and either ground truth labels or noisy labels estimated in an unsupervised manner.

weights end-to-end.

We evaluate and report results using the TIMIT [13] and Buckeye [14] speech corpora and find our model eclipses previous state-of-the-art performance on both datasets in the supervised and unsupervised settings. Unsupervised training is accomplished by furnishing a noisy label-set with the predictions of a separate model [5] that was trained in an unsupervised fashion using contrastive predictive coding [15] on a next frame prediction task. In supervised training, we also explore the label efficiency of our approach by sweeping over the amount of labeled training data used find that the model surpasses previous SotA performance with as little as 10% of the training set. With this work we demonstrate a successful application of self-supervised pre-training to the phoneme boundary detection task and offer a new SotA benchmark in the unsupervised and text-independent supervised settings. In addition, in later sections we describe ambiguity and inconsistencies in the commonly used evaluation protocol and offer a resolution in the form of two distinct evaluation schemes.

## 2. RELATED WORK

### 2.1. Phoneme Boundary Detection

Phoneme boundary detection has been explored using a variety of different model types and under various levels of supervision. In the text-independent, supervised setting, recent work revolves around the usage of recurrent neural network models. RNNs have been used as binary predictors [2] and

feature learners for a subsequent structured prediction task [1]. In text-dependent phoneme segmentation, probabilistic models such as HMMs have been applied [3], and recently a multi-task learning framework using pre-trained model features was proposed [4]. In the unsupervised setting, signal processing based approaches were initially dominant [16, 17], but recent research has focused on learning-based methods. [18] proposed a nonparametric Bayesian approach to unsupervised phonetic segmentation and clustering, and more recently the noise contrastive estimation principle has been applied to optimize the similarity of adjacent frames while making distant frames dissimilar [5]. Other work has applied contrastive learning at multiple levels by jointly optimizing both phoneme and word segmentation models [6].

### 2.2. Self-Supervised Pre-Training

Self-supervised pre-training has seen great success in numerous speech processing tasks. Borrowing ideas from research in natural language processing and computer vision, self-supervised models such as wav2vec2.0 [7] and HuBERT [8] are trained to reconstruct masked input from unmasked representations. The resulting internal representations obtained by these and other training objectives have been successfully applied to downstream tasks including ASR [7, 8, 9], emotion recognition [10, 11], and speaker verification [11, 12], among others.

**Table 1.** Results obtained in the fully supervised setting. \* Indicates application of the strict evaluation framework and  $\infty$  denotes author reported scores. The NA placeholder is used where results are not available. Bolded values indicate highest score for the specific metric and dataset.

Data	Model	Precision	Precision*	Recall	Recall*	F1	F1*	R-Value	R-Value*
Buckeye	Lin et al. [4] $\infty$	88.49	NA	90.33	NA	89.40	NA	90.90	NA
	Kreuk et al. [1] $\infty$	85.40	NA	89.12	NA	87.23	NA	88.76	NA
	W2V2 finetune	<b>94.01</b>	<b>90.56</b>	93.08	<b>90.28</b>	<b>93.54</b>	<b>90.42</b>	<b>94.41</b>	<b>91.81</b>
	HuBERT finetune	93.83	89.81	<b>93.11</b>	<b>90.28</b>	93.47	90.05	94.37	91.51
	W2V2 readout	93.38	89.14	92.74	89.66	93.00	89.40	93.99	90.96
	HuBERT readout	93.37	89.30	92.95	89.94	93.16	89.62	94.13	91.15
TIMIT	Lin et al. [4] $\infty$	93.42	NA	95.96	NA	94.67	NA	95.18	NA
	Kreuk et al. [1] $\infty$	94.03	NA	90.46	NA	92.22	NA	92.79	NA
	Kreuk et al. [1]	92.94	92.14	92.31	89.26	92.63	90.68	93.66	91.71
	W2V2 finetune	96.90	<b>94.35</b>	<b>96.30</b>	<b>93.91</b>	<b>96.60</b>	<b>94.13</b>	<b>97.04</b>	<b>94.96</b>
	HuBERT finetune	<b>96.93</b>	94.31	96.09	93.68	96.51	94.00	96.92	94.83
	W2V2 readout	96.67	93.75	95.56	92.65	96.11	93.20	96.55	94.10
	HuBERT readout	96.50	93.23	95.93	93.47	96.21	93.35	96.71	94.33

### 3. PROBLEM STATEMENT

In phoneme segmentation the input is a raw speech waveform  $x \in \mathcal{X}$  represented as  $x = (x_0, x_1, \dots, x_N)$  where each  $x_i$  is a single floating point value representing relative pressure in the transmission medium. Typically,  $x$  will be pre-processed and temporally down-sampled by some transformation  $f_x : \mathcal{X} \rightarrow \mathcal{Z}$  to produce  $f_x(x) = z = (z_1, z_2, \dots, z_T)$  where  $T \ll N$  and  $z_i \in \mathbb{R}^d$ . Here,  $z \in \mathbb{R}^{T \times d}$  can be thought of as representing a series of acoustic feature frames and  $T$  now encodes the temporal resolution we desire to make predictions with.

Each input speech sample is paired with a label sequence of time stamps  $y = (y_1, y_2, \dots, y_K)$  where each  $y_i$  indicates the presence of one boundary and is represented as a single floating point value encoding the time units relative to the beginning of the utterance. Similar to the down-sampling of  $x$ , one might choose to bin  $y$  such that each element is converted to units of acoustic feature frames. We call this representation  $\bar{y} = (\bar{y}_1, \dots, \bar{y}_K)$  and note that  $K$ ,  $N$ , and  $T$  may vary across input / label pairs.

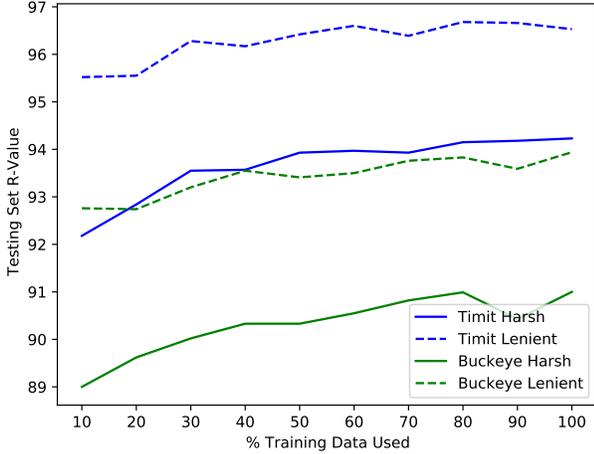
Automatic phoneme segmentation thus asks for a prediction,  $\hat{y} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_K)$  that closely matches the ground truth label  $\bar{y}$ . Classically, the closeness of a reference and predicted segmentation is evaluated with the precision, recall, F1, and R-Value metrics [19]. Section 5 describes these quantities as well as their nuances in detail.

### 4. MODEL DESCRIPTION

Inspired by the success of self-supervised learning in numerous speech processing tasks we adopt pre-trained wav2vec2.0 and HuBERT model checkpoints to compose the backbone of a frame-wise binary classifier. Both pre-trained models share a similar architecture. For our purposes we consider only the encoder, which we denote by the function composition  $g \circ f$ . Elaborating on the component functions of this composition,  $f : \mathcal{X} \rightarrow \mathcal{Z}$  is commonly referred to as the convolutional feature extractor, which processes raw waveform input and outputs a time series of latent speech representations. Thus,  $f$  acts like the previously defined  $f_x$ ; however,  $f$  is not strictly a pre-processing step since it is learned during end-to-end training of wav2vec2.0 and HuBERT. Meanwhile,  $g : \mathcal{Z} \rightarrow \mathcal{C}$  is known as the context network, which applies learned attention masks to synthesize a context-aware representation  $c_i \in \mathcal{C}$  from each  $z_i \in \mathcal{Z}$ .  $g$  is itself a compositional function built from a cascade of  $n$  transformer self-attention blocks. Thus, we can also write  $g = g^n \circ g^{n-1} \circ \dots \circ g^1$ . Note that functions  $f$  and  $g$  may be initialized by either wav2vec2.0 or HuBERT.

We develop two separate classification model formulations built on-top of the pre-trained network backbone. The first case, which we call *fine-tune* mode, appends a single linear projection layer,  $h^{ft}$ , to the output of the pre-trained model. As the name suggests, in this setting, the entire pre-trained model and added projection receive gradient updates, and the model can be formalized as the function composition  $f \circ g \circ h^{ft}$ .

The second case, called *readout* mode, is depicted in Figure 2. Here, we freeze the pre-trained model and apply



**Fig. 3.** Supervised model with wav2vec2.0 backbone in *read-out* mode trained from scratch on incrementally larger fractions of labeled data. Vertical axis shows testing set R-Value performance.

learned, layer-specific convolutions,  $h^{c,1}, h^{c,2}, \dots, h^{c,n}$ , to feature representations extracted from each  $g^i$ . The outputs are then summed and passed through a final series of convolutions and perceptron layers, denoted  $h^{r,o}$ . Empirically, we discovered that applying a learned weight parameter to each layer’s processed features before computing the summation improved performance; however, we omit these terms in the following expression for simplicity. Denoting the outputs of each  $g^i$  as  $c^i$ , the *readout* model can be formalized as  $h^{r,o}(\sum_i h^{c,i}(c^i))$ .

Both models output a series of frame-wise binary labels,  $\hat{y}_b$ , where a 1 is interpreted as the occurrence of a boundary. Given a training set of input utterances  $\mathcal{S} = \{x^i, y_b^i\}_{i=1}^m$ , loss is computed and models are updated according to the binary cross entropy (BCE) objective function in Equation 1. Here, the term  $w^*$  is a strictly positive weight value assigned to the loss associated with frames where the reference ground truth indicates the presence of a boundary.

$$\mathcal{L}_{BCE} = \sum_{i=0}^m \sum_{j \in y_b^i} w^* y_{b,j}^i \log(\hat{y}_{b,j}) + (1 - y_{b,j}^i) \log(1 - \hat{y}_{b,j}) \quad (1)$$

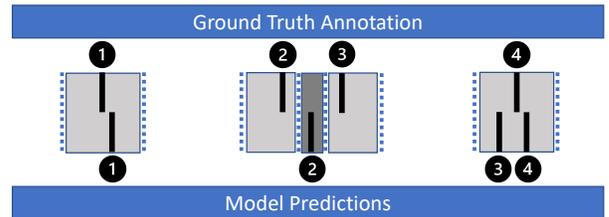
Optimization of the above objective is explored in the supervised and unsupervised settings. In the former case, we rely on the time-aligned transcriptions that come as part of the TIMIT and Buckeye corpora to construct supervised targets. In the latter case, we perform inference over the TIMIT and Buckeye training sets using the unsupervised model provided by Kreuk et al. [5]. These predicted labels then serve as the supervisory signal used during training with BCE loss.

## 5. EVALUATION

Previous work has articulated the challenges of conducting representative evaluations of phoneme segmentation results [19]. These issues are further confounded by the wide range of (prediction to label) temporal tolerance levels found in the literature for true positives identification. Here, we consider a 20 millisecond tolerance window on either side of a ground truth label for true positive calling, which is consistent with recent work [1, 5, 6, 4].

We report our results in terms of precision, recall, F1, and R-Value according to their definition in [19]. These metrics and their interpretation are widely cited in the phoneme segmentation task literature; however, based on a combination of published code review and attempts to reproduce results we believe there remains meaningful ambiguity in the calculation of precision, recall and their derivative quantities (e.g., F1 and R-Value). By elaborating on these definitions and their implications below, we hope to align the community around shared standards for reporting phoneme boundary detection results.

When computing precision the primary source of ambiguity revolves around interpreting multiple positive boundary predictions falling within the tolerance window of a ground truth boundary. For recall, the parallel situation arises where a single predicted positive boundary falls within the tolerance window of more than one ground truth boundary. See Figure 4 for a visual illustration of the ambiguous situations arising during evaluation.



**Fig. 4.** Illustration of ambiguities during phoneme segmentation evaluation. Vertical black stripes indicate ground truth (top) and predicted (bottom) boundaries. The light gray regions correspond to ground truth boundary tolerance windows and the dark gray region shows where two tolerance windows overlap. Predicted and ground truth boundaries 1 match. Ground truth boundaries 2, 3 both match predicted boundary 2 while predicted boundaries 3 and 4 both match and ground truth 4.

Without loss of generality, consider a boundary predictor operating at 50Hz (i.e. predictions correspond to the occurrence of a boundary in a 20 millisecond window). While computing the precision of this model’s predictions with a 20 millisecond tolerance window, one could encounter an isolated ground truth boundary at frame  $p$  and model boundary

**Table 2.** Results obtained in the unsupervised setting. A noisy label-set was furnished using publicly available checkpoints from an unsupervised segmentation model [5].

Data	Model	Precision	Precision*	Recall	Recall*	F1	F1*	R-Value	R-Value*
Buckeye	Bhati et al. [6] <sup>∞</sup>	76.53	NA	78.72	NA	77.61	NA	80.72	NA
	Kreuk et al. [5] <sup>∞</sup>	75.78	NA	76.86	NA	76.31	NA	79.69	NA
	Kreuk et al. [5]	77.17	72.21	79.71	75.55	78.42	73.85	81.39	77.28
	W2V2 finetune	82.15	75.56	<b>85.13</b>	<b>79.47</b>	83.61	77.47	85.81	80.33
	HuBERT finetune	83.09	76.62	84.47	78.75	83.77	<b>77.67</b>	86.11	80.79
	W2V2 readout	<b>84.24</b>	<b>77.92</b>	82.88	77.41	83.55	<b>77.67</b>	85.92	<b>80.95</b>
	HuBERT readout	83.35	75.29	84.68	79.37	<b>84.01</b>	77.28	<b>86.31</b>	80.13
TIMIT	Bhati et al. [6] <sup>∞</sup>	84.63	NA	86.04	NA	85.33	NA	87.44	NA
	Kreuk et al. [5] <sup>∞</sup>	83.89	NA	83.55	NA	83.71	NA	86.02	NA
	Kreuk et al. [5]	85.27	81.42	83.48	76.53	84.36	78.90	86.57	81.71
	W2V2 finetune	88.93	82.16	<b>88.60</b>	80.83	88.76	81.49	90.40	84.18
	HuBERT finetune	89.05	82.07	88.44	80.70	88.75	81.38	90.37	84.08
	W2V2 readout	90.69	<b>84.92</b>	86.78	78.52	88.69	81.59	89.90	83.69
	HuBERT readout	<b>90.98</b>	82.44	88.48	<b>81.18</b>	<b>89.71</b>	<b>81.81</b>	<b>90.98</b>	<b>84.45</b>

predictions at  $p - 1$ ,  $p$ ,  $p + 1$ . Different approaches to this calculation could result in a three-fold difference in performance, and the situation would be exacerbated by an increase in the predictor’s frame rate. In practice, the statistics of English language phoneme duration and presentation render a three-fold performance difference highly unlikely; however, others have reported differences of up to 5% [20], and our results consistently show deviations of 3-4% in the supervised setting and 5-7% in the unsupervised setting.

We then wish to delineate a *strict* and *lenient* evaluation scheme for phoneme boundary detection where the *strict* scheme prohibits double counting and the *lenient* scheme allows it. Specifically, while computing the hit rate [19] of an automated phoneme segmenter, in the *strict* scheme once a ground truth boundary is matched by a predicted boundary the ground truth is removed from consideration for matching additional model predictions. On the other hand, in the *lenient* scheme the same ground truth boundary may match multiple predicted boundaries so long as they fall within the tolerance window. Further, in the *lenient* scheme, the hit rate used for computing precision and that for recall may differ since more than one predicted boundary is allowed to match a the same ground truth and visa versa.

We denote results following the *strict* scheme with a \* and then define F1\* and R-Value\* as those metrics computed with their *strict* counterparts P\* (precision\*) and R\* (recall\*). Our code reviews and efforts to reproduce past results indicate that previous SotA methods use the *lenient* scheme. For parity, our results tables below include both *strict* and *lenient* scores for our models. In some cases where we were able to reproduce previous published results we also add *strict* scores.

In other cases, it was not possible to verify the exact evaluation framework used by some authors. However, all these papers explicitly describe sharing evaluation methodology with the aforementioned previous SotA, against which they benchmark their model performance. Accordingly, we assume they also evaluate performance using the *lenient* framework.

## 6. EXPERIMENTS

### 6.1. Datasets

We used the TIMIT [13] and Buckeye [14] speech corpora to train and evaluate the *fine-tune* and *readout* models. For TIMIT, we used the standard train/test split and sampled 10% of the training data for model validation. For Buckeye, we followed previous work [1, 5, 2] in our training, validation, and testing set construction. First, we split the corpus at the speaker level, reserving 80%, 10%, 10% for training, validation, and testing, respectively. In addition, long recordings were split during non-vocal noise and silence into shorter continuous speech segments such that each segment starts and ends with no more than 20 milliseconds of non-speech.

### 6.2. Experimental Setup

Experiments conducted with HuBERT used the base architecture and those with wav2vec2.0 used the small architecture. Both model checkpoints were pre-trained on Librispeech [21] and collected from Fairseq [22]. We explored the effectiveness of larger model architectures (e.g. wav2vec2.0 large, HuBERT large/x-large) but found they offered no boost on fi-

nal performance metrics. For our unsupervised experiments, we used model checkpoints made available with the code accompanying [5] to bootstrap labels for TIMIT and Buckeye.

All models were trained on an NVIDIA Quadro RTX 8000 with a batch size of 16 for 50 epochs. The Adam optimizer was used with a learning rate of  $1e-3$  and  $1e-4$  while training in *readout* and *fine-tune* mode, respectively. Models were regularly evaluated during training using the validation set’s R-Value\* and the best performing model was saved for testing.

In *readout* mode the layer specific convolutions,  $h^{c,i}$  were defined with a kernel size of 9, stride of 1, and 768 input and channels. The output architecture  $h^{r,o}$  is a depth five convolution stack with a shared kernel size of 3 and stride of 1 followed by a linear projection. As we mentioned previously, in this setting we also added a parameter to learn a weighted sum of the layer specific features before application of  $h^{r,o}$ .

Throughout our experiments we explored various values of  $w^*$  - the loss weight applied to frames labeled as boundary positive. In all supervised experiments,  $w^*$  was ultimately set to 1 for the entire duration of model training. We made anecdotal observations that setting  $1 < w^* < 2.5$  tended to speed up model convergence; however,  $w^*$  had to be subsequently turned down and training continued to obtain the best performance metrics. In the unsupervised setting, we found that, relative to ground truth labels, the noisy labels scored substantially lower in recall than precision. Acknowledging the need then to incentive positive predictions, we swept values of  $w^*$  and obtained optimal validation performance using  $w^* = 1.4$ .

### 6.3. Results

In Table 1 we report results for our models in the fully supervised setting. We also include reported scores from Lin et al. [4], Kreuk et al. [1], which stand as previous benchmark results in text-dependent and text-independent phoneme segmentation, respectively. Another result we include is our attempt at reproducing Kreuk et al.’s [1] results for TIMIT - here we are able to share both the *harsh* and *lenient* evaluations. We were unable to reproduce comparable scores for Buckeye using the model from [1]. Altogether, results indicate that the best of our four models - composed through a selection of a backbone pre-trained network and *fine-tune* or *readout* mode - eclipse previous SotA in every metric category for both TIMIT and Buckeye. With few exceptions, all four of our models outpace previous SotA, and we emphasize that our top performing model, which was trained in the text-independent setting, surpasses the performance of SotA text-dependent [4].

Figure 3 highlights the small amount of labeled training data required to surpass previous SotA performance. Results reported in this figure come from experiments with a *readout* mode model trained with a wav2vec2.0 back-bone. For both

TIMIT and Buckeye we obtain R-Value SotA using only 10% of the labeled data from the respective training sets.

Table 2 reports results for models in the unsupervised setting along with other previous SotA results. As in the supervised case, our best performing model achieves a new SotA result for both TIMIT and Buckeye in every metric category. Notably, wherein the supervised setting a typical deviation between the *lenient* and *harsh* schemes is in the 2-3% range, in the unsupervised setting we observe deviations of, in some cases, more than 8%. As the Kreuk et al. [5] and Bhati et al. [6] unsupervised models reported here perform inference through a peak-picking algorithm over a learned representation, it is possible that over prediction near boundaries stems from the difficulty of enforcing temporally precise transitions in the learned representation. Similarly, as our models are trained using a noisy label-set bootstrapped from [5], our model is liable to the same failure modes.

During experiments with noisy (unsupervised) label-sets, we explored the impact of multiple self-training loops to refine the labels and improve final model test performance. Ultimately, we observed marginal gains that did not inspire a deep exploration of how bootstrapped labels could be refined in an unsupervised fashion. In fact, in *fine-tune* mode, performance declined after multiple self-training loops. Incidentally, throughout our experiments in the unsupervised setting, *readout* models tended to perform better than their fully fine-tuned counterparts. Relevant metrics observed during training indicated that the more expressive fine-tuned models were much more liable to over-fit label noise than their *readout* mode counterparts.

## 7. DISCUSSION

Here we introduced a new model formulation based on self-supervised pre-training and transfer learning to perform phoneme boundary detection in the supervised and unsupervised settings. We empirically demonstrate that our formulation sets a new SotA benchmark for both settings on standard datasets used for the task - the TIMIT and Buckeye speech corpora. Additionally, we bring to the community’s attention a need for shared implementation strategies for key evaluation metrics and define two evaluation frameworks that can be used to alleviate future ambiguity.

We believe there are several promising directions for future work. First, an exploration of regularization and self-training strategies to improve noisy label-sets will likely push unsupervised results further than we have been able to. Second, in the supervised setting we obtained excellent performance even with small amounts of training data. We are optimistic then that low resource languages can benefit from self-supervised pre-training for phoneme boundary detection. Finally, our model formulation may be, with minimal modifications, well-suited to alternate speech segmentation tasks.

## 8. REFERENCES

- [1] Felix Kreuk, Yaniv Sheena, Joseph Keshet, and Yossi Adi, "Phoneme boundary detection using learnable segmental features," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 8089–8093.
- [2] Joerg Franke, Markus Mueller, Fatima Hamlaoui, Sebastian Stueker, and Alex Waibel, "Phoneme boundary detection using deep bidirectional lstms," in *Speech Communication; 12. ITG Symposium*. VDE, 2016, pp. 1–5.
- [3] Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger, "Montreal forced aligner: Trainable text-speech alignment using kaldii," in *Interspeech*, 2017, vol. 2017, pp. 498–502.
- [4] Binghuai Lin and Liyuan Wang, "Learning acoustic frame labeling for phoneme segmentation with regularized attention mechanism," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 7882–7886.
- [5] Felix Kreuk, Joseph Keshet, and Yossi Adi, "Self-supervised contrastive learning for unsupervised phoneme segmentation," 2020.
- [6] Saurabhchand Bhati, Jesús Villalba, Piotr Żelasko, Laureano Moro-Velazquez, and Najim Dehak, "Segmental contrastive predictive coding for unsupervised word segmentation," 2021.
- [7] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 12449–12460, 2020.
- [8] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhota, Ruslan Salakhutdinov, and Abdelrahman Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [9] Cheng Yi, Jianzhong Wang, Ning Cheng, Shiyu Zhou, and Bo Xu, "Applying wav2vec2. 0 to speech recognition in various low-resource languages," *arXiv preprint arXiv:2012.12121*, 2020.
- [10] Omar Mohamed and Salah A Aly, "Arabic speech emotion recognition employing wav2vec2. 0 and hubert based on baved dataset," *arXiv preprint arXiv:2110.04425*, 2021.
- [11] Yingzhi Wang, Abdelmoumene Boumadane, and Abdelwahab Heba, "A fine-tuned wav2vec 2.0/hubert benchmark for speech emotion recognition, speaker verification and spoken language understanding," *arXiv preprint arXiv:2111.02735*, 2021.
- [12] Zhengyang Chen, Sanyuan Chen, Yu Wu, Yao Qian, Chengyi Wang, Shujie Liu, Yanmin Qian, and Michael Zeng, "Large-scale self-supervised speech representation learning for automatic speaker verification," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 6147–6151.
- [13] John S Garofolo, "Timit acoustic phonetic continuous speech corpus," *Linguistic Data Consortium, 1993*, 1993.
- [14] Mark A Pitt, Laura Dilley, Keith Johnson, Scott Kiesling, William Raymond, Elizabeth Hume, and Eric Fosler-Lussier, "Buckeye corpus of conversational speech (2nd release)," *Columbus, OH: Department of Psychology, Ohio State University*, pp. 265–270, 2007.
- [15] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, "Representation learning with contrastive predictive coding," 2018.
- [16] James Glass, "A probabilistic framework for segment-based speech recognition," *Computer Speech and Language*, vol. 13, pp. 137–152, 2003.
- [17] Sorin Dusan and Lawrence Rabiner, "On the relation between maximum spectral transition positions and phone boundaries," in *Interspeech*, 2006.
- [18] Chia-ying Lee and James Glass, "A nonparametric Bayesian approach to acoustic model discovery," in *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2012.
- [19] Okko Johannes Räsänen, Unto Kalervo Laine, and Toomas Altsaar, "An improved speech segmentation quality measure: the r-value," in *Tenth Annual Conference of the International Speech Communication Association*. Citeseer, 2009.
- [20] Okko Räsänen, "Speech segmentation and clustering methods for a new speech recognition architecture," *helsinki university of technology*, 2007.
- [21] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: An asr corpus based on public domain audio books," *ICASSP*, 2015.
- [22] Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael

Auli, “fairseq: A fast, extensible toolkit for sequence modeling,” in *Proceedings of NAACL-HLT 2019: Demonstrations*, 2019.