

# I SEE WHAT YOU HEAR: A VISION-INSPIRED METHOD TO LOCALIZE WORDS

Mohammad Samragh, Arnav Kundu, Ting-Yao Hu, Aman Chadha\*, Ashish Srivastava\*, Minsik Cho, Oncel Tuzel, Devang Naik

Apple

## ABSTRACT

This paper explores the possibility of using visual object detection techniques for word localization in speech data. Object detection has been thoroughly studied in the contemporary literature for visual data. Noting that an audio can be interpreted as a 1-dimensional image, object localization techniques can be fundamentally useful for word localization. Building upon this idea, we propose a lightweight solution for word detection and localization. We use bounding box regression for word localization, which enables our model to detect the occurrence, offset, and duration of keywords in a given audio stream. We experiment with LibriSpeech and train a model to localize 1000 words. Compared to existing work [1], our method reduces model size by 94%, and improves the F1 score by 6.5%.

**Index Terms**— voice event detection, word detection, word localization

## 1. INTRODUCTION

Recent advancements in automatic speech recognition (ASR) technologies have made human machine interaction seamless and natural [2, 3, 4]. ASR systems are designed to accurately comprehend any speech data, hence, they are often computationally expensive, power hungry, and memory intensive. Acoustic models with limited-vocabulary can only recognize certain words, but they offer computationally efficient solutions [5, 6, 7]. In this paper we focus on the latter.

Limited-vocabulary models are important for several reasons. First, users often interact with their devices using simple commands [8], and recognizing these commands may not necessarily require an ASR model. Second, limited-vocabulary models are needed to detect trigger phrases, e.g., “Alexa, OK Google, hey Siri”, which indicate that a user wants to interact with the ASR model. A limited vocabulary model should **(a)** accurately recognize if certain words are spoken, and **(b)** precisely locate the occurrence time of the words. The latter is rather important due to privacy and efficiency reasons, as an ASR model should only be queried/executed when users intend to interact with it. Additionally, a limited-vocabulary model with localization capabilities can improve the accuracy of an ASR model by providing noise-free payloads to it.

A plethora of existing work focus on keyword detection [5, 6, 7], yet, efficient and accurate word localization needs more investigation. In computer vision, object localization has been solved using bounding box detection [9, 10, 11]. Since an audio can be interpreted as a 1-D image, similar techniques can be used in principle to localize words in an audio. The first effort in this track is SpeechYolo [1], which shows the great potential of using visual object detection techniques for word localization.

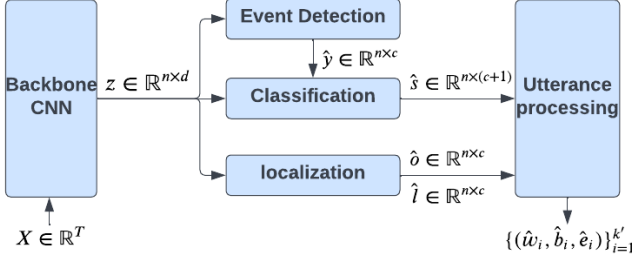
This paper presents an alternative vision-based word localizer. In our design, we pay attention to several important properties: **(a)** having small memory footprint, **(b)** the ability to process streaming audio, **(c)** accurate word detection and **(d)** accurate word localization. To achieve these goals, we propose customized metrics that indicate the presence or absence of words in the input audio. We then devise appropriate loss functions and train our model to perform three tasks simultaneously on streaming audio: detection, classification, and localization. We experiment with LibriSpeech and train a model to localize 1000 words. Compared to SpeechYolo, our model is 94% smaller, is capable of processing streaming audio, and achieves 6.5% better F1 score.

## 2. RELATED WORK

The task of detecting limited vocabulary has been studied in the keyword spotting literature [5, 6, 7], where the goal is to detect if a keyword is uttered in a segmented audio. Our problem definition is more challenging as we expect our model to precisely localize keywords in a streaming audio. In addition, the majority of existing literature target a small vocabulary, e.g., 12-36 words in Google Speech Commands [12], whereas we show scalability to 1000 words in our experiments.

Post processing approaches can be utilized for word localization. In DNN-HMM models [13], for instance, the sequence of predicted phonemes can be traced to find word boundaries. Another example is [14] where word scores are generated by a CNN model for streaming audio segments, and the location of the words can be estimated accordingly. Since the above models are not directly trained for localization, their localization performance is likely not optimal. A more precise localization can be obtained by forced alignment after transcribing the audio using an ASR model [15], or by coupling an ASR model with a CTC decoder [16]. However, these solutions are computationally expensive.

\* contributed when employed by Apple



**Fig. 1:** high-level overview of our word localization model.

Recently, principles from Yolo object detection [9] have been used for speech processing, which incorporate both detection and accurate localization into the model training phase. YOHO [17] can localize audio categories, e.g., it can distinguish music from speech. A more related problem definition to ours is studied in SpeechYolo [1], which shows great potential of vision-based localization techniques by localizing words in segments of 1-second audio. In this paper, we present a more carefully designed word localizer that is capable of processing streaming audio. Additionally, we show that our design yields better detection and localization accuracy with smaller memory footprint.

### 3. PROBLEM FORMULATION

We aim to detect words in the lexicon of  $c$ -words, indexed by  $w \in \{1, \dots, c\}$ . Let  $\{(w_1, b_1, e_1), \dots, (w_k, b_k, e_k)\}$  be ground-truth events in an utterance  $X \in \mathbb{R}^T$ . Each event  $(w, b, e)$  contains the word label  $w \in \{1 \dots c\}$ , the event beginning time  $b$ , and the event ending time  $e$ . Our goal is to train  $f$  parameterized by  $\theta$  that predicts proposals:

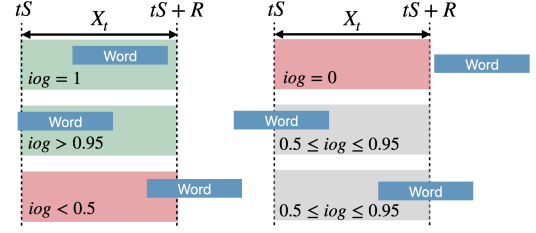
$$f(\theta, X) = \{(\hat{w}_k, \hat{b}_k, \hat{e}_k)\}_{k=1}^{k'}. \quad (1)$$

The recognizer model  $f(\theta, X)$  should be trained such that the  $k'$  predicted events match the  $k$  ground truth events.

### 4. METHODOLOGY

The overall flow of our word localization is shown in Figure 1. The CNN backbone converts the input audio into a feature matrix  $z$ . The rest of the modules use this feature matrix to detect events (encoded by  $\hat{y}$ ), classify event types (encoded by  $\hat{s}$ ), and predict the event offset  $\hat{o}$  and length  $\hat{l}$ . Finally,  $(\hat{s}, \hat{o}, \hat{l})$  are processed by the utterance processing module and events are proposed as  $\{(\hat{w}_i, \hat{b}_i, \hat{e}_i)\}_{i=1}^{k'}$ .

**Backbone.** The CNN model receives an utterance  $X \in \mathbb{R}^T$  and converts it into a feature matrix  $z \in \mathbb{R}^{n \times d}$ . The rows of  $z$  correspond to utterance segments  $X[tS : tS + R]$ , which we denote by  $X_t$  in the remainder of the paper for notation simplicity. Here,  $R$  is the length of each segment (a.k.a., the network’s receptive field) and  $S$  is the shift between two consecutive segments (a.k.a. the network’s stride). In an utterance of length  $T$ , there are  $n = \lfloor \frac{T-R}{S} \rfloor + 1$  total segments.



**Fig. 2:** Examples of positive (green highlight), negative (red highlight), and “don’t care” samples (grey highlight).

**Event detection.** The ground-truth event detection label is a binary matrix  $y_{n \times c}$ , where  $y_{t,w}$  specifies whether  $X_t$  contains the  $w$ -th word in the lexicon. To assign these hard-labels, we compute the intersection over ground-truth (iog) metric. Let  $(w, b, e)$  be a ground-truth event. We compute:

$$iog_{t,w} = \frac{\text{overlap}[(tS, tS + R), (b, e)]}{e - b} \quad (2)$$

We threshold the  $iog$  metric and assign labels accordingly:

- if  $iog_{t,w} > 0.95$ , word  $w$  is almost perfectly contained in  $X_t$ ; In this case we have  $y_{t,w} = 1$ .
- if  $iog_{t,w} < 0.5$ , word  $w$  is not contained in  $X_t$ ; In this case we have  $y_{t,w} = 0$ .
- if  $0.5 \leq iog_{t,w} \leq 0.95$ , word  $w$  is partially contained in and partially outside  $X_t$ . The classification label is “don’t care” in this case.

Figure 2 illustrates several examples for the above cases.

**Event detection loss.** We compute event probabilities as  $\hat{y}_{n \times c} = \text{sigmoid}(W^{\text{detection}} \cdot z)$ . Our goal here is to encode the presence of word  $w$  in  $X_t$  as  $\hat{y}_{t,w}$ . If  $X_t$  contains word  $w$ ,  $\hat{y}_{t,w}$  should be close to 1. To enforce this behaviour, we define the positive loss:

$$L_{pos}(y, \hat{y}) = \frac{\sum_{t=1}^n \sum_{w=1}^c BCE(y_{t,w}, \hat{y}_{t,w}) \cdot I(y_{t,w}, 1)}{\sum_{t=1}^n \sum_{w=1}^c I(y_{t,w}, 1)} \quad (3)$$

where the numerator is the sum of binary cross entropy (BCE) loss over all elements with a ground truth label of 1, and the denominator is a normalizer that counts the number of positive labels. When word  $w$  is not in  $X_t$ ,  $\hat{y}_{t,w}$  should be close to 0. To enforce this behaviour, we define the negative loss:

$$L_{neg}(y, \hat{y}) = \frac{\sum_{t=1}^n \sum_{w=1}^c BCE(y_{t,w}, \hat{y}_{t,w}) \cdot I(y_{t,w}, 0)}{\sum_{t=1}^n \sum_{w=1}^c I(y_{t,w}, 0)} \quad (4)$$

**Localization loss.** to predict event begin and end times, we adopt the CenterNet approach from visual object detection literature [11]. Let  $c_t = \frac{2tS+R}{2S}$  be the center of  $X_t$ , and  $(b, e)$  be the beginning and end of an event. We define the offset of the event as  $o = \frac{b+e}{2S} - c_t$ , and the length of the event as  $l = \frac{e-b}{R}$ . If our model can predict  $o$  and  $l$  accurately, we can calculate  $b, e$ . We generate offset prediction  $\hat{o}_{n \times c} = W^{\text{offset}} \cdot z$

and length prediction  $\hat{l}_{n \times c} = W^{length} \cdot z$ . During training we minimize the normalized L1 distance between the ground-truth and predicted values. Equation 5 shows the offset loss function  $L_o$ . The length loss function  $L_l$  is defined similarly.

$$L_o(o, \hat{o}, y) = \frac{\sum_{t=1}^n \sum_{w=1}^c |o_{t,w} - \hat{o}_{t,w}| \cdot I(y_{t,w}, 1)}{\sum_{t=1}^n \sum_{w=1}^c I(y_{t,w}, 1)} \quad (5)$$

The predictions  $(\hat{y}, \hat{o}, \hat{l})$  defined up to here can be used to generate region proposals for possible events. We initially applied non-maximum suppression (NMS), a technique widely used in image object localization [18], to select the best non-overlapping proposals but many proposals were wrong. Our investigations unveiled two reasons for this matter:

- **collision:**  $X_t$  may contain multiple words, thus,  $\hat{y}_{t,w}$  might be non-zero for multiple  $w$ . This makes it unlikely for NMS to make a correct proposal.
- **confusion:** even if there is only a single event within  $X_t$ , it might get confused, e.g., our model might confuse a ground-truth word “two” with “too” or “to”.

**Classifier.** To address the collision and confusion issues stated above, we propose to train a classifier that predicts:

$$\hat{s}_{n \times c+1} = \text{softmax}([W^{classifier} \cdot z] \otimes \text{mask}(\hat{y})) \quad (6)$$

Here,  $\hat{s}_{t,w}$  predicts the probability that the  $X_t$ : **(a)** contains the  $w$ -th word in the lexicon, for  $w \leq c$ , or **(b)** does not contain any of the words in the lexicon, for  $w = c + 1$ . The  $\otimes$  operator in Equation 6 denotes element-wise multiplication, and  $\text{mask}(\hat{y})$  is a binary tensor:

$$\text{mask}_{t,w} = \begin{cases} 1 & \text{if } \hat{y}_{t,w} \geq 0.5 \text{ or } w = c + 1 \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

With the above masking formulation, our Softmax layer does not need to solve a  $c + 1$  classification problem. Instead, it solves the problem only for the positive classes proposed by  $\hat{y}$  and the negative class. We train the softmax classifier using the cross-entropy loss:

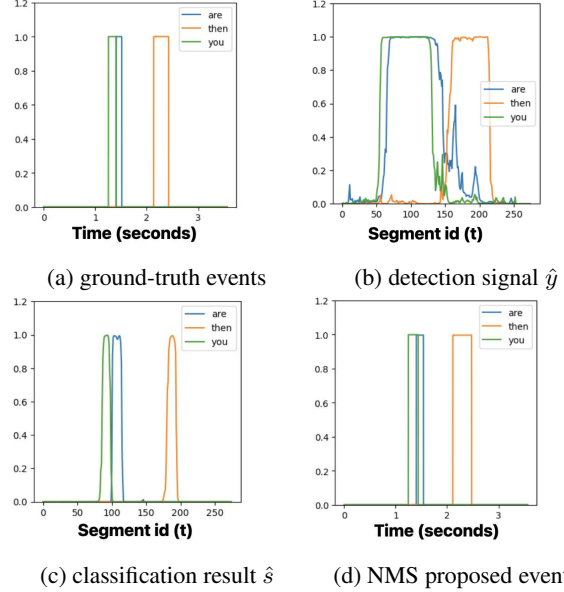
$$L_s(s, \hat{s}) = \frac{1}{n} \sum_{t=1}^n CE(s_t, \hat{s}_t) \quad (8)$$

where  $s$  is the ground-truth label and  $CE$  is the cross-entropy loss. The ground truth  $s_t \in \mathbb{R}^{c+1}$  is a one-hot vector that indexes either one of the  $c$  word classes or the negative class. In cases that  $X_t$  contains more than one ground-truth event,  $s_t$  indexes the event with smallest offset  $o$ . In essence, the Softmax classifier either rejects the events suggested by  $\hat{y}$  or selects one event (which has maximum probability) from them.

**Total loss.** The training loss is given by:

$$L = L_{pos} + L_{neg} + L_o + L_l + L_s \quad (9)$$

**Utterance processing (inference only).** Once the model is trained, we can query it with an audio and receive proposed events. Figure 3-(a) shows the ground-truth events for an example audio, and the non-zero columns of the predicted  $\hat{y}$  are



**Fig. 3:** Example input and corresponding outputs.

plotted in Figure 3-(b). The detection  $\hat{y}$  is then used in the classifier (Eq. 6) to compute  $\hat{s}$ , illustrated in Figure 3-(c). At frame  $t$ , if the maximum prediction  $\max(\hat{s}_t)$  is larger than some threshold  $\lambda$ , an event is proposed as:

$$\{event : (\arg \max(\hat{s}_t), \hat{b}_t, \hat{e}_t), score : \max(\hat{s}_t)\} \quad (10)$$

$$\hat{b}_t = S(c_t + \hat{o}_t) - \frac{\hat{l}_t}{2}R, \quad \hat{e}_t = \hat{b}_t + \hat{l}_tR$$

where  $\hat{b}_t$  and  $\hat{e}_t$  are the estimated event begin and event end times. The extracted events are a set of overlapping windows in the time domain, each of which has a score. To suppress repetitive proposals, we use NMS. The NMS output is illustrated in Figure 3-(d); as seen, the proposed events are quite accurate compared to the ground-truth events.

## 5. EVALUATIONS

**Architecture.** We use the BCResNet architecture [19] to implement the backbone in Figure 1. The input to the model is raw audio sampled at  $16KHz$ . The receptive field is  $R = 13200$  (825 ms) and stride is  $S = 160$  (10 ms). The layers and output dimensions per layer are shown in Table 1 when the raw audio length is equal to  $R$ . The entire CNN backbone encodes each segment  $X_t$  as a 128-dimensional vector.

**Dataset.** Similar to prior work [1], we use the Montreal Forced Aligner to extract the start and end time of words in LibriSpeech dataset. The lexicon is similarly chosen as 1000 words that appear the most in the training set.

**Training.** We train the network for 100 epochs with the Adam optimizer [20] and a Cosine Annealing scheduler that starts from a learning rate of 0.001 and gradually reduces it to 0.0001 by the end of training. During training, we randomly cut a small portion from the beginning of each utterance to

**Table 1:** Model architecture used in our experiments. Here, “dim” is layer output dimension (ignoring batch dimension, assuming the input audio length is equal to the receptive field), “k” is kernels size, “s” is stride, and “d” is dilation.

	Module	dim	k	s	d	
<b>Input</b>	-	[13200]	-	-	-	
<b>Feature Extractor</b>	FBank	[1, 40, 81]	400	160	1	
<b>BC-ResNet</b>	Conv2d	[256, 20, 77]	(5, 5)	(2,1)	(1,1)	
	Transition	[128, 20, 75]	3	(1,1)	(1,1)	
	Normal	[128, 20, 73]	3	(1,1)	(1,1)	
	Transition	[192, 10, 69]	3	(2,1)	(1,2)	
	Normal	[192, 10, 65]	3	(1,1)	(1,2)	
	Transition	[256, 5, 57]	3	(2,1)	(1,4)	
	Normal	[256, 5, 49]	3	(1,1)	(1,4)	
	Normal	[256, 5, 41]	3	(1,1)	(1,4)	
	Normal	[256, 5, 33]	3	(1,1)	(1,4)	
	Transition	[320, 5, 17]	3	(1,1)	(1,8)	
	Normal	[320, 5, 1]	3	(1,1)	(1,8)	
	Conv2d	[128, 1, 1]	(5, 1)	(1,1)	(1,1)	
	<b>Feature (z)</b>	-	[1, 128]	-	-	-

make the model agnostic to shifts. No other data augmentation is applied during training. The batch size in our experiments is 32 per GPU, and we use Pytorch distributed data parallelism to train our model on 8 NVIDIA-V100 GPUs.

### 5.1. Word detection and localization

For evaluation, we run NMS on events that have  $score \geq \lambda$  to obtain the proposed events. We then compute true positives (TP) as the number of proposed events that overlap with a ground-truth event of the same class. If a ground-truth event is not predicted, we count it as a false negative (FN). If a proposed event is not in the ground truth events or the predicted class is wrong we count the event as a false positive. We then compute precision, recall, F1-score, actual accuracy, and average IOU the same way SpeechYolo does [1]. Table 2 summarizes the performance comparison between our method and SpeechYolo. Here, “Ours-L” represents the model in Table 1 and “Ours-S” is the same network with half number of features per layer. We report the SpeechYolo performance numbers from the original paper<sup>1</sup>. Our large model is 94% smaller than SpeechYolo. It can process arbitrary-long audios compared to 1-second segments in SpeechYolo. Our F1-score, IOU, and actual accuracy are also consistently higher.

The better performance of our work is due to the fact that we customize the underlying detector design (Figure 1) specifically to overcome challenges in acoustic modeling. To illustrate the effect of different aspects of our design, we perform an ablation study with the large model in Table 3. In summary,  $\hat{s}$  improves precision,  $\hat{y}$  improves recall, and  $(\hat{o}, \hat{l})$  improve the localization capability (IOU).

<sup>1</sup>We were not able to reproduce their results as the authors do not provide pre-processed inputs in their github repository

**Table 2:** Comparison of our trained models with SpeechYolo. In each column, metrics that outperform SpeechYolo are bold. The decision threshold  $\lambda$  is separately tuned for SpeechYolo and our work ( $\lambda = 0.95$ ) to maximize the F1-scores.

Data	Method	Model Size	Prec.	Recall	F1	Actual Acc.	IOU
test_clean	SY [1]	108 MB	0.836	0.779	0.807	0.774	0.843
	Ours-L	<b>6.2 MB</b>	<b>0.863</b>	<b>0.880</b>	<b>0.872</b>	<b>0.873</b>	<b>0.857</b>
	Ours-S	<b>2.1 MB</b>	<b>0.852</b>	0.770	<b>0.809</b>	0.759	<b>0.855</b>
test_other	SY [1]	108 MB	0.697	0.553	0.617	-	-
	Ours-L	<b>6.2 MB</b>	<b>0.764</b>	<b>0.713</b>	<b>0.738</b>	<b>0.704</b>	<b>0.850</b>
	Ours-S	<b>2.1 MB</b>	<b>0.777</b>	0.549	<b>0.643</b>	<b>0.531</b>	<b>0.849</b>

**Table 3:** Effect of the predicted signals on performance.

$\hat{y}$	$\hat{s}$	$(\hat{o}, \hat{l})$	Precision	Recall	IOU
Yes	No	Yes	0.338	0.871	0.807
No	Yes	Yes	0.835	0.652	0.848
Yes	Yes	No	0.892	0.522	0.396
Yes	Yes	Yes	0.863	0.880	0.857

### 5.2. keyword spotting

In the next analysis, we compare our method in keyword spotting where the goal is to accurately identify a limited set of words, i.e., 20 words defined in Table 2 of [14] and used by speechYolo. The evacuation metric here is Term Weight Value (TWV) [21] defined as follows:

$$TWV(\lambda) = 1 - \frac{1}{K} [P_{miss}(k, \lambda) + \beta P_{FA}(k, \lambda)] \quad (11)$$

where  $k \in \{1, \dots, K\}$  refers to the keywords,  $P_{miss}(k, \lambda)$  is the probability of missing the  $k$ -th keyword, and  $P_{FA}(k, \lambda)$  is the probability of falsely detecting the  $k$ -th keyword in 1 second of audio data. Here,  $\beta = 999.9$  is a constant that severely penalizes the TWV score for high false alarm rates. Table 4 compares our keyword spotting performance with SpeechYolo, where the  $\lambda$  is tuned for each keyword such that the maximum TWV value (MTWV) is achieved. As seen, our method outperforms SpeechYolo.

**Table 4:** MTWV scores for 20-keyword spotting application. Bold numbers are the ones that outperform SpeechYolo.

Partition	SY [1]	Ours-L	Ours-S
test_clean	0.74	<b>0.80</b>	0.74
test_other	0.38	<b>0.64</b>	<b>0.53</b>

## 6. CONCLUSION

This paper proposes a solution to limited-vocabulary word detection and localization in speech data. We devise model components that ensure a high precision, recall, and localization score. We then define the loss functions required to train the model components. We showed in our experiments that, compared to existing work, our model is more accurate, smaller in size, and capable of processing arbitrary length audio.

## 7. REFERENCES

- [1] Yael Segal, Tzeviya Sylvia Fuchs, and Joseph Keshet, “Speechyolo: Detection and localization of speech objects,” *Proc. Interspeech 2019*, pp. 4210–4214, 2019.
- [2] Yu Zhang, James Qin, Daniel S Park, Wei Han, Chung-Cheng Chiu, Ruoming Pang, Quoc V Le, and Yonghui Wu, “Pushing the limits of semi-supervised learning for automatic speech recognition,” *arXiv preprint arXiv:2010.10504*, 2020.
- [3] Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng Chiu, James Qin, Ruoming Pang, and Yonghui Wu, “W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training,” in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2021, pp. 244–250.
- [4] Qiantong Xu, Alexei Baevski, Tatiana Likhomanenko, Paden Tomasello, Alexis Conneau, Ronan Collobert, Gabriel Synnaeve, and Michael Auli, “Self-training and pre-training are complementary for speech recognition,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 3030–3034.
- [5] Oleg Rybakov, Natasha Kononenko, Niranjana Subrahmanya, Mirkó Visontai, and Stella Laurenzo, “Streaming keyword spotting on mobile devices,” *arXiv preprint arXiv:2005.06720*, 2020.
- [6] Raziel Alvarez and Hyun-Jin Park, “End-to-end streaming keyword spotting,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6336–6340.
- [7] Raphael Tang and Jimmy Lin, “Deep residual learning for small-footprint keyword spotting,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5484–5488.
- [8] Raphael Tang, Karun Kumar, Ji Xin, Piyush Vyas, Wenyan Li, Gefei Yang, Yajie Mao, Craig Murray, and Jimmy Lin, “Temporal early exiting for streaming speech commands recognition,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 7567–7571.
- [9] Joseph Redmon and Ali Farhadi, “Yolov3: An incremental improvement,” *arXiv*, 2018.
- [10] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg, “Ssd: Single shot multibox detector,” in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [11] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl, “Objects as points,” *arXiv preprint arXiv:1904.07850*, 2019.
- [12] Pete Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *arXiv preprint arXiv:1804.03209*, 2018.
- [13] Ashish Shrivastava, Arnav Kundu, Chandra Dhir, Devang Naik, and Oncel Tuzel, “Optimize what matters: Training dnn-hmm keyword spotting model using end metric,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 4000–4004.
- [14] Dimitri Palaz, Gabriel Synnaeve, and Ronan Collobert, “Jointly learning to locate and classify words using convolutional networks,” in *Interspeech*, 2016, pp. 2741–2745.
- [15] Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger, “Montreal forced aligner: Trainable text-speech alignment using kaldii,” 08 2017, pp. 498–502.
- [16] Alex Graves, “Connectionist temporal classification,” in *Supervised sequence labelling with recurrent neural networks*, pp. 61–93. Springer, 2012.
- [17] Satvik Venkatesh, David Moffat, and Eduardo Reck Miranda, “You only hear once: a yolo-like algorithm for audio segmentation and sound event detection,” *Applied Sciences*, vol. 12, no. 7, pp. 3293, 2022.
- [18] Alexander Neubeck and Luc Van Gool, “Efficient non-maximum suppression,” in *18th International Conference on Pattern Recognition (ICPR’06)*. IEEE, 2006, vol. 3, pp. 850–855.
- [19] Byeongeun Kim, Simyung Chang, Jinkyu Lee, and Dooyong Sung, “Broadcasted residual learning for efficient keyword spotting,” *arXiv preprint arXiv:2106.04140*, 2021.
- [20] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [21] Jonathan G Fiscus, Jerome Ajot, John S Garofolo, and George Doddington, “Results of the 2006 spoken term detection evaluation,” in *Proc. sigir*, 2007, vol. 7, pp. 51–57.