

# Speech Augmentation Based Unsupervised Learning for Keyword Spotting

Jian Luo<sup>1</sup>, Jianzong Wang<sup>1\*</sup>, Ning Cheng<sup>1</sup>, Haobin Tang<sup>1,2</sup>, Jing Xiao<sup>1</sup>

<sup>1</sup>*Ping An Technology (Shenzhen) Co., Ltd.*

<sup>2</sup>*University of Science and Technology of China*

**Abstract**—In this paper, we investigated a speech augmentation based unsupervised learning approach for keyword spotting (KWS) task. KWS is a useful speech application, yet also heavily depends on the labeled data. We designed a CNN-Attention architecture to conduct the KWS task. CNN layers focus on the local acoustic features, and attention layers model the long-time dependency. To improve the robustness of KWS model, we also proposed an unsupervised learning method. The unsupervised loss is based on the similarity between the original and augmented speech features, as well as the audio reconstructing information. Two speech augmentation methods are explored in the unsupervised learning: speed and intensity. The experiments on Google Speech Commands V2 Dataset demonstrated that our CNN-Attention model has competitive results. Moreover, the augmentation based unsupervised learning could further improve the classification accuracy of KWS task. In our experiments, with augmentation based unsupervised learning, our KWS model achieves better performance than other unsupervised methods, such as CPC, APC, and MPC.

**Index Terms**—Speech Augmentation, Unsupervised Learning, Keyword Spotting

## I. INTRODUCTION

Keyword Spotting (KWS) is a useful speech application in real-world scenarios. KWS aims at detecting a relatively small set of pre-defined keywords in an audio stream, which usually exists on the interactive agents. The KWS systems usually have two kinds of applications: Firstly, it can detect the startup commands, such as “hey Siri” or “OK, Google”, providing explicit cues for interactions. Secondly, KWS can help to detect some sensitive words to protect the privacy of the speaker. Therefore, highly accurate and robust KWS systems can be of great significance to real speech applications [1]–[3].

Recently, extensive literature research on KWS has been published [4]–[6]. As a traditional solution, keyword/filler Hidden Markov Model (HMM) has been widely applied to KWS tasks, and remains competitive results [7]. In this generative approach, an HMM model is trained for each keyword, while another HMM model is trained from not-keyword speech segments. At inference, the Viterbi decoding is required, which might be computationally expensive depending on the HMM topology. In recent years, deep learning models have gained popularity on the KWS task, which show better performance than traditional approaches. Google proposed to use Deep Neural Networks (DNN) to predict sub-keyword targets. It uses the posterior processing method to generate the final

confidence score, and outperforms the HMM-based system [8]. In contrast, Convolutional Neural Networks (CNN) is more attractive, because DNN ignores the input topology, but audio features could have a strong dependency in time or frequency domains [9]–[11]. However, there is a potential drawback that CNN might not model much contextual information. Also, Recurrent Neural Networks (RNN) with Connectionist Temporal Classification (CTC) loss was also investigated for KWS. However, the limitation of RNN is that it directly models the speech features without learning local structure between successive time series and frequency steps [12]. There are also some works that combined CNN and RNN to improve the accuracy of KWS. For example, Convolutional Recurrent Neural Networks (CRNN) and Gated Convolutional Long Short-Term Memory (LSTM), achieved better performance than that of only using CNN or RNN [13]. In recent years, many researchers focus on the transformer-based models with self-attention mechanism. As a typical model, Bidirectional Encoder Representations from Transformer (BERT) has been proven to be an effective model in many Natural Language Processing (NLP) tasks [14]–[16]. The transformer-based models have also obtained much application in Automatic Speech Recognition (ASR) tasks [17], [18]. In this work, we introduced transformer to the network architecture of KWS. We think that transformer encoder has great advantage on the speech representation, and established a CNN-Attention based network to deal with the KWS task. The CNN helps network to learn the local feature, and the self-attention mechanism of transformer focuses on the long-time information.

The above supervised approaches have acquired good performance, but these models require a lot of labeled datasets. Obviously, for KWS task, the negative samples could be more procurable than positive samples, meaning that the positive samples might not be obtained easily. Especially when the keyword changes, it requires much time to collect the positive target samples, and the existing models might not easily transfer to other KWS models. In this paper, we focus on the unsupervised learning approach to alleviate this problem. The unsupervised learning mechanism allows the neural network to be trained on unlabeled datasets. With unsupervised learning, the performance of downstream task could be improved with limited labeled datasets. Unsupervised learning has made great success in the audio, image and text tasks [19]. In speech area, researchers also proposed some unsupervised pre-training algorithms [20]–[22]. Contrastive Predictive Coding

\*Corresponding author: Jianzong Wang, jzwang@188.com

(CPC) is one of those unsupervised approaches, extracting speech representation by predicting future information [23]. Apart from CPC, the Autoregressive Predictive Coding (APC) is another pre-training model, which also gets comparable results on phoneme classification and speaker verification tasks [24]. Meanwhile, Masked Predictive Coding (MPC) designs a Masked Language Model (MLM) objective in the unsupervised pre-training, and enables the model to incorporate context from both directions [25]. Based on these unsupervised learning methods, lots of unlabeled audio data can be used to obtain a better audio representation and this representation can be applied to the follow-up tasks through fine-tuning mechanism. For a robust KWS system, it should deal with different styles of speech in real-world applications. Speed and volume are major variations of the speech style. Unlike traditional unsupervised learning focuses on the general audio representation, we proposed an augmentation based approach. Our approach is to improve the model performance on KWS task with different speed and intensity situations. We designed an unsupervised loss based on the distance between the original and augmented speech, as well as the audio reconstructing information for auxiliary training. We think that speech utterances with the same keyword but at different speeds or volumes should have similar high-level feature representations for KWS tasks.

This paper investigated unsupervised speech representative methods to conduct KWS task. The unsupervised learning methods could utilize a lot of unlabeled audio datasets to improve the performance of downstream KWS task when labeled data are limited. In addition, speech augmentation based unsupervised representation might help the network to learn the speech information in various speech styles, and get a more robust performance. In summary, our major contributions of this work are the followings:

- Propose a CNN-Attention architecture for keyword spotting task, having competitive results on Google Speech Commands V2 Dataset.
- Design an unsupervised loss based on the Mean Square Error (MSE) to measure the distance between the original and augmented speech.
- Define a speech augmentation based unsupervised learning approach, utilizing the similarity between the bottleneck layer feature, as well as the audio reconstructing information for auxiliary training.

The rest of the paper is organized as follows. Sec. II highlights the related prior works about data augmentation, unsupervised learning, and other methodologies of KWS tasks. Sec. III describes the proposed model architecture and augmentation based unsupervised learning loss. Sec. IV reports the experimental results compared with other pre-training methods. We also discuss relationship between pre-training steps and performance of downstream KWS tasks. In Sec. V, we conclude with the summary of the paper and future works.

## II. RELATED WORK

Data augmentation is a common strategy to enlarge the training set of speech applications, such as Automatic Speech Recognition (ASR) and Keyword Spotting (KWS). The work [26] studied the vocal tract length perturbation method to improve the performance of ASR systems. The work [27] investigated a speed-perturbation technique to change the speed of the audio signal. Noisy audio signals have been used in [28], corrupting clean training speech with noise signal, to improve the robustness of the speech recognizer. SpecAugment [29] is a spectral-domain augmentation whose effect is to mask bands of frequency and/or time axes. SpecAugment is also explored further on large scale dataset in [30]. WavAugment [31] combines pitch modification, additive noise and reverberation to increase the performance of Contrastive Predictive Coding (CPC). In this work, we apply the speed and volume perturbation in our speech augmentation method.

Although supervised learning has been the major approach in keyword spotting area, current supervised learning models require large amounts of labeled data. Those high quality labeled datasets require substantial effort and are hardly available for the less frequently used languages. For this reason, recently there has been a great surge of interest in weakly supervised solutions that use datasets with few human annotations. Noisy student training, a semi-supervised learning method was proposed to ASR [32] and later used for robust keyword spotting [33]. There also have been related researches investigating the use of unsupervised methods to perform keyword spotting [34]–[36]. [34] proposed a self-organizing speech recognizer, and minimal transcriptions are used to train a grapheme-to-sound-unit converter. [35] presented a prototype KWS system that doesn't need manually transcribed data to train the acoustic model. In [36], the authors proposed an unsupervised learning framework without transcription. A GMM model is used to label keyword samples and test utterances by Gaussian posteriorgram. After that, segmental dynamic time warping (SDTW) gives a relevant score, and ranks the score to figure out the output. The feasibility and effectiveness of these results encourage us to introduce unsupervised learning framework to the task of keyword spotting.

Google Speech Commands V2 Dataset, is a well-studied and benchmarked dataset for novel ideas in KWS. A lot of previous works perform experiments on this dataset. [37] introduced a convolutional recurrent network with attention on multiple KWS tasks. MatchboxNet [38] is a deep residual network composed from 1D time-channel separable convolution, batch-norm layers, ReLU and dropout layers. Inspired by [37] and [38], EdgeCRNN [39] was proposed, an edge-computing oriented model of acoustic feature enhancement for keyword spotting. Recently, [40] combined a triplet loss-based embedding and a variant of K-Nearest Neighbor (KNN) for classification. We also evaluated our speech augmentation based unsupervised learning method on this dataset, and compared with other unsupervised approaches, including CPC [23], APC [24] and MPC [25].

### III. PROPOSED METHOD

#### A. KWS Model Architecture

The keyword spotting task could be described as a sequence classification task. The keyword spotting network maps an input audio sequence  $X = (x_0, x_1, \dots, x_T)$  to a limited of keyword classes  $Y \in y_{1:S}$ . In which,  $T$  is the number of audio frames and  $S$  is the number of classes. Our proposed model architecture for keyword spotting is shown in Fig 1. The network contains five parts: (1) CNN Block, (2) Transformer Block, (3) Feature Selecting Layer, (4) Bottleneck Layer, and (5) Project Layer.

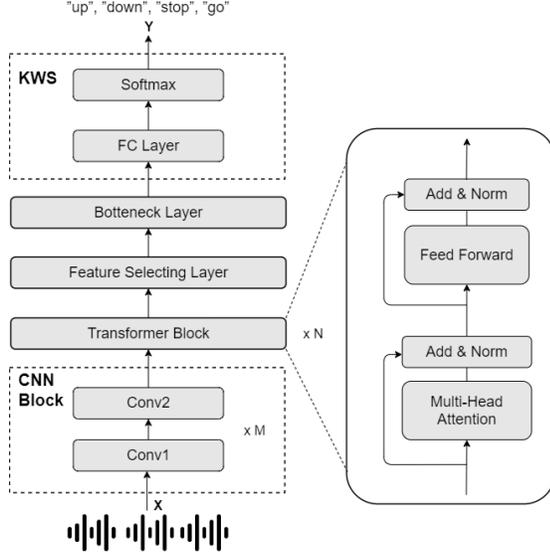


Fig. 1. The Architecture of our CNN-Attention model for keyword spotting task. The network is composed of CNN layers, self-attention layers, feature selecting layer, bottleneck layer, and project layer. In the feature selecting layer, the last few frames are selected. Finally, the project layer maps the features to predict the keyword classification.

The CNN block consists of several 2D-convolutional layers, handling the local variance on time and spectrum axes.

$$E_{cnn} = 2DConv_{\times N}(X) \quad (1)$$

In which,  $N$  is the number of convolutional layers. Then, the CNN output  $E_{cnn}$  is inputted to the transformer block, to capture long-time information with self-attention mechanism.

$$E_{tran} = SelfAttention_{\times M}(E_{cnn}) \quad (2)$$

In which,  $M$  is the number of self-attention layers. After transformer block, we designed a feature selecting layer to extract keyword information from sequence  $E_{tran}$ .

$$E_{feat} = Concat(E_{tran}[T-r, T]) \quad (3)$$

In feature selecting layer, we firstly collect last  $r$  frames of  $E_{tran}$ . And then, we concatenate all the collected frames together, into one feature vector  $E_{feat}$ . After feature selecting layer, we use a bottleneck layer and a project layer, projecting the hidden states to the predicted classification classes  $\tilde{Y}$ .

$$E_{bn} = FC_{bn}(E_{feat}) \quad (4)$$

$$\tilde{Y} = FC_{proj}(E_{bn}) \quad (5)$$

Finally, the the cross-entropy (CE) loss for supervised learning and model fine-tuning is calculated via predicted classes  $\tilde{Y}$  and ground truth classes  $Y$ .

$$\mathcal{L}_{ce} = CE(Y, \tilde{Y}) \quad (6)$$

#### B. Augmentation Method

Data augmentation are the most common used methods to promote the robustness and performance of the model in speech tasks. In this work, speed and volume based augmentation are investigated in the unsupervised learning of keyword spotting. For a given audio sequence  $X$ , we denote it as the amplitude  $A$  and time index  $t$ .

$$X = A(t) \quad (7)$$

For speed augmentation, we set a speed ratio  $\lambda_{speed}$  to adjust the speed of  $X$ .

$$X^{aug} = A(\lambda_{speed}t) \quad (8)$$

For volume augmentation, we also set an intensity ratio  $\lambda_{volume}$  to change the volume of  $X$ .

$$X^{aug} = \lambda_{volume}A(t) \quad (9)$$

With different ratios  $\lambda_{speed}$  and  $\lambda_{volume}$ , we could obtain multiple speech sequence pairs  $(X, X^{aug})$ , to train the audio representation network with unsupervised learning. We think that speech utterances at different speed or volume should have similar high-level feature representation for KWS tasks.

#### C. Unsupervised Learning Loss

The overall architecture of augmentation based unsupervised learning is shown in Fig 2. Similar to other unsupervised methods, the proposed approach also consists of two stages: (1) pre-training on unsupervised data, and (2) fine-tuning on supervised KWS data. In the pre-training stage, the bottleneck feature was obtained through training the unlabeled speech. In fine-tuning stage, the extracted bottleneck features are used for KWS prediction.

In the pre-training stage, the pair speech data  $(X, X^{aug})$  are inputted into the CNN-Attention models respectively, but with the same model parameters. Because  $X^{aug}$  comes from  $X$ , our designed unsupervised methods expect that  $X$  and  $X^{aug}$  will output similar high-level bottleneck features. It means that no matter how fast or how loud a speaker says, the content of the speech is the same. Thus, the optimization of network needs to reflect the similarity of  $X$  and  $X^{aug}$ . We choose the Mean Square Error (MSE)  $\mathcal{L}_{sim}$  to measure the distance between the output of  $X$  and  $X^{aug}$ .

$$\mathcal{L}_{sim} = \frac{1}{U_{bn}} \sum_{u=0}^{U_{bn}} |E_{bn}(u) - E_{bn}^{aug}(u)|^2 \quad (10)$$

Where  $U_{bn}$  represents the dimension of the bottleneck feature vector.  $E_{bn}$  and  $E_{bn}^{aug}$  are the output of bottleneck layer of original speech  $X$  and augmented speech  $X^{aug}$  respectively.

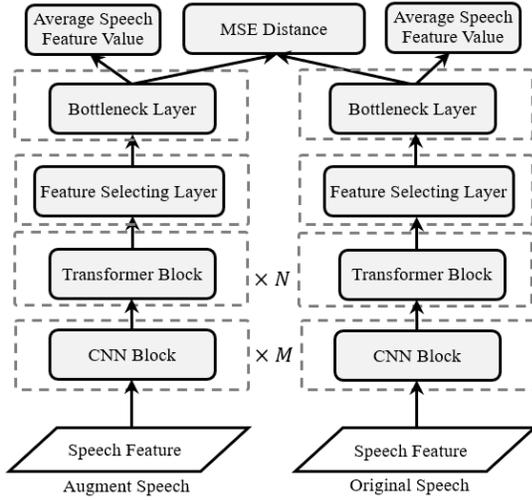


Fig. 2. The proposed speech augmentation based audio unsupervised learning method. In the pre-training stage, the pair of original and augmented speech will be inputted into the network separately but with the same model parameters. The network will output the average speech feature values and the bottleneck feature. The two bottleneck features are calculated by MSE loss, since the augmented and original speech should output similar high-level features for keyword spotting.

In addition, the designed network has another branch for auxiliary training, which predicts the average feature of the input speech segment. This branch guides the network to learn the intrinsic feature of the speech utterance. We firstly compute the average vector of the input Fbank vector  $X$  alongside the time axis  $t$ . Then, we use another reconstructing layer attached to the bottleneck layer, to reconstruct the average Fbank vector  $\tilde{X}$ . We also use MSE loss  $\mathcal{L}_x$  to calculate the similarity between these two audio vectors alongside the feature dimension  $U_x$ .

$$\begin{aligned} \mathcal{X} &= \frac{1}{T} \sum_T (X) \\ \tilde{X} &= \text{FC}_{\text{reconstruct}}(E_{bn}) \\ \mathcal{L}_x &= \frac{1}{U_x} \sum_{u=0}^{U_x} |\mathcal{X}(u) - \tilde{X}(u)|^2 \end{aligned} \quad (11)$$

In which,  $U_x$  represents the dimension of Fbank feature vector, and  $\mathcal{X}$  denotes the average vector of  $X$ . Similarly, the loss  $\mathcal{L}_x^{aug}$  between the augmented average audio  $\mathcal{X}^{aug}$  and

ured feature  $\tilde{X}^{aug}$  could be defined as follows:

$$\mathcal{L}_x^{aug} = \frac{1}{U_x} \sum_{u=0}^{U_x} |\mathcal{X}^{aug}(u) - \tilde{X}^{aug}(u)|^2 \quad (12)$$

Therefore, the final loss function  $\mathcal{L}_{ul}$  of the unsupervised learning (UL) consists of the above three losses  $\mathcal{L}_{sim}$ ,  $\mathcal{L}_x$ , and  $\mathcal{L}_x^{aug}$ .

$$\mathcal{L}_{ul} = \lambda_1 \mathcal{L}_{sim} + \lambda_2 \mathcal{L}_x + \lambda_3 \mathcal{L}_x^{aug} \quad (13)$$

Where  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$  are factor ratio of each loss component.

In fine-tuning stage, the branch of average feature prediction is removed. A project layer and a softmax layer are added after the bottleneck layer to make the KWS prediction. In the fine-tuning, the parameters of original network could be fixed or updated. In our experiments, we found that updating all the parameters could help to improve the performance. Thus, we choose to update all parameters in the fine-tuning stage.

## IV. EXPERIMENTS

In this section, we evaluated the proposed method in keyword spotting tasks. We implemented our CNN-Attention model with supervised training and compared it with Google's model. We also made an ablation study, to explore the effect of speed and volume augmentation on unsupervised learning. What's more, other unsupervised learning methods are compared with our approach, including CPC, APC, MPC. When implementing these approaches, we used the network and hyperparameters in their publications, but all experimental tricks were not leveraged [23]–[25]. We also discuss the impact of different pre-training steps on the performance and convergence of downstream KWS task.

### A. Datasets

We used Google's Speech Commands V2 Dataset [41] for evaluating the proposed models. The dataset contains about 106000 one-second or more long utterances. Total 30 short words were recorded by thousands of different people, as well as background noise such as pink noise, white noise, and human-made sounds. The KWS task is to discriminate among 12 classes: "yes", "no", "up", "down", "left", "right", "on", "off", "stop", "go", unknown, or silence. The dataset was split into training, validation, and test sets, with 80% training, 10% validation, and 10% test. This results in about 37000 samples for training, and 4600 each for validation and testing. We

TABLE I  
MODEL CONFIGURATIONS

Unit Name	Hyperparameters
#CNN Blocks	$M = 2$ layers, $3 \times 3$ kernel, $2 \times 2$ stride, 32 channels
#Transformer Block	$N = 2$ layers, dimension = 320, 4 head, feedforward = 1024
#Feature Selecting Layer	Last $r = 2$ frames, $2 \times 320$ dimension
#Bottleneck Layer	one FC layer, 800 dimension
#Project Layer	one FC layer, 12 dimension softmax
#Reconstruct Layer	one FC layer, 40 dimension softmax
#Factor Ratio	$\lambda_1 = 0.9$ , $\lambda_2 = 0.05$ , $\lambda_3 = 0.05$

TABLE II  
RESULTS COMPARISON OF KWS MODEL, CLASSIFICATION ACCURACY (%)

Model Name	Supervised Training Data	Dev	Eval
Sainath and Parada (Google)	Speech Commands	-	84.7
CNN-Attention (ours)	Speech Commands	86.4	85.3
<b>CNN-Attention + volume &amp; speed augment (ours)</b>	Speech Commands	<b>87</b>	<b>85.7</b>

TABLE III  
ABLATION STUDY, THE EFFECT OF SPEED AND VOLUME AUGMENTATION, CLASSIFICATION ACCURACY (%)

Model Name	Pre-training Data	Fine-tuning Data	Dev	Eval
CNN-Attention + volume pre-training	Speech Commands	Speech Commands	86.1	85.9
CNN-Attention + speed pre-training	Speech Commands	Speech Commands	87.8	86.9
<b>CNN-Attention + volume &amp; speed pre-training</b>	Speech Commands	Speech Commands	<b>87.9</b>	<b>87.2</b>
CNN-Attention + volume pre-training	Librispeech-100	Speech Commands	86.3	86.0
CNN-Attention + speed pre-training	Librispeech-100	Speech Commands	87.9	87.9
<b>CNN-Attention + volume &amp; speed pre-training</b>	Librispeech-100	Speech Commands	<b>88.2</b>	<b>88.1</b>

used the real noisy data HuNonspeech<sup>1</sup> to corrupt the original speech. In the experiments, the Aurora4 tools were used to implement this strategy<sup>2</sup>. Each utterance will be randomly corrupted by public 100 kinds of noise in HuNonspeech. Each utterance has a level of 0-20dB Signal Noise Ratio (SNR), and all datasets have an average 10dB SNR.

Similar to other unsupervised methods, a large unlabeled corpus, 100 hours of Librispeech [42] clean speech were also leveraged to pre-train the network by unsupervised learning. Firstly, the long utterances were split up into 1 second segments, keeping consistent with Speech Commands datasets. Nextly, the clean segments were also mixed with noisy HuNonspeech data by Aurora 4 tools, and the corrupted mechanism was as same as the Speech Commands.

### B. Experimental Setups

The acoustic features were 40-dimensional log-mel filter-bank with 30ms frame length and 10ms frame shift. The detailed hyperparameters of our proposed network were shown in Table I. For training the KWS model, all of the matrix weights are initialized with random uniform initialization, and the bias parameters are initialized with the constant value 0.1. In our experiments, we trained all the networks with Adam optimizer for 30k steps with a batchsize 200 until the loss becomes little change. In addition, the factor ratios of loss  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are set to 0.9, 0.05, 0.05 respectively.

To demonstrate the effectiveness of our proposed model, we investigated several other approaches for comparison. For supervised learning, we used Sainath and Parada’s model by Google [43] as the baseline model. The Google blog post released the Sainath and Parada’s model implemented by TensorFlow. For unsupervised learning, we compared our method with other pre-training models:

- **Contrastive Predictive Coding (CPC) [23]:** Through an unsupervised mechanism by utilizing next step prediction, CPC learns representations from high-dimensional signal.

The CPC network mainly contains a non-linear encoder and an autoregressive decoder. An input sequence is embedded to a latent space, producing a context representation. Targeting at predicting future observations, the density ratio is established to maximize the mutual information between future observations and current context representation.

- **Autoregressive Predictive Coding (APC) [24]:** APC also belongs to the family of predictive models. APC directly optimizes L1 loss between input sequence and output sequence. APC has proved an effective method in recent language model pre-training task and speech representation.
- **Masked Predictive Coding (MPC) [25]:** Inspired by BERT, MPC uses Masked Language Model (MLM) structure to perform predictive coding on Transformer based models. Similar to BERT, 15% of feature frames in each utterance are chosen to be masked during the pre-training procedure. Among these chosen frames, 80% are replaced with zero vectors, 10% are replaced with random positions, and the rest remain unchanged. L1 loss is computed between masked input features and encoder output at corresponding position. Dynamic masking was also adopted where the masking pattern is generated when a sequence is fed into the model.

### C. Results

Table II lists the experimental results of supervised learning with Speech Commands dataset. We firstly implemented the Google’s Sainath and Parada model by the original TensorFlow recipes, achieving the accuracy of 84.7%. Secondly, our CNN-Attention model is implemented by supervised loss  $\mathcal{L}_{ce}$  without any augmented data and achieved 0.6% higher accuracy than Google’s model. It is proved that our designed CNN-Attention architecture is effective for KWS task. Finally, after adding speed and volume augmentation to speech, we got a higher accuracy. It corresponds with the existing research that augmented dataset is helpful for improving the performance

<sup>1</sup><http://web.cse.ohio-state.edu/pnl/corpus/HuNonspeech/>

<sup>2</sup><http://aurora.hsnr.de/index-2.html>

TABLE IV  
 COMPARED WITH OTHER UNSUPERVISED LEARNING METHODS, CLASSIFICATION ACCURACY (%)

Model Name	Pre-training Data	Fine-tuning Data	Dev	Eval
Contrastive Predictive Coding (CPC) [23]	Speech Commands	Speech Commands	87.6	86.9
Autoregressive Predictive Coding (APC) [24]	Speech Commands	Speech Commands	87.2	86.5
Masked Predictive Coding (MPC) [25]	Speech Commands	Speech Commands	87.0	86.7
<b>CNN-Attention + volume &amp; speed pre-training (ours)</b>	Speech Commands	Speech Commands	<b>87.9</b>	<b>87.2</b>
Contrastive Predictive Coding (CPC) [23]	Librispeech-100	Speech Commands	87.8	87.4
Autoregressive Predictive Coding (APC) [24]	Librispeech-100	Speech Commands	87.7	87.5
Masked Predictive Coding (MPC) [25]	Librispeech-100	Speech Commands	87.9	87.0
<b>CNN-Attention + volume &amp; speed pre-training (ours)</b>	Librispeech-100	Speech Commands	<b>88.2</b>	<b>88.1</b>

of the model. It also inspires our motivation for building augmentation based unsupervised learning methods.

To analyze the effect of speed and volume augmentation on unsupervised learning, we also made an ablation study in our experiments. The experimental results are shown in Table III. The volume pre-training model means that the augmented speech pairs  $(X, X^{aug})$  only contain the intensity augment data. Meanwhile, the speed pre-training model is trained only by speed augmented pairs. For better investigation, we pre-trained the model with two datasets by unsupervised learning loss  $\mathcal{L}_{ul}$ . The results indicate that speed augmented unsupervised learning has better performance than intensity based augmented pre-training. With both volume and speed augmentation, we could achieve better classification accuracy than only with single augmentation method. In addition, large datasets pre-training (Librispeech-100) results in better performance than small datasets (Speech Commands). Our proposed augmentation based unsupervised method (Eval 87.2% in Table III) also promotes the accuracy of adding augmentation to supervised training (Eval 85.7% in Table II) even with the same training data.

After that, we established the CPC, APC, MPC and made the comparison with these unsupervised learning methods. As depicted in Table IV, CPC achieves better performance than APC and MPC. Our augmentation based approach outperforms all of the other unsupervised methods on both two pre-training datasets (Speech Commands and Librispeech-100). The comparison demonstrated that our proposed augmentation based unsupervised learning is capable of extracting the speech information, and is an effective approach for KWS tasks.

#### D. Pre-training Analysis

More pre-training steps usually help to improve the performance of downstream tasks. To get a better understanding of our unsupervised approach, we also conducted experiments with different pre-training steps. The 5K, 10K, 20K, 30K pre-training steps were used for making this comparison. The performance of different steps is plotted in Fig 3.

We show the model training of supervised learning with these different steps of pre-training. Our experiments demonstrated that more pre-training steps are not only helpful for achieving better performance but also making downstream KWS task converge faster. Unsupervised learning with 30K steps has the highest classification accuracy and the fastest

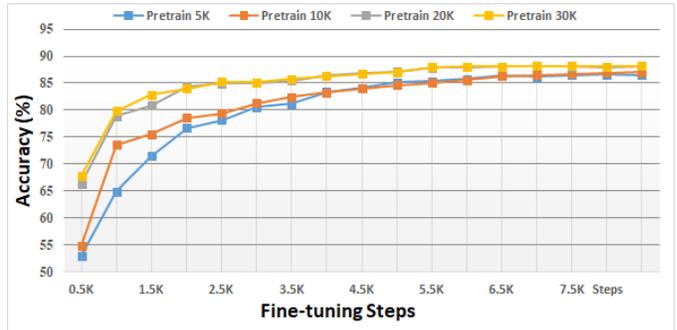


Fig. 3. The results comparison with different pre-training steps. Different pre-training steps of unsupervised learning result in different accuracy performance and fine-tuning convergence. In our experiments, pre-training 30K steps have the highest classification accuracy, and fastest convergence.

convergence. It also should be noted that the difference between 20K and 30K was very close, meaning that the pre-training steps are enough to obtain the desired performance.

#### V. CONCLUSION

This paper investigated unsupervised learning method for keyword spotting task. We designed a CNN-Attention architecture and achieved competitive results on the Speech Commands dataset. In addition, we proposed a speech augmentation based unsupervised learning approach for KWS. Our method uses speed and intensity augmentation to establish training pairs, and pre-trains the network via the similarity loss between the speech pair and the speech reconstructed loss. In our experiments, the proposed unsupervised approach could further improve the model performance, and outperform other unsupervised methods, such as CPC, APC and MPC. We also found that more pre-training steps are not only helpful for better performance but also for faster convergence. In future works, we are interested in applying the augmentation based unsupervised learning approach to other speech tasks, such as speaker verification and speech recognition.

#### VI. ACKNOWLEDGEMENT

This paper is supported by the Key Research and Development Program of Guangdong Province under grant No. 2021B0101400003. Corresponding author is Jianzong Wang from Ping An Technology (Shenzhen) Co., Ltd (jzwang@188.com).

## REFERENCES

- [1] B. Li, T. N. Sainath, A. Narayanan, J. Caroselli, M. Bacchiani, A. Misra, I. Shafraan, H. Sak, G. Pundak, K. K. Chin *et al.*, “Acoustic modeling for google home.” in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2017.
- [2] J. Luo, J. Wang, N. Cheng, G. Jiang, and J. Xiao, “End-to-end silent speech recognition with acoustic sensing,” in *IEEE Spoken Language Technology Workshop (SLT)*, 2021.
- [3] J. Schalkwyk, D. Beeferman, F. Beaufays, B. Byrne, C. Chelba, M. Cohen, M. Kamvar, and B. Strope, “your word is my command: Google search by voice: A case study,” *Springer Advances in speech recognition*, 2010.
- [4] X. Wang, S. Sun, C. Shan, J. Hou, L. Xie, S. Li, and X. Lei, “Adversarial examples for improving end-to-end attention-based small-footprint keyword spotting,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019.
- [5] A. Rosenberg, K. Audhkhasi, A. Sethy, B. Ramabhadran, and M. Picheny, “End-to-end speech recognition and keyword search on low-resource languages,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017.
- [6] C. Shan, J. Zhang, Y. Wang, and L. Xie, “Attention-based end-to-end models for small-footprint keyword spotting,” in *arXiv preprint:1803.10916*, 2018.
- [7] M.-C. Silaghi, “Spotting subsequences matching an hmm using the average observation probability criteria with application to keyword spotting,” in *The Association for the Advancement of Artificial Intelligence (AAAI)*, 2005.
- [8] G. Chen, C. Parada, and G. Heigold, “Small-footprint keyword spotting using deep neural networks,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2014.
- [9] R. Tang and J. Lin, “Deep residual learning for small-footprint keyword spotting,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018.
- [10] J. Luo, J. Wang, N. Cheng, and J. Xiao, “Unidirectional memory-self-attention transducer for online speech recognition,” in *IEEE International Conference on Acoustics Speech and Signal Processing Proceedings (ICASSP)*, 2021.
- [11] M. Xu and X.-L. Zhang, “Depthwise Separable Convolutional ResNet with Squeeze-and-Excitation Blocks for Small-Footprint Keyword Spotting,” in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2020.
- [12] M. Sun, A. Raju, G. Tucker, S. Panchapagesan, G. Fu, A. Mandal, S. Matsoukas, N. Strom, and S. Vitaladevuni, “Max-pooling loss training of long short-term memory networks for small-footprint keyword spotting,” in *IEEE Spoken Language Technology Workshop (SLT)*, 2016.
- [13] S. O. Arik, M. Kliegl, R. Child, J. Hestness, A. Gibiansky, C. Fougner, R. Prenger, and A. Coates, “Convolutional recurrent neural networks for small-footprint keyword spotting,” in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2017.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems (NIPS)*, 2017.
- [15] X. Jia, J. Wang, Z. Zhang, N. Cheng, and J. Xiao, “Large-scale transfer learning for low-resource spoken language understanding,” in *IEEE Conference of the International Speech Communication Association (INTERSPEECH)*, 2020.
- [16] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *arXiv preprint:1810.04805*, 2018.
- [17] S. Karita, N. Chen, T. Hayashi, T. Hori, H. Inaguma, Z. Jiang, M. Someki, N. E. Y. Soplin, R. Yamamoto, X. Wang *et al.*, “A comparative study on transformer vs rnn in speech applications,” in *IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, 2019.
- [18] J. Luo, J. Wang, N. Cheng, G. Jiang, and J. Xiao, “Multi-quartznet: Multi-resolution convolution for speech recognition with multi-layer feature fusion,” in *IEEE Spoken Language Technology Workshop (SLT)*, 2021.
- [19] S. Gidaris, P. Singh, and N. Komodakis, “Unsupervised representation learning by predicting image rotations,” in *arXiv preprint:1803.07728*, 2018.
- [20] Y.-C. Wang, S. Venkataramani, and P. Smaragdis, “Self-supervised learning for speech enhancement,” in *Proceedings of the 37-th International Conference on Machine Learning (ICML)*, 2020.
- [21] J. Luo, J. Wang, N. Cheng, and J. Xiao, “Dropout regularization for self-supervised learning of transformer encoder speech representation,” in *IEEE Conference of the International Speech Communication Association (INTERSPEECH)*, 2021.
- [22] S. Schneider, A. Baevski, R. Collobert, and M. Auli, “wav2vec: Unsupervised pre-training for speech recognition,” in *arXiv preprint:1904.05862*, 2019.
- [23] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” in *arXiv preprint:1807.03748*, 2018.
- [24] Y.-A. Chung, W.-N. Hsu, H. Tang, and J. Glass, “An unsupervised autoregressive model for speech representation learning,” in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2019.
- [25] D. Jiang, X. Lei, W. Li, N. Luo, Y. Hu, W. Zou, and X. Li, “Improving transformer-based speech recognition using unsupervised pre-training,” in *arXiv preprint:1910.09932*, 2019.
- [26] N. Jaitly and G. E. Hinton, “Vocal tract length perturbation (vtlp) improves speech recognition,” in *ICML Workshop on Deep Learning for Audio, Speech and Language*, 2013.
- [27] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, “Audio augmentation for speech recognition,” in *Sixteenth annual conference of the international speech communication association (INTERSPEECH)*, 2015.
- [28] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh, S. Sengupta, A. Coates *et al.*, “Deep speech: Scaling up end-to-end speech recognition,” in *arXiv preprint:1412.5567*, 2014.
- [29] D. S. Park, W. Chan, Y. Zhang, C.-C. Chiu, B. Zoph, E. D. Cubuk, and Q. V. Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *arXiv preprint:1904.08779*, 2019.
- [30] D. S. Park, Y. Zhang, C.-C. Chiu, Y. Chen, B. Li, W. Chan, Q. V. Le, and Y. Wu, “SpecAugment on large scale datasets,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020.
- [31] E. Khairtonov, M. Rivière, G. Synnaeve, L. Wolf, P.-E. Mazaré, M. Douze, and E. Dupoux, “Data augmenting contrastive learning of speech representations in the time domain,” in *2021 IEEE Spoken Language Technology Workshop (SLT)*, 2021.
- [32] D. S. Park, Y. Zhang, Y. Jia, W. Han, C.-C. Chiu, B. Li, Y. Wu, and Q. V. Le, “Improved noisy student training for automatic speech recognition,” in *arXiv preprint:2005.09629*, 2020.
- [33] H.-J. Park, P. Zhu, I. L. Moreno, and N. Subrahmanya, “Noisy student-teacher training for robust keyword spotting,” in *arXiv preprint:2106.01604*, 2021.
- [34] A. Garcia and H. Gish, “Keyword spotting of arbitrary words using minimal speech resources,” in *IEEE International Conference on Acoustics Speech and Signal Processing Proceedings (ICASSP)*, 2006.
- [35] P. Li, J. Liang, and B. Xu, “A novel instance matching based unsupervised keyword spotting system,” in *Second International Conference on Innovative Computing, Information and Control (ICICIC)*, 2007.
- [36] Y. Zhang and J. R. Glass, “Unsupervised spoken keyword spotting via segmental dtw on gaussian posteriorgrams,” in *IEEE Workshop on Automatic Speech Recognition & Understanding (ASRU)*, 2009.
- [37] D. C. de Andrade, S. Leo, M. L. D. S. Viana, and C. Bernkopf, “A neural attention model for speech command recognition,” in *arXiv preprint:1808.08929*, 2018.
- [38] S. Majumdar and B. Ginsburg, “Matchboxnet: 1d time-channel separable convolutional neural network architecture for speech commands recognition,” in *arXiv preprint:2004.08531*, 2020.
- [39] Y. Wei, Z. Gong, S. Yang, K. Ye, and Y. Wen, “Edgcrnn: an edge-computing oriented model of acoustic feature enhancement for keyword spotting,” in *Journal of Ambient Intelligence and Humanized Computing*, 2021.
- [40] R. Vygón and N. Mikheylovskiy, “Learning efficient representations for keyword spotting with triplet loss,” in *arXiv preprint:2101.04792*, 2021.
- [41] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *arXiv preprint:1804.03209*, 2018.
- [42] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An asr corpus based on public domain audio books,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015.
- [43] T. N. Sainath and C. Parada, “Convolutional neural networks for small-footprint keyword spotting,” in *Conference of the International Speech Communication Association (INTERSPEECH)*, 2015.