



# Exploring the Impact of Back-End Network on Wav2vec 2.0 for Dialect Identification

Qibao Luo, Ruohua Zhou \*

School of Electrical and Information Engineering, Beijing University of Civil Engineering and Architecture, Beijing, China

zhouruohua@bucea.edu.cn

## Abstract

This paper explores the wav2vec 2.0 model for dialect identification, focusing on the impact of the back-end network during fine-tuning. Prior research has typically used wav2vec 2.0 as a frame-level feature extractor, followed by a simple back-end consisting of a pooling layer and a fully connected layer. In contrast, we employ multi-scale aggregation and a graph neural network to design a more sophisticated back-end that implicitly exploit phoneme sequence information and significantly improves system performance. We evaluate our system on the dialect identification task of the Oriental Language Recognition Challenge 2020 (AP20-OLR). Experimental results demonstrate that our system outperforms the state-of-the-art baseline by a relative reduction of 50% in  $C_{avg}$ . We also verify the effectiveness of our proposed back-end network, which results in a relative reduction of 54% in  $C_{avg}$ . Our findings highlight the importance of incorporating a more effective back-end network for improved dialect identification performance when using the wav2vec 2.0 model.

**Index Terms:** dialect identification, wav2vec 2.0, fine-tuning, multi-scale aggregation, graph neural network

## 1. Introduction

Dialect identification (DID) is the process of automatically determining the dialect category from a speech sample. DID is considered a more challenging task than language identification (LID) because it involves identifying different dialects within the same language family. The two most effective approaches to LID are the acoustic-phonetic approach and the phonotactic approach [1]. The acoustic-phonetic approach is based on phonetic differences between languages, and it is assumed that phonetic characteristics can be extracted from acoustic signals. Conventional acoustic-phonetic approaches typically use acoustic features and a universal-background model-based GMM (GMM-UBM) to model the acoustic-phonetic distribution of a language [2]. In addition, SVM and i-vector techniques have also been explored to improve LID performance [3, 4]. Modern DNN-based acoustic LID systems use deep network to extract utterance-level embedding for LID, achieving high performance [5, 6]. However, they typically require a large amount of labeled training data, which is usually scarce in dialect language.

Based on phonological research, each language has its own set of lexical-phonological rules that govern the combinations of different phonemes and determine permissible phone sequences. While phonemes can be shared across languages, the statistics of their sequential patterns differ greatly between languages. Therefore, using phoneme sequence information can

significantly improve LID performance. Early phonotactic approaches typically use one or more phone recognizers as the front-end and phone n-gram modeling for target languages as the back-end [7, 8]. Recently, the phoneme-aware acoustic LID systems combining acoustic features and phoneme information significantly improve LID performance [9, 10]. As they involve phoneme-related tasks such as automatic speech recognition and phoneme classification, the phonetic transcription of speech is required.

[11] proposed the phonetic and phonotactic LID (PHO-LID) method, which incorporates phonetic and phonotactic information hierarchically via a CNN-Trans encoder without the use of phoneme annotations. In the PHO-LID model, a self-supervised phoneme segmentation task and a LID task share a CNN module, which encodes both language identity and sequential phonemic information.

Inspired by [11], we propose a new method for DID that implicitly incorporates phoneme sequence information by designing a more sophisticated back-end for the wav2vec 2.0 pre-training model. Wav2vec 2.0 [12] is a recently proposed framework for self-supervised learning of representations from raw audio data. Recent studies have explored wav2vec 2.0 to improve LID performance [13, 14], typically using it as a frame-level feature extractor, followed by a simple back-end consisting of a pooling layer and a fully connected layer. While the statistics pooling layer can obtain utterance-level embedding from the frame-level features, it cannot exploit the long-term dependency that may contain phoneme sequence patterns. Therefore, we employ multi-scale aggregation to implicitly exploit phoneme sequence information from the frame-level features obtained by the wav2vec 2.0 model. Additionally, we apply a graph neural network to better fuse the temporal and spectral aggregation features. Experiments verify that the proposed back-end significantly improves system performance.

## 2. Related work

This section describes the wav2vec 2.0. Figure 1 shows the process of pre-training and fine-tuning.

### 2.1. Pre-training

Wav2vec2.0 is a transformer-based model [15]. The model consists of three sub-modules, including a feature encoder, transformer module, and quantization module. The feature encoder is a multi-layer CNN that maps the raw audio  $X_{1:L}$  to a latent speech representation  $Z_{1:N}$ . Then, the transformer module contextualizes the masked representation to generate a contextual representation  $C_{1:N}$ . Finally, the quantization module discretizes the latent speech representation  $Z_{1:N}$  into a trainable codebook  $Q_{1:N}$ .

\* Ruohua Zhou is the Corresponding author.

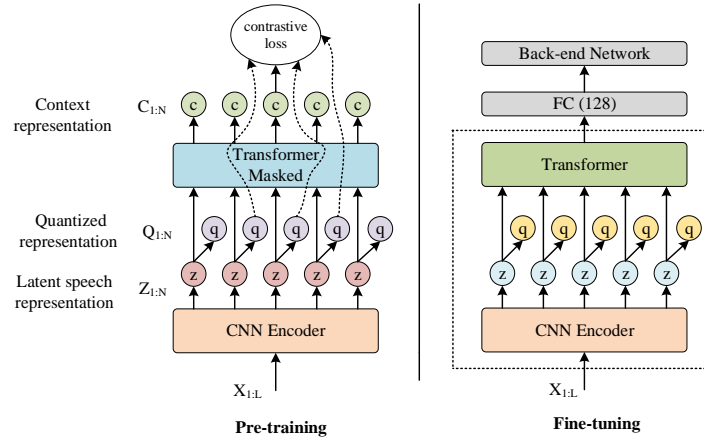


Figure 1: *Pre-training and fine-tuning process of the wav2vec 2.0 model.*

During pre-training, the model learns the representation of speech audio by solving contrastive tasks to improve the model’s performance. By calculating the contrastive loss for each masked time step  $n$ , the degree to which the target  $q_n$  is recognized from a set of distractors is measured given the corresponding contextual vector  $c_n$ .

## 2.2. Fine-tuning

The wav2vec 2.0 XLS-R (0.3B) [16] model was used for all the work described in this paper. The XLS-R pre-trained model uses data from many languages, but the target dialect we need to recognize is not included. In tasks such as speaker recognition and emotion recognition, researchers have explored fine-tuning the pre-trained model [17, 18]. The fine-tuning process is shown on the right side of Figure 1. We added a fully connected layer on top of the wav2vec encoder to reduce the representation dimension. In addition, to extract higher-level features, we added a back-end network after the fully connected layer. During fine-tuning, no masking was applied to the features, and the cross-entropy objective function was used to minimize the training loss.

## 3. Multi-scale aggregation graph neural network

In this section, we will introduce the architecture of the multi-scale aggregation graph neural network with a fine-tuned back-end. It includes the multi-scale aggregation module, the graph neural network module, and the attention aggregation layer. Phoneme sequences can provide important discriminative information for language identification. Even with the same phoneme sequence, different temporal dependencies can be presented due to changes in speaking rate. Multi-scale mechanism can capture temporal dependencies of different lengths and obtain richer phoneme sequence information. Additionally, Graph neural network can model the mutual dependency between the time and frequency domains. Attention aggregation layer is used to extract more relevant spectral and temporal representations.

An overview of the complete architecture is presented in Figure 2. In our experiments, we discovered that the best results were achieved by using a combination of Rectified Linear Units (ReLU) and batch normalization (BN) [19] in the multi-scale

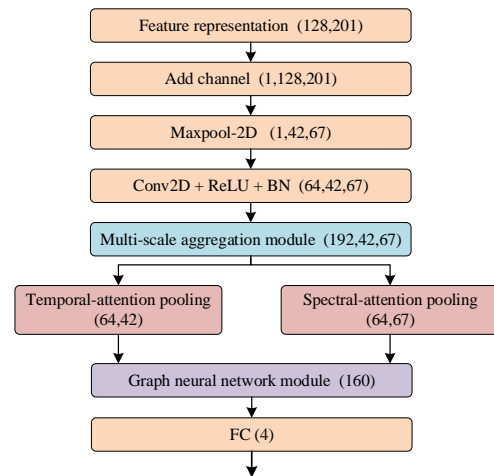


Figure 2: *Multi-scale aggregation graph neural network.*

aggregation module, and a combination of BN and Scaled Exponential Linear Units (SeLU) in other modules. The following sections will introduce the main modules of the network.

### 3.1. Multi-scale aggregation module

Our multi-scale aggregation module is derived from a partial structure of ECAPA-TDN [20], as shown in Figure 3. However, we replaced the original one-dimensional convolution with a two-dimensional convolution in this module. Specifically, the multi-scale aggregation module consists of two parts: the SE-Res2Block and the multi-scale aggregation.

#### 3.1.1. SE-Res2Block

We used the 2D Squeeze and Excitation (SE) [21] and Res2Net [22] to construct the SE-Res2Block (2D) module. In addition, we utilized dilated convolutions with different rates to increase the receptive field without reducing the size of the feature map. SE is primarily used to learn the correlations between feature channels by employing another neural network to obtain the importance of each feature channel.

Res2Block is a module that introduces hierarchical connections into residual units. It replaces the original convolu-

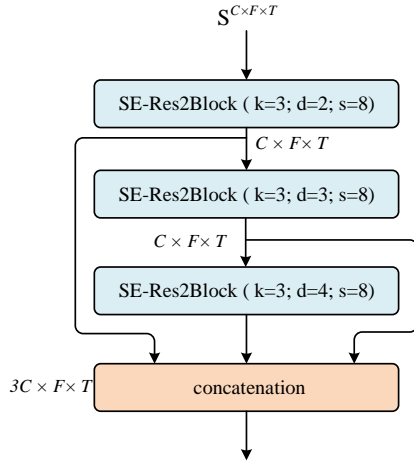


Figure 3: Multi-scale aggregation module.  $C, F,$  and  $T$  represent channel number, frequency dimension, and time dimension, respectively.  $k$  denotes the size of the convolution kernel,  $d$  denotes the dilation factor, and  $s$  denotes the group size.

tion kernel with a set of filters that include dilated convolutions to hierarchically connect different filter groups. Specifically, Res2Block first uses a  $1 \times 1$  convolution kernel to adjust the input and output channel numbers, and then divides the output features into  $s$  groups. For each feature map, except for the first group of channels, a dilated convolution operation is performed to receive information from all previous feature maps. Finally, all output features from  $s$  groups are connected along the channel dimension, and feature maps with channel attention are obtained through Conv2D layers and SE modules.

### 3.1.2. Multi-scale aggregation

According to the studies in [23, 24], aggregating features from different layers can improve the accuracy of the model in speaker verification tasks. Considering the hierarchical structure of neural networks, we concatenate the bottleneck features of all SE-Res2Block blocks and output the multi-scale aggregated features.

## 3.2. Graph neural network

According to [25], an end-to-end integrated spectrotemporal graph attention network was proposed and achieved success in the field of speech anti-spoofing. We applied the graph neural network to our model, as shown in Figure 4. Specifically, this structure includes a graph attention network, graph pooling, heterogeneous stacked graph attention layers (HS-GAL), maximum graph operation (MGO), and output layer. Firstly, by using the graph attention network and graph pooling, we learned spectral and temporal representations from the input to construct the spectral input graph  $G_s$  ( $G_s \in R^{N_s \times d_s}$ ) and the temporal input graph  $G_t$  ( $G_t \in R^{N_t \times d_t}$ ). Here,  $N_s$  and  $N_t$  are the sets of nodes in the spectral and temporal graphs, and  $d$  is the feature dimension of each node. Next, we project  $G_s$  and  $G_t$  into another latent space with a common dimension  $d_{st}$  to construct a combined heterogeneous graph  $G_{st}$ .  $G_{st}$  has  $N_s + N_t$  nodes, where each node in  $G_s$  is connected by edges to every node in  $G_t$ .

HS-GAL consists of heterogeneous attention and stack node. Ours MGO comprises two parallel branches, each of

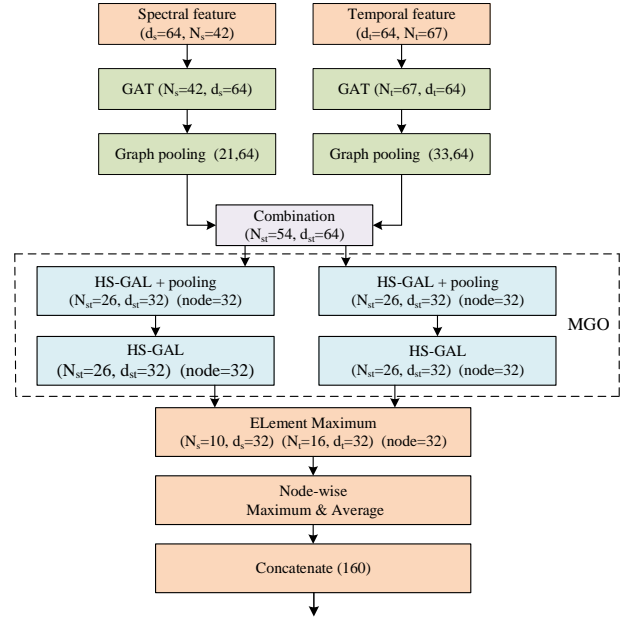


Figure 4: Integrated spectrotemporal graph attention network.

which includes two HS-GAL and one graph pooling layer. We then apply  $G_{st}$  to MGO and apply element-wise maximum operation to the outputs of MGO branches to generate another heterogeneous graph  $G_{ST}$ . HS-GAL in each branch of MGO shares a common stack node. Then, we apply maximum and average pooling to spectral and temporal nodes of  $G_{ST}$ , add the stack node, and concatenate these five different features as the output of the graph neural network. More details can be found in [25].

## 3.3. Attention aggregation layer

Attention statistic pooling [26] has shown positive effects on speaker recognition and fraud detection tasks [27, 28]. In order to extract target-related spectral and temporal representations and obtain discriminative information, we introduced a 2D attention aggregation layer between the multi-scale aggregation module and the graph neural network module. This attention aggregation layer uses 2D convolutional layers to generate a 2D attention weight matrix and calculates attention weights on both the temporal and frequency dimensions separately to reweight the features, instead of using traditional 1D attention. Our implementation is the same as in [28].

# 4. Experimental setup

## 4.1. Datasets

Task 2 of AP20-OLR is an open-set dialect identification task [29], which provides three Chinese dialect datasets, namely Hokkien, Sichuanese, and Shanghainese, for model training and evaluation. It is worth noting that the pre-trained model XLS-R does not include our target dialects. We use the same fine-tuning data and challenges as in the competition. Furthermore, the test set includes three interfering dialects, namely Mandarin, Malay, and Thai. Details of the original sources of these datasets are shown in Table 1.

Table 1: *Datasets used in our systems*

Data composition	
Training set	AP16-OL7, AP17-OL3, AP17-OLR-test, AP18-OLR-test, AP19-OLR-dev, AP20-OLR-dialect
validation set	AP19-OLR-dev-task3, AP19-OLR-zero
Test set	AP20-OLR-dialect-test

## 4.2. Implementation details

Data augmentation is an effective method to alleviate overfitting problems and improve model generalization. We used the RawBoost data augmentation tool [30] to add convolutional and additive noise to the existing fine-tuning data. In this experiment, audio data was cropped or padded to 4-second segments with a sampling rate of 16kHz. During fine-tuning, we used a batch size of 16 and the standard Adam optimizer [31] with a learning rate of  $10^{-6}$ . The output dimension of the fully connected layer during fine-tuning was set to 128, as shown on the right side of Figure 1. The specific parameter settings of the proposed back-end network are shown in Figure 2, Figure 3, and Figure 4. All models were trained for 20 epochs on a single NVIDIA V100S GPU.

## 4.3. Evaluation protocol

We followed the rules of the AP20-OLR challenge and used  $C_{avg}$  and  $EER$  as evaluation metrics for the dialect identification system [29]. Under open testing conditions, all interfering languages were treated as one unknown languages.  $C_{avg}$  is defined as the average cost performance between test languages, with a prior probability  $P_{target} = 0.5$  representing the likelihood of the target language.

## 5. Result and analysis

Table 2: *Performance comparison with and without back-end network during fine-tuning.*

System	$C_{avg}$	$EER(\%)$
Wav2vec2.0 + mean pooling	0.0645	6.63
Wav2vec2.0 + proposed back-end	<b>0.0298</b>	<b>3.78</b>

We used the wav2vec 2.0 pre-trained model as the feature extractor and applied two fine-tuning settings: one added a back-end multi-scale aggregation graph neural network, and the other added a simple back-end network, using only average pooling and fully connected classification layers, similar to [13] and [14]. As shown in Table 2, the addition of the multi-scale aggregation graph neural network in the back-end during fine-tuning resulted in a significant improvement in our system performance. Compared to the system with a simple back-end network, we achieved a relative improvement of 54% in  $C_{avg}$  and 43% in  $EER$ .

In Table 3, we summarize the performance of some top systems and our proposed system. To our knowledge, using the self-supervised wav2vec 2.0 front-end combined with back-end networks has resulted in the lowest  $C_{avg}$  in the OLR2020 dialect task report. However, we must acknowledge that the reported results in the literature were obtained using fixed training data,

Table 3: *Comparison with Top Systems. In this table, Winning team refers to the 1st place winner of the OLR Challenge 2020. DK-TDNN and Conformer systems are advanced systems published in Interspeech 2021.*

System	$C_{avg}$	$EER(\%)$
Winning team[32]	0.0738	11.97
DK-TDNN[33]	0.0670	6.52
Conformer[34]	0.0594	8.95
Wav2vec2.0 + proposed back-end	<b>0.0298</b>	<b>3.78</b>

while the results reported in this paper were obtained using a pre-trained model. Nevertheless, the improvement in  $C_{avg}$  relative to the best result is as high as 50%, indicating that using a pre-trained model combined with back-end networks can lead to significant performance improvements. We will continue to conduct ablation experiments for the various modules introduced in Section 3 and provide an overview of these results in Table 4.

Table 4: *Ablation Study of Multi-Scale Aggregation Graph Neural Network Architecture.*

System	$C_{avg}$	$EER(\%)$
Wav2vec2.0 + proposed back-end	<b>0.0298</b>	<b>3.78</b>
w/o multi-scale aggregation module	0.0363	4.18
w/o graph neural network module	0.0372	3.86
w/o attention aggregation layer	0.0446	5.37

According to the results in Table 4, it can be seen that each module in our proposed multi-scale aggregation graph neural network architecture has a significant impact on performance. When we remove the multi-scale aggregation module, the  $C_{avg}$  value and  $EER$  increase by 22% and 10%, respectively, indicating that this module plays a significant role in extracting features at different scales for dialect identification. Secondly, removing the graph neural network module results in a relative increase of 25% in  $C_{avg}$ . This may be because the graph neural network can better utilize the structural information of speech data. Finally, removing the attention aggregation layer results in a significant increase in both  $C_{avg}$  and  $EER$ , indicating that this module plays an important role in connecting the multi-scale aggregation module and the graph neural network module.

## 6. Conclusions

This paper explores the wav2vec 2.0 model for dialect identification, focusing on the impact of the back-end network during fine-tuning. The proposed method employs multi-scale aggregation and a graph neural network to design a more sophisticated back-end that implicitly exploit phoneme sequence information and significantly improves system performance. On the dialect task of the AP20-OLR dataset, we achieved state-of-the-art performance. Compared to the best-known results, the proposed system achieved relative reductions of 50% in  $C_{avg}$  and 42% in  $EER$ . Our research emphasizes the significance of integrating a more efficient back-end network to enhance the performance of dialect identification while utilizing the wav2vec 2.0 model.

## 7. References

- [1] H. Li, B. Ma, and K. A. Lee, "Spoken Language Recognition: From Fundamentals to Practice," *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1136–1159, May 2013.
- [2] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker Verification Using Adapted Gaussian Mixture Models," *Digital Signal Processing*, vol. 10, no. 1, pp. 19–41, Jan. 2000.
- [3] W. M. Campbell, J. P. Campbell, D. A. Reynolds, E. Singer, and P. A. Torres-Carrasquillo, "Support vector machines for speaker and language recognition," *Computer Speech & Language*, vol. 20, no. 2, pp. 210–229, Apr. 2006.
- [4] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-End Factor Analysis for Speaker Verification," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, May 2011.
- [5] D. Snyder, D. Garcia-Romero, A. McCree, G. Sell, D. Povey, and S. Khudanpur, "Spoken Language Recognition using X-vectors," in *The Speaker and Language Recognition Workshop (Odyssey 2018)*, Jun. 2018, pp. 105–111.
- [6] X. Miao, I. McLoughlin, and Y. Yan, "A New Time-Frequency Attention Mechanism for TDNN and CNN-LSTM-TDNN, with Application to Language Identification," in *Interspeech 2019*, Sep. 2019, pp. 4080–4084.
- [7] M. A. Zissman, "Comparison of four approaches to automatic language identification of telephone speech," *IEEE Transactions on speech and audio processing*, vol. 4, no. 1, p. 31, Jan. 1996.
- [8] R. Tong, B. Ma, H. Li, E. S. Chng, and K.-A. Lee, "Target-aware language models for spoken language recognition," in *Proc. Interspeech 2009*, Sep. 2009, pp. 200–203.
- [9] S. Ling, J. Salazar, Y. Liu, and K. Kirchhoff, "BERTphone: Phonetically-aware Encoder Representations for Utterance-level Speaker and Language Recognition," in *The Speaker and Language Recognition Workshop (Odyssey 2020)*, Nov. 2020, pp. 9–16.
- [10] M. Zhao, R. Li, S. Yan, Z. Li, H. Lu, S. Xia, Q. Hong, and L. Li, "Phone-Aware Multi-task Learning and Length Expanding for Short-Duration Language Recognition," in *2019 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, Nov. 2019, pp. 433–437.
- [11] H. Liu, L. P. Garcia Perera, A. Khong, S. Styles, and S. Khudanpur, "PHO-LID: A Unified Model Incorporating Acoustic-Phonetic and Phonotactic Information for Language Identification," in *Interspeech 2022*, Sep. 2022, pp. 2233–2237.
- [12] A. Baeviski, Y. Zhou, A. Mohamed, and M. Auli, "Wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 12 449–12 460.
- [13] Z. Fan, M. Li, S. Zhou, and B. Xu, "Exploring wav2vec 2.0 on Speaker Verification and Language Identification," in *Interspeech 2021*, Aug. 2021, pp. 1509–1513.
- [14] A. Tjandra, D. G. Choudhury, F. Zhang, K. Singh, A. Conneau, A. Baeviski, A. Sela, Y. Saraf, and M. Auli, "Improved Language Identification Through Cross-Lingual Self-Supervised Learning," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2022, pp. 6877–6881.
- [15] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [16] A. Babu, C. Wang, A. Tjandra, K. Lakhota, Q. Xu, N. Goyal, K. Singh, P. von Platen, Y. Saraf, J. Pino, A. Baeviski, A. Conneau, and M. Auli, "XLS-R: Self-supervised Cross-lingual Speech Representation Learning at Scale," in *Proc. Interspeech 2022*, Sep. 2022, pp. 2278–2282.
- [17] N. Vaessen and D. A. Van Leeuwen, "Fine-Tuning Wav2Vec2 for Speaker Recognition," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2022, pp. 7967–7971.
- [18] L.-W. Chen and A. Rudnicky, *Exploring Wav2vec 2.0 Fine-Tuning for Improved Speech Emotion Recognition*, Oct. 2021.
- [19] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," in *Proceedings of the 32nd International Conference on Machine Learning*, Jun. 2015, pp. 448–456.
- [20] B. Desplanques, J. Thienpondt, and K. Demuyne, "ECAPA-TDNN: Emphasized Channel Attention, Propagation and Aggregation in TDNN Based Speaker Verification," in *Interspeech 2020*, Oct. 2020, pp. 3830–3834.
- [21] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [22] S.-H. Gao, M.-M. Cheng, K. Zhao, X.-Y. Zhang, M.-H. Yang, and P. Torr, "Res2Net: A New Multi-Scale Backbone Architecture," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 2, pp. 652–662, 2021.
- [23] J. Lee and J. Nam, "Multi-level and multi-scale feature aggregation using pretrained convolutional neural networks for music auto-tagging," *IEEE signal processing letters*, vol. 24, no. 8, pp. 1208–1212, 2017.
- [24] Z. Gao, Y. Song, I. McLoughlin, P. Li, Y. Jiang, and L.-R. Dai, "Improving aggregation and loss function for better embedding learning in end-to-end speaker verification system," in *Interspeech 2019*, Sep. 2019, pp. 361–365.
- [25] J.-w. Jung, H.-S. Heo, H. Tak, H.-j. Shim, J. S. Chung, B.-J. Lee, H.-J. Yu, and N. Evans, "AASIST: Audio Anti-Spoofing Using Integrated Spectro-Temporal Graph Attention Networks," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2022, pp. 6367–6371.
- [26] K. Okabe, T. Koshinaka, and K. Shinoda, "Attentive Statistics Pooling for Deep Speaker Embedding," in *Interspeech 2018*, Sep. 2018, pp. 2252–2256.
- [27] Y. Jung, S. M. Kye, Y. Choi, M. Jung, and H. Kim, "Improving Multi-Scale Aggregation Using Feature Pyramid Module for Robust Speaker Verification of Variable-Duration Utterances," in *Interspeech 2020*, Oct. 2020, pp. 1501–1505.
- [28] H. Tak, M. Todisco, X. Wang, J. weon Jung, J. Yamagishi, and N. Evans, "Automatic Speaker Verification Spoofing and Deepfake Detection Using Wav2vec 2.0 and Data Augmentation," in *Proc. The Speaker and Language Recognition Workshop (Odyssey 2022)*, Jun. 2022, pp. 112–119.
- [29] Z. Li, M. Zhao, Q. Hong, L. Li, Z. Tang, D. Wang, L. Song, and C. Yang, "AP20-OLR Challenge: Three Tasks and Their Baselines," in *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, Dec. 2020, pp. 550–555.
- [30] H. Tak, M. Kamble, J. Patino, M. Todisco, and N. Evans, "Rawboost: A Raw Data Boosting and Augmentation Method Applied to Automatic Speaker Verification Anti-Spoofing," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, May 2022, pp. 6382–6386.
- [31] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR*, 2015.
- [32] J. Li, B. Wang, Y. Zhi, Z. Li, L. Li, Q. Hong, and D. Wang, "Oriental language recognition (olr) 2020: Summary and analysis," *arXiv*, 2021.
- [33] T. Kong, S. Yin, D. Zhang, W. Geng, X. Wang, D. Song, J. Huang, H. Shi, and X. Wang, "Dynamic Multi-Scale Convolution for Dialect Identification," in *Interspeech 2021*, Aug. 2021, pp. 3261–3265.
- [34] D. Wang, S. Ye, X. Hu, S. Li, and X. Xu, "An End-to-End Dialect Identification System with Transfer Learning from a Multilingual Automatic Speech Recognition Model," in *Interspeech 2021*, Aug. 2021, pp. 3266–3270.