SPECTRAL CLUSTERING-AWARE LEARNING OF EMBEDDINGS FOR SPEAKER DIARISATION

Evonne P.C. Lee, Guangzhi Sun, Chao Zhang, Philip C. Woodland

Cambridge University Engineering Dept., Trumpington St., Cambridge, CB2 1PZ U.K.

{epcl2,gs534,cz277,pcw}@eng.cam.ac.uk

ABSTRACT

In speaker diarisation, speaker embedding extraction models often suffer from the mismatch between their training loss functions and the speaker clustering method. In this paper, we propose the method of spectral clustering-aware learning of embeddings (SCALE) to address the mismatch. Specifically, besides an angular prototypical (AP) loss, SCALE uses a novel affinity matrix loss which directly minimises the error between the affinity matrix estimated from speaker embeddings and the reference. SCALE also includes *p*percentile thresholding and Gaussian blur as two important hyperparameters for spectral clustering in training. Experiments on the AMI dataset showed that speaker embeddings obtained with SCALE achieved over 50% relative speaker error rate reductions using oracle segmentation, and over 30% relative diarisation error rate reductions using automatic segmentation when compared to a strong baseline with the AP-loss-based speaker embeddings.

Index Terms— speaker diarisation, speaker embedding, spectral clustering, foundation model, wav2vec 2.0

1. INTRODUCTION

Speaker diarisation is the task of finding 'who spoke when' in an audio stream with multiple speakers. It has many applications including speech recognition and information retrieval *etc*. In a typical diarisation pipeline, audio data that contains speech is divided into fixed-length segments (e.g. 2 seconds), also known as windows. A speaker embedding is then extracted for each window using deep neural networks trained for speaker classification, known as d-vectors [1–5], and a speaker label is assigned via clustering algorithms. In particular, spectral clustering [6] has been widely adopted in many recent diarisation systems [7–9], which identifies groups of nodes based on the graph affinity matrix computed from speaker embeddings in a fully unsupervised process.

Using neural network-derived d-vectors in a diarisation system often suffers from the mismatch between training and clustering since d-vectors are not trained to discriminate the relative speaker differences across multiple utterances which is particularly important for unsupervised clustering. To mitigate such a mismatch, metric learning losses such as the angular prototypical (AP) loss [10–12] and the angular margin prototypical (AMP) loss [12] have been applied when fine-tuning a d-vector extraction model for speaker verification. Besides, work has been performed on supervised clustering methods [13, 14], or adaptive hyper-parameters for spectral clustering have been adopted [15]. Recently, representations from pre-trained speech foundation models [16–18], such as wav2vec 2.0, have achieved superior performance on speaker diarisation [19, 20].

However, a similar mismatch still exists since the pre-trained models are often fine-tuned by speaker classification [20, 21].

To resolve the mismatch between the speaker embedding extraction model and the downstream unsupervised speaker clustering, in this paper, we propose the spectral clustering-aware learning of embeddings (SCALE)¹ method for fine-tuning the pre-trained representations for speaker diarisation. Specifically, an AP loss is used to fine-tune the wav2vec 2.0 model to encourage the representations to discriminate the relative speaker differences across multiple utterances. To simulate spectral clustering during training, an extra affinity matrix (AM) loss is used, which minimises the mean squared errors between the reference and hypothesis affinity matrices and helps learn the structure of the speaker embedding space. Furthermore, the two key steps in spectral clustering, namely the Gaussian blurring and p-th percentile thresholding, are also accounted for during training in both the AP loss and the AM loss. Speaker diarisation experiments were performed on the widely used AMI meeting corpus [22], as well as a combination of AMI and VoxCeleb1+2 [23, 24] datasets. Consistent and statistically significant improvements were achieved using SCALE on both settings compared to both the speaker classification baseline and a strong AP loss baseline.

The rest of the paper is organised as follows. Section 2 gives an overview of related work. Section 3 explains about the speaker diarisation system. In Section 4, SCALE is presented. The experimental setup is given in Section 5 and the results are presented in Section 6.

2. RELATED WORK

2.1. Wav2vec 2.0 for speech tasks

Wav2vec 2.0 [17] learns generic speech representations via selfsupervised learning on large amounts of unlabelled speech data. It is trained to identify the correct quantised latent representation from a set of distractors, where the distractors are sampled uniformly from other masked time steps of the same utterance. Despite its success in various speech tasks, fine-tuning the model has only achieved moderate performance on speaker-related tasks, including speaker recognition [25] and verification [26], and speaker diarisation [19, 20], due to the mismatch between training and inference [27, 28].

2.2. Loss functions for speaker diarisation

For speaker-related tasks, metric-learning-based loss functions have demonstrated competitive results [10]. Due to the constraints from the unsupervised clustering present later in the pipeline, it is preferable to extract speaker embeddings with small intra-speaker distances and large inter-speaker distances. Work in [29] proposed the

Guangzhi Sun is supported by Cambridge Trust

¹Code available at https://github.com/epcl2/scale

triplet loss which required a triplet to be selected meticulously as the performance relied on the "difficulty" of the negative exemplars. The prototypical loss was proposed in [30], where a query embedding was pushed away from the centroid of all negative samples based on the squared Euclidean distance in a mini-batch. Instead of using only one utterance from each speaker as the query, in the generalised end-to-end loss [31], every utterance in the mini-batch served as a query. Similar to the prototypical loss, the angular prototypical (AP) loss [10] used only one utterance from each class as the query, with a cosine similarity-based metric. Work in [32] used a contrastive self-supervised learning approach for text-independent speaker verification, where they adopted the AP loss.

3. SPEAKER DIARISATION PIPELINE

The full speaker diarisation pipeline in this paper comprises several stages, including voice activity detection (VAD), change point detection (CPD), speaker embedding extraction and the spectral clustering stage. Neural VAD and CPD were used, which followed [8]. The VAD detects audio that contains speech and CPD splits speech segments into speaker-homogeneous segments. Each speaker-homogeneous segment is split into multiple windows of the same length, and speaker embeddings are extracted for each window using a fine-tuned wav2vec 2.0 encoder. Finally, spectral clustering is performed at the window level and all windows in the same speaker-homogeneous segment are assigned to the same speaker.

3.1. Fine-tuning with angular prototypical loss

The wav2vec 2.0 encoder is fine-tuned with an extra output layer based on the AP loss. For each mini-batch, N utterances, $\mathbf{u}_1^1, \ldots, \mathbf{u}_N^a$ from N distinct speakers are randomly selected and act as the anchor utterances. Another N utterances, $\mathbf{u}_1^1, \ldots, \mathbf{u}_N^a$, are selected from the same N speakers associated with the anchor utterances, which act as the positive utterances. For each pair of utterance $i \ (1 \le i \le N), \mathbf{u}_i^a$ and \mathbf{u}_i^p have the same speaker identity but $\mathbf{u}_i^a \ne \mathbf{u}_i^p$. All anchor utterances from other speakers $j \ (j \ne i)$ in the same mini-batch serve as the negative samples for speaker *i*. The anchor and positive embeddings \mathbf{e}_i^a and \mathbf{e}_i^p are derived from the penultimate layer of the model (*i.e.* the final wav2vec 2.0 encoder layer). The AP loss, \mathcal{L}_{AP} , used to optimise the model is then

$$\mathcal{L}_{AP} = -\frac{1}{N} \sum_{i=1}^{N} \log \frac{\exp(\mathbf{S}_{i,i})}{\sum_{j=1}^{N} \exp(\mathbf{S}_{i,j})}$$
(1)

where $S_{i,j}$ is an element in the similarity matrix S, defined as the scaled cosine similarity between each pair of embeddings. That is,

$$\mathbf{S}_{i,j} = \begin{cases} w \cdot \sin(\mathbf{e}_i^{\mathrm{a}}, \mathbf{e}_j^{\mathrm{p}}) + b & \text{if } i = j \\ w \cdot \sin(\mathbf{e}_i^{\mathrm{a}}, \mathbf{e}_j^{\mathrm{a}}) + b & \text{otherwise} \end{cases}$$
(2)

where $\sin(\mathbf{x}, \mathbf{y}) = (\cos(\mathbf{x}, \mathbf{y}) + 1)/2$, and w and b are trainable scalar values. In $\mathbf{S} = (\mathbf{S}_{i,j}) \in \mathbb{R}^{N \times N}$, the similarity score of diagonal terms is computed between the anchor embedding and the positive embedding; the similarity score of the non-diagonal terms is computed between anchor embeddings to simulate the situation during spectral clustering. With the AP loss, the model explicitly minimises the distance among embeddings from the same speaker across different utterances (intra-speaker distances) and maximises the distance among embeddings from different speaker distances). This differs from the pre-training of wav2vec 2.0 where the distractors were sampled from the same utterance.

3.2. Spectral clustering

In spectral clustering, the affinity matrix **A** is first constructed, where \mathbf{A}_{ij} is the cosine similarity between the embeddings of the *i*-th and *j*-th window for $i \neq j$. A modified version of spectral clustering based on [3] was used, where a sequence of refinement steps was applied to de-noise the affinity matrix. Two crucial steps, Gaussian blur and row-wise thresholding are highlighted here.

(i) Gaussian blur: A Gaussian kernel with standard deviation σ is used to smooth the data. Since window-level clustering is used, neighbouring embeddings are often derived from the same utterance, and hence should have similar values in the affinity matrix. Gaussian blur preserves this property among neighbouring windows.

(ii) Row-wise thresholding: Any element whose value ranked less than a particular row's *p*-th percentile is set to be zero, which "zeroes out" the affinities between embeddings from two distinct speakers. Ideally, the similarity scores of the embeddings belonging to different speakers should be below this threshold so that they can be "thresholded-out", whereas those of the embeddings belonging to the same speakers should be above this threshold and reserved.

4. SPEAKER CLUSTERING WITH SCALE

In Section 3, the model was fine-tuned without accounting for the clustering stage. To further reduce the mismatch between training and clustering, three steps are introduced to training, i.e. the AM loss, absolute thresholding and relative thresholding with Gaussian blurring. SCALE modifies the loss function such that the model is aware of the clustering stage and the associated refinement steps.

4.1. Affinity matrix loss

The AM loss encourages the hypothesis affinity matrix \mathbf{A} , constructed using real speaker embeddings to be close to the reference affinity matrix with 1 for pairs from the same speaker and 0 for pairs from different speakers. By sampling N pairs of utterances for N distinct speakers in each mini-batch (as in Section 3.1), the reference affinity matrix is an identity matrix \mathbf{I} . \mathcal{L}_{AM} , the AM loss calculating the mean squared errors between \mathbf{I} and \mathbf{A} , is

$$\mathcal{L}_{\rm AM} = \frac{1}{N^2} \sum_{i=1}^{N} \sum_{j=1}^{N} (\mathbf{I}_{i,j} - \mathbf{A}_{i,j})^2$$
(3)

For \mathcal{L}_{AM} , **A** is equivalent to **S** in Section 3.1 with w = 1 and b = 0. The final form of loss is a weighted combination of \mathcal{L}_{AP} and \mathcal{L}_{AM} as below.

$$\mathcal{L} = (1 - \alpha) \mathcal{L}_{AP} + \alpha \mathcal{L}_{AM} \tag{4}$$

where α is between 0 and 1. For more than two speakers, although the ideal affinity matrix where embeddings from different speakers are opposite to each other cannot be achieved, the main aim for the AM loss is to increase the cosine distance for different speakers. Threhsolding in the next subsection resolves this issue.

4.2. Absolute thresholding

For the row-wise thresholding step in spectral clustering, any value less than the threshold was set close to zero. Hence, it is important to encourage pairs of embeddings from different speakers to have a similarity score below the threshold during training. Therefore, a threshold, t, between 0 and 1 was used during training. The threshold can be applied for both the AM and AP losses. For the AM loss, a mask matrix $\mathbf{M} = (\mathbf{M}_{i,j}) \in \mathbb{R}^{N \times N}$ is created to only calculate

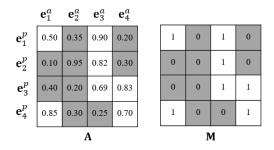


Fig. 1. An example of **A** with its corresponding mask **M**. The numbers in **A** are the cosine similarity scores between the speaker embeddings of the corresponding columns and rows. The AM loss is only calculated for cells that are not greyed out.

the positive pairs whose similarity scores are lower than t and the negative pairs whose similarity higher than t. That is,

$$\mathbf{M}_{i,j} = \begin{cases} 1 & \text{if } i = j, \mathbf{A}_{i,j} \leq t \\ 1 & \text{if } i \neq j, \mathbf{A}_{i,j} \geq t \\ 0 & \text{otherwise} \end{cases}$$
(5)

Thus the AM loss becomes

$$\mathcal{L}_{AM} = \frac{1}{\sum_{i,j} \mathbf{M}_{i,j}} \sum_{i=1}^{N} \sum_{j=1}^{N} \mathbf{M}_{i,j} (\mathbf{I}_{i,j} - \mathbf{A}_{i,j})^{2}.$$
 (6)

The thresholding concept with t = 0.8 is illustrated in Fig. 1. This enables the loss to focus only on the difficult positive and negative pairs as the "easy negative pairs" would have been thresholded to 0. Similarly, thresholding can be imposed on the AP loss, where a mask can be created such that the cross-entropy is calculated for all positive pairs and the difficult negative pairs, as the positive pairs are always needed.

4.3. Relative thresholding with Gaussian blurring

As relative thresholding was used in spectral clustering, to closely simulate the row-wise thresholding during clustering, the Gaussian blur and relative thresholding are integrated into speaker embedding extraction model training. This improves the estimation of the actual threshold by using Gaussian-blurred matrices. To achieve this, the affinity matrix is first blurred to determine the relative threshold. The threshold which is in general different for different rows is then used on the unblurred affinity matrix to avoid gradient annihilation. The detailed steps are given below and shown in Figure 2.

- 1. Compute the affinity matrix.
- Replace all diagonal elements with 1 since, during clustering, diagonal elements are the maximum of the row.
- 3. Perform Gaussian blurring on the affinity matrix.
- 4. Determine the relative threshold of each row by multiplying *t* with the value of the diagonal element.
- 5. Use the relative threshold to create a single mask matrix for both of the AP and AM losses.

5. EXPERIMENTAL SETUP

5.1. Data

The dataset used to fine-tune and evaluate the final speaker diarisation performance was the AMI meeting corpus [22], which consists

	\mathbf{e}_1^a	\mathbf{e}_2^a	\mathbf{e}_3^a	\mathbf{e}_4^a										
\mathbf{e}_1^p	0.50	0.35	0.90	0.20		0.74	0.56	0.77	0.36	0.59	1	0	1	0
\mathbf{e}_2^p	0.10	0.95	0.82	0.30		0.36	0.78	0.79	0.45	0.62	0	0	1	0
\mathbf{e}_3^p	0.40	0.20	0.69	0.83		0.39	0.37	0.82	0.81	0.66	0	0	0	1
\mathbf{e}_4^p	0.85	0.30	0.25	0.70		0.65	0.34	0.46	0.85	0.68	1	0	0	0
	Α				blurred A			Rel. Thold.			М			

Fig. 2. A and **A** with Gaussian blur applied. The row-wise relative thresholds are obtained as described in step 4, which are then applied to **A**, giving the mask **M**. Rel. Thold. stands for Relative Threshold.

of meeting recordings with 4-5 speakers per meeting. The training set contains 135 meetings with 155 speakers. The development (Dev) and evaluation (Eval) sets followed [8]. Moreover, the joint VoxCeleb1 and VoxCeleb2 data were used as an intermediate finetuning stage followed by the final fine-tuning on the AMI training set in a two-stage fine-tuning setup. The input features of all models were beamformed raw waveform with BeamformIt [33].

5.2. System specifications

During training, a 2-second window was sampled from each speech segment as input. During inference, speaker embeddings were extracted from the penultimate layer of the embedding extraction model. An average pooling layer was added on top of the wav2vec 2.0 model to produce an embedding for each window. This was followed by two additional fully connected layers, the first layer projected the embedding to the desired dimension (128D) and the second one was used to classify each input into a speaker label. The wav2vec 2.0 encoder was frozen for the first 10% of fine-tuning steps and a triangular learning rate scheduler was adopted. For two-stage fine-tuning, SCALE was only applied to the stage on AMI as there was a mismatch between the nature of the VoxCeleb and AMI data. For the baseline model, the angular Softmax (A-Softmax) loss [34, 35] was used for speaker classification, with m = 1.

During inference, each speaker-homogeneous segment was split into 2-second windows with a 1-second overlap. Speaker embeddings were obtained for each window and spectral clustering described in Sections 3.2 and 4 were applied. All windows in a segment were assigned the same speaker label, following [8]. The hyperparameters for spectral clustering (the *p*-th percentile and standard deviation for Gaussian blur) were tuned on the Dev set by grid search and applied to the Eval set directly. Three baselines were used, including a non-wav2vec *TDNN baseline* following [8], a *classification baseline* using A-Softmax loss for fine-tuning wav2vec 2.0 and the strongest *AP loss baseline* that used the AP loss for fine-tuning.

5.3. Evaluation

Both reference and automatic segments were used for evaluation. Automatic segments were found using the same VAD and CPD in [8]. For the reference segmentation, the SER was scored with a 0.25-second collar on both sides of the segment without overlap. For the automatic segmentation, DER which is the sum of the SER, the missed speech (MS) and false alarm (FA), was reported. To avoid scoring against a large amount of long silence in the original reference, diarisation results with automatic segmentation were scored against the modified reference following [8]. In addition, a meetingby-meeting sign test was performed to show the statistical significance of any improvements where appropriate.

System	SER (%)		
System	Dev	Eval	
	Dev	Eval	
TDNN [8]	14.3	15.4	
Classification with A-Softmax	9.6	19.3	
AP loss	9.9	17.1	
AP loss + AM loss	8.8	16.4	
AP loss + AM loss (Abs. Thold.)	7.5	15.8	
AP loss (Abs. Thold.) + AM loss (Abs. Thold.)	6.6	13.2	
AP loss (Rel. Thold.) + AM loss (Rel. Thold.)	6.9	12.0	

Table 1. SER on AMI Dev and Eval sets using systems with SCALE fine-tuned on AMI only. TDNN is a non-wav2vec baseline while all the other systems were fine-tuned wav2vec 2.0. Abs. Thold. refers to the absolute thresholding (Sec. 4.2) and Rel. Thold. refers to the relative thresholding (Sec. 4.3) for all tables.

6. RESULTS AND DISCUSSION

6.1. With one-stage fine-tuning

In this part of the experiment, wav2vec 2.0 was directly fine-tuned on AMI with SCALE and the results are shown in Table 1. For SCALE, α was set to 0.5 for the AP loss, and t was set to 0.8 and 0.95 for the absolute and relative thresholding respectively. For relative thresholding with Gaussian blur, σ was set to 1. The two wav2vec 2.0 baselines performed better on the Dev set while worse on the Eval set compared to the TDNN baseline²

With the AM loss term added, the Dev and Eval sets SERs reduced, demonstrating the effectiveness of the AM loss term. Next, when absolute thresholding was applied to the AM loss, reductions in SERs on both Dev and Eval sets were found and the Eval set SER was similar to the TDNN baseline while the Dev set DER is clearly lower. Finally, by applying the SCALE method with absolute thresholding on both AM and AP losses, the system achieved 27% relative SER reductions on both the Dev and Eval sets compared to using only the AP loss.

Applying thresholding on both losses was more effective than only applying it on the AP loss, as the training process put more emphasis on the harder pairs. When using a relative threshold, the Dev SER was slightly higher than when using an absolute threshold, but the Eval SER was lower. A meeting-level sign test showed that the improvement achieved by SCALE compared to the classification baseline was statistically significant, and there was no statistical difference between absolute and relative thresholding results.

6.2. With two-stage fine-tuning

In the two-stage fine-tuning, SCALE was applied to the final finetuning on AMI, where $\alpha = 0.5$ and t = 0.9 were used for both absolute and relative thresholding. For relative thresholding, $\sigma = 0.5$ was used with Gaussian blur. The results are shown in Table 2. As before, fine-tuning with the AP loss resulted in better diarisation performance than with the classification loss. When applying SCALE with thresholding on both the AM and AP losses, the model performed much better than the baselines. Meanwhile, the gap between the Dev and Eval SER became smaller, and it was also discovered that embeddings with SCALE had a much smaller SER change during hyperparameter tuning for spectral clustering. Therefore, SCALE achieved better robustness to spectral clustering hy-

System	SER (%)		
	Dev	Eval	
Classification with A-Softmax	9.0	14.7	
AP loss	7.0	9.6	
AP loss + AM loss	6.1	9.7	
AP loss + AM loss (Abs. Thold.)	6.4	11.2	
AP loss (Abs. Thold.) + AM loss (Abs. Thold.)	3.3	4.7	
AP loss (Rel. Thold.) + AM loss (Rel. Thold.)	3.4	4.9	

Table 2. SER on AMI Dev and Eval sets using systems first finetuned on VoxCeleb1+2, and then on AMI with SCALE. All systems used fine-tuned wav2vec 2.0.

System	DER (%)		
	Dev	Eval	
TDNN [8]	12.6	15.6	
AP loss	14.1	17.1	
AP loss (Abs. Thold.) + AM loss (Abs. Thold.)	9.9	10.1	
AP loss (Rel. Thold.) + AM loss (Rel. Thold.)	9.9	11.2	

Table 3. DERs on AMI Dev and Eval sets using automatic segmentation. MS + FA rates was 5.2% and 4.9% on the Dev and Eval set.

perparameters. Overall, the best system was achieved using absolute thresholding on both losses, resulting in a relative SER reduction of 53% and 51% on Dev and Eval sets respectively when compared to the AP loss baseline. As before, the improvements were found to be statistically significant, and there was no statistical difference between the two thresholding methods.

6.3. With automatic segmentation

To investigate the effectiveness of using SCALE to fine-tune a wav2vec 2.0 in the full speaker diarisation pipeline, experiments were performed with automatic segmentation obtained after VAD and CPD, with results shown in Table 3. The MS and FA rates obtained after VAD and CPD were 1.2% and 4.0% respectively for the Dev set, and 1.3% and 3.6% respectively for the Eval set. The system trained with AP loss only, and SCALE systems with absolute and relative thresholding on both losses were investigated, which were the same as those in Table 2. Compared to the AP loss, using SCALE with absolute thresholding achieved the best performance among systems, with a relative reduction of 30% an the Dev set and 40% on the Eval set.

7. CONCLUSIONS

This paper proposed SCALE, spectral clustering-aware learning of embedding framework to fine-tune the pre-trained speaker representations for speaker diarisation. SCALE used a combination of the angular prototypical (AP) and the affinity matrix (AM) losses to learn the structure of the embedding space, and to reduce the mismatch between the speaker embeddings and spectral clustering. SCALE also includes the *p*-percentile thresholding and Gaussian blur of spectral clustering into training to further reduce the mismatch. Experiments on the AMI data with SCALE achieved a relative SER reduction of 53% and 51% on the Dev and Eval sets respectively, compared to the AP loss baseline using oracle segmentation. SCALE also achieved 30% and 40% relative DER reduction on the Dev and Eval sets respectively with automatic segmentation.

²A clear comparison with the literature is hard as various setups are used, e.g. VBx[36] gives good error rates but cannot be directly compared.

References

- E. Variani, X. Lei, E. McDermott, I. L. Moreno, and J. Gonzalez-Dominguez, "Deep neural networks for small footprint text-dependent speaker verification," in *Proc. ICASSP*, Florence, 2014.
- [2] G. Heigold, I. Moreno, S. Bengio, and N. Shazeer, "End-toend text-dependent speaker verification," in *Proc. ICASSP*, Shanghai, 2016.
- [3] Q. Wang, C. Downey, L. Wan, P. A. Mansfield, and I. L. Moreno, "Speaker diarization with LSTM," in *Proc. ICASSP*, Calgary, 2018.
- [4] D. Garcia-Romero, D. Snyder, G. Sell, D. Povey, and A. Mc-Cree, "Speaker diarization using deep neural network embeddings," in *Proc. ICASSP*, New Orleans, 2017.
- [5] D. Snyder, D. Garcia-Romero, D. Povey, and S. Khudanpur, "Deep neural network embeddings for text-independent speaker verification.," in *Proc. Interspeech*, 2017.
- [6] U. von Luxburg, "A tutorial on spectral clustering," *Statistics and Computing*, vol. 17, pp. 395–416, 2007.
- [7] H. Ning, M. Liu, H. Tang, and T. Huang, "A spectral clustering approach to speaker diarization," in *Proc. ICSLP*, 2006.
- [8] G. Sun, C. Zhang, and P. C. Woodland, "Combination of deep speaker embeddings for diarisation," *Neural Networks*, vol. 141, pp. 372–384, Sep. 2021.
- [9] T. J. Park, K. J. Han, M. Kumar, and S. Narayanan, "Autotuning spectral clustering for speaker diarization using normalized maximum eigengap," *IEEE Signal Processing Letters*, vol. 27, pp. 381–385, 2019.
- [10] J. S. Chung, J. Huh, S. Mun, M. Lee, H.-S. Heo, S. Choe, C. Ham, S. Jung, B.-J. Lee, and I. Han, "In defence of metric learning for speaker recognition," in *Proc. Interspeech*, Shanghai, 2020.
- [11] A. Brown, J. Huh, A. Nagrani, J. S. Chung, and A. Zisserman, "Playing a part: Speaker verification at the movies," in *Proc. ICASSP*, Toronto, 2021.
- [12] D. V. Thanh, T. P. Viet, and T. N. T. Thu, "Deep speaker verification model for low-resource languages and Vietnamese dataset," in *Proc. PACLIC*, Shanghai, 2021.
- [13] Y. Fujita, N. Kanda, S. Horiguchi, Y. Xue, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with self-attention," in *Proc. ASRU*, Singapore, 2019.
- [14] Q. Li, F. L. Kreyssig, C. Zhang, and P. C. Woodland, "Discriminative neural clustering for speaker diarisation," in *Proc. SLT*, Shenzhen, 2021.
- [15] W. Xia, H. Lu, Q. Wang, A. Tripathi, Y. Huang, I. L. Moreno, and H. Sak, "Turn-to-diarize: Online speaker diarization constrained by transformer transducer speaker turn detection," in *Proc. ICASSP*, Singapore, 2022.
- [16] S. Schneider, A. Baevski, R. Collobert, and M. Auli, "Wav2vec: Unsupervised pre-training for speech recognition," in *Proc. Interspeech*, Graz, 2019.
- [17] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "Wav2vec 2.0: A framework for self-supervised learning of speech representations," in *Proc. NeurIPS*, Vancouver, 2020.
- [18] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 3451–3460, 2021.
- [19] S.-W. Yang, P.-H. Chi, Y.-S. Chuang, C.-I. J. Lai, K. Lakhotia, Y. Y. Lin, A. T. Liu, J. Shi, X. Chang, G.-T. Lin, T.-H.

Huang, W.-C. Tseng, K.-t. Lee, D.-R. Liu, Z. Huang, S. Dong, S.-W. Li, S. Watanabe, A. Mohamed, and H.-y. Lee, "SUPERB: Speech processing universal performance benchmark," in *Proc. Interspeech*, Brno, 2021.

- [20] X. Zheng, C. Zhang, and P. C. Woodland, "Tandem multitask training of speaker diarisation and speech recognition for meeting transcription," in *Proc. Interspeech*, 2022.
- [21] L. Zhang, H. Zhao, Q. Meng, Y. Chen, M. Liu, and L. Xie, "Beijing ZKJ-NPU speaker verification system for voxceleb speaker recognition challenge 2021," *arXiv preprint arXiv:2109.03568*, 2021.
- [22] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal, G. Lathoud, M. Lincoln, A. Lisowska, I. McCowan, W. Post, D. Reidsma, and P. Wellner, "The AMI meeting corpus: A pre-announcement," in *Proc. MLMI*, Bethesda, 2006.
- [23] A. Nagrani, J. S. Chung, and A. Zisserman, "VoxCeleb: A large-scale speaker identification dataset," in *Proc. Inter*speech, Stockholm, 2017.
- [24] J. S. Chung, A. Nagrani, and A. Zisserman, "VoxCeleb2: Deep speaker recognition," in *Proc. Interspeech*, 2018.
- [25] N. Vaessen and D. A. Van Leeuwen, "Fine-tuning wav2vec2 for speaker recognition," in *Proc. ICASSP*, Singapore, 2022.
- [26] Z. Fan, M. Li, S. Zhou, and B. Xu, "Exploring wav2vec 2.0 on speaker verification and language identification," in *Proc. Interspeech*, Brno, 2021.
- [27] Z. Chen, S. Chen, Y. Wu, Y. Qian, C. Wang, S. Liu, Y. Qian, and M. Zeng, "Large-scale self-supervised speech representation learning for automatic speaker verification," in *Proc. ICASSP*, Singapore, 2022.
- [28] S. Chen, Y. Wu, C. Wang, Z. Chen, Z. Chen, S. Liu, J. Wu, Y. Qian, F. Wei, J. Li, *et al.*, "UniSpeech-SAT: Universal speech representation learning with speaker aware pre-training," in *Proc. ICASSP*, Singapore, 2022.
- [29] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in *Proc. CVPR*, Boston, 2015.
- [30] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Proc. NeurIPS*, Long Beach, 2017.
- [31] L. Wan, Q. Wang, A. Papir, and I. L. Moreno, "Generalized end-to-end loss for speaker verification," in *Proc. ICASSP*, Calgary, 2018.
- [32] H. Zhang, Y. Zou, and H. Wang, "Contrastive self-supervised learning for text-independent speaker verification," in *Proc. ICASSP*, Toronto, 2021.
- [33] X. Anguera, C. Wooters, and J. Hernando, "Acoustic beamforming for speaker diarization of meetings," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 7, pp. 2011–2022, 2007.
- [34] Z. Huang, S. Wang, and K. Yu, "Angular softmax for shortduration text-independent speaker verification," in *Proc. Interspeech*, Hyderabad, 2018.
- [35] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "SphereFace: Deep hypersphere embedding for face recognition," in *Proc. CVPR*, Honolulu, 2017.
- [36] F. Landini, J. Profant, M. Diez, and L. Burget, "Bayesian HMM clustering of x-vector sequences (VBx) in speaker diarization: Theory, implementation and analysis on standard tasks," *Computer Speech & Language*, vol. 71, p. 101 254, 2022.