



Multi-scale Speaker Diarization with Dynamic Scale Weighting

Tae Jin Park, Nithin Rao Koluguri, Jagadeesh Balam and Boris Ginsburg

NVIDIA

{taejinp, nkoluguri, jbalam, bginsburg}@nvidia.com

Abstract

Speaker diarization systems are challenged by a trade-off between the temporal resolution and the fidelity of the speaker representation. By obtaining a superior temporal resolution with an enhanced accuracy, a multi-scale approach is a way to cope with such a trade-off. In this paper, we propose a more advanced multi-scale diarization system based on a multi-scale diarization decoder. There are two main contributions in this study that significantly improve the diarization performance. First, we use multi-scale clustering as an initialization to estimate the number of speakers and obtain the average speaker representation vector for each speaker and each scale. Next, we propose the use of 1-D convolutional neural networks that dynamically determine the importance of each scale at each time step. To handle a variable number of speakers and overlapping speech, the proposed system can estimate the number of existing speakers. Our proposed system achieves a state-of-art performance on the CALLHOME and AMI MixHeadset datasets, with 3.92% and 1.05% diarization error rates, respectively.

Index Terms: speaker diarization, multi-scale

1. Introduction

Speaker diarization is a task of partitioning an input audio stream into speaker-homogeneous segments allowing audio segments to be associated with speaker labels. Speaker diarization requires decisions on relatively short segments ranging from a few tenths of a second to several seconds. In speaker embedding extraction, to obtain high-quality speaker representation vectors, the temporal resolution should be sacrificed by taking a long speech segment. Thus, speaker diarization systems always face a trade-off between two quantities, i.e., the temporal resolution and quality of the representation. In the early versions of speaker diarization, the speaker homogeneous variable-length Mel-frequency cepstral coefficients (MFCCs) segments were generated through the detection of speaker change points [1].

Owing to the increased popularity of i-vectors [2] and x-vectors [3] in the field of speaker diarization, a uniform segmentation approach has been widely used, in which a fixed speaker segment length is applied to extract the speaker representation vector per segment. Because fixed-length segments lead to stable speaker representations by controlling the length factor in speaker embedding extraction, speaker diarization systems based on a uniform segmentation method [2, 4] have shown a competitive performance. However, the uniform segmentation approach has innate limitations in terms of the temporal resolution because reducing the segment length leads to a decrease in accuracy [5]. In a uniform segmentation setting, the shortest temporal resolution is limited to the hop-length during the segmentation process. Moreover, without a post processing approach, because each segment is assigned to only one speaker, the clustering approaches lack an overlap-aware diarization capability.

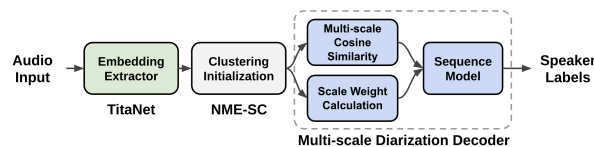


Figure 1: Data-flow of the proposed multi-scale speaker diarization system.

Recent speaker diarization systems such as an end-to-end approach [6] or target-speaker voice activity detection (TS-VAD) [7] employ feature frame-level speaker labels. These systems are based on sequence models such as long short-term memory (LSTM) [7, 8] or a Transformer-style encoder-decoder [6]. Frame-level sequence model based systems enable a far superior temporal resolution of approximately 0.01 s. Another significant benefit of these methods is overlap-aware diarization where the neural network model outputs more than one speaker label output. The multi-dimensional sigmoid output enables an overlap detection without employing an additional resegmentation module on top of the diarization system.

However, sequence model based methods also have a few drawbacks. First, because these sequence models are trained with a fixed number of speakers, sequence model based diarization systems often lack the capability to handle a flexible number of speakers [7] or estimate their numbers [6]. Second, according to the recent results of a diarization challenge [9, 10], the performances of sequence model based end-to-end approaches [6, 8, 11] still lag behind that of a state-of-art modular diarization [7, 12]. The hidden Markov model (HMM) based clustering approach proposed in [12] does not have the aforementioned downsides, having the capabilities of speaker counting and overlap detection and achieving a state-of-art performance. Thus, we compared our proposed system with the system proposed in [12] on the same datasets.

Our proposed multi-scale diarization decoder (MSDD) tackles the problems inherent to previous studies. The proposed MSDD was designed to support the following features: overlap-aware diarization, an improved temporal resolution, and a flexible number of speakers. As shown in Fig.1, the proposed system comprises three components overall, i.e., a pretrained speaker embedding extractor (TitaNet, [13]), a multi-scale speaker clustering module, and the MSDD model.

Although we apply a multi-scale clustering idea from a previous study [5], we do not estimate the static session level multi-scale weights. Instead, we propose the MSDD approach, which takes advantage of the initialized clusters by comparing the extracted speaker embeddings with the cluster-average speaker representation vectors. The weight of each scale at each time step is determined through a scale weighting mechanism where the scale weights are calculated from a 1-D convolutional neural network (CNN) applied to the input speaker embeddings and the cluster-average embedding based on the clustering results. Finally, based on the weighted cosine similarity vectors, an LSTM-based sequence model estimates the labels of each

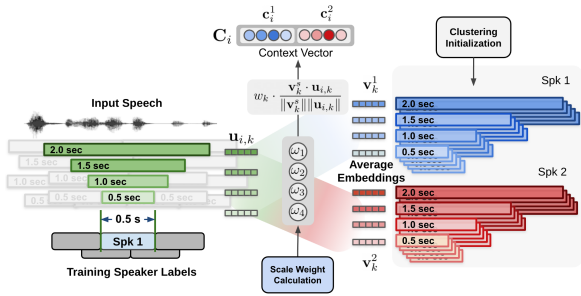


Figure 2: Cosine similarity values from each speaker and each scale are weighted by the scale weights to form a weighted cosine similarity vector.

speaker. As described in [7], using the clustering result as an initialization allows the proposed system to be free from permutations issue [6] and turns the diarization problem into binary classification problem at each time step. Unlike the systems described in [6, 7] we only train and test two-speaker models to handle a flexible number of speakers. Therefore, the speaker labels are predicted by taking the average of the sigmoid outputs from the multiple pairs. The 2-D dimensional output enables an overlap-aware diarization for the given input stream.

In the experimental result section, we demonstrate that the proposed multi-scale approach is more accurate than our previous single scale (SS), clustering based diarization system. Furthermore, we show that our proposed system achieves new state-of-art results on the CALLHOME and AMI *MixHeadset* datasets.

2. Multi-scale Clustering

2.1. Multi-scale segmentation

The largest difference between a conventional single scale clustering method and multi-scale method is in the way in which the speech stream is segmented. The proposed system uses the same multi-scale embedding extraction and clustering mechanism as in [5], except that the scale weights are computed differently in this work. We employ a uniform segmentation scheme that originally appeared in [2, 4], and such segmentation is applied for multiple scales. Fig. 2 shows an example of a multi-scale segmentation. As shown in Fig. 2, we refer to the finest scale, i.e., 0.5 s, as the base scale because the speaker label estimation is applied at the finest scale. We denote the number of scales as K .

After segmentation is applied for all scales, grouping among the segments at each scale is then applied. The grouping process is conducted by assigning the segments from each scale for the corresponding base scale segment. The segments from the lower temporal resolution scales (in this example, 2.0, 1.5, and 1.0 s in length) are selected and grouped by measuring the distance between the centers of the segments and choosing the closest ones. By grouping the segments as in Fig. 2, each group will generate one weighted cosine similarity value when a weighted affinity matrix is calculated. A detailed description of multi-scale clustering can be found in [5]. We use the speaker embedding extractor model proposed in [13], which generates 192-dimensional embedding.

2.2. Clustering for initialization

Because we previously proposed a multi-scale system in [5], we employ the auto-tuning spectral clustering approach proposed in [14], which is referred to as normalized maximum eigengap spectral clustering (NME-SC). From the clustering result, we

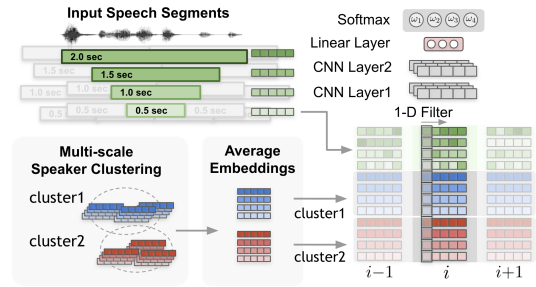


Figure 3: Scale weights calculated from a 1-D CNN. The 1-D filter captures the context from the input embeddings and cluster-average embeddings.

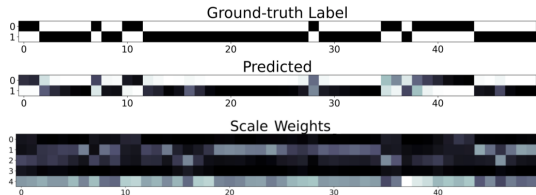


Figure 4: Example plot of target labels, prediction, and scale weights ($K=5$). Note that scale weights vary at each time step.

obtain the estimated cluster label per each base scale segment and the estimated number of speakers S .

Based on the clustering results, we take the average of all the speaker embeddings obtained from the initial clustering result as in Eq. (1). We obtain cluster-average embedding vectors \mathbf{v}_k^s from the initial clustering result as follows:

$$\mathbf{v}_k^s = \frac{1}{N_k} \sum_{i=1}^{N_k} \mathbf{v}_{i,k}^s, \quad (1)$$

where k is the scale index, N_k is the number of k -th scale embeddings during the session and s is the cluster index (speaker index) from the clustering result.

3. Multi-scale Diarization Decoder

3.1. 1-D CNN for Dynamic Scale Weights

To dynamically adjust the scale weights during the inference phase, the proposed diarization system should look into the input speaker and reference embeddings simultaneously. Fig. 3 shows how scale weights are calculated in the case of $K=4$ scales for $S=2$ speakers. We stack the embeddings from the input signal and the average embeddings from the clustering results of two speakers ($s=1,2$) in the following manner:

$$\mathbf{D}_i = [\mathbf{u}_{i,0}; \dots; \mathbf{u}_{i,K-1}; \mathbf{v}_0^1; \dots; \mathbf{v}_{K-1}^1; \mathbf{v}_0^2; \dots; \mathbf{v}_{K-1}^2]^\top, \quad (2)$$

where $\mathbf{u}_{i,k}$ are column-wise input speaker embeddings, and \mathbf{v}_k^s are cluster-average speaker embeddings, as described in Section 3.2. Thus, the CNN input becomes a $3K \times N_e$ matrix \mathbf{D}_i , where N_e is the embedding dimension, and K is the number of scales.

To estimate the scale weight, we propose 1-D CNN with 1-D filters. The motivation behind using a 1-D filter is to compare the embedding vectors bin-by-bin such that the learned weights in the filters focus only on the difference between other embeddings for each bin. Two CNN layers are followed by two linear layers and the softmax layer, i.e.,

$$w_k = \frac{e^{f(z_k)}}{\sum_{k=0}^{K-1} e^{f(z_k)}}, \quad (3)$$

where k is the scale index, and $f(z_k)$ is the output of the linear layer (see Fig. 3). The scale weight values w_k in Eq. (3) are multiplied with the cosine similarity in an element-wise manner and create the context vectors for each speaker \mathbf{c}_i^s . We concatenate these two speaker-specific context vectors for each i -th time step into the global context vector $\mathbf{C}_i = [\mathbf{c}_i^1; \mathbf{c}_i^2]$ of length $2K$ (see Fig. 2).

3.2. Scale Weighting Mechanism

Although clustering based diarization shows a competitive performance, it lacks the ability to assign multiple speakers at the same time step, which is needed for overlap-aware diarization. In addition, the clustering-based multi-scale diarization requires a separate training step for multi-scale weights or a high-dimensional grid search. This negatively affects the accuracy of the model for unseen domains during training. Moreover, the performance of a clustering-based multi-scale diarization system is heavily dependent on the scale weights, as shown in [5].

Hence, we derived a scale weighting mechanism that dynamically calculates the scale weights at every time step when a trained sequence model estimates the speaker labels. In the scale weighting system, during the training process, the scale weights are multiplied to the cosine similarity values to create a context vector that is fed to the sequence model (LSTM). The scale weighting mechanism is described in Fig. 2. Let $\mathbf{u}_{i,k}$ be the k -th scale embedding of the input speech signal and i be the time-step index. In addition, \mathbf{v}_k^s is the k -th scale average embedding vector of cluster s from the initial clustering result in Eq. (1), and let $w_{k,i}$ be the k -th scale weight at the i -th time-step index. The scale weight $w_{k,i}$ is applied to the cosine distance between the inputs \mathbf{v}_k^s and $\mathbf{u}_{i,k}$, i.e.,

$$\mathbf{c}_i^s[k] = w_k \cdot \frac{\mathbf{v}_k^s \cdot \mathbf{u}_{i,k}}{\|\mathbf{v}_k^s\| \|\mathbf{u}_{i,k}\|}, \quad (4)$$

where \mathbf{c}_i^s is the context vector for the s -th speaker (see Fig.5).

3.3. Sequence model for speaker label estimation

The weighted cosine similarity vectors are fed to the sequence model. In our proposed system, we employ a two-layer Bi-LSTM. Fig. 5 shows how the context vector is fed to the LSTM for estimating the speaker-wise label. Note that the output layer goes through the sigmoid layer such that the output value ranges from zero to 1 independently from the neighboring values. We use the binary cross-entropy loss to train the model. The binary ground truth label is generated by calculating the overlap between the ground truth speaker timestamps and the base scale segment: If the overlap is greater than 50%, the segment is assigned a value of 1 (see Fig.4).

The speaker label inference for the sessions with more than two speakers is obtained by selecting pairs from the estimated number of speakers and averaging the prediction for the corresponding speaker from the two-speaker models. For example, in the three speaker cases in which speakers A, B, and C exist, we average the two results for A from (A, B) and (A, C). Letting S be the total number of speakers in a session, then the proposed model infers $\binom{S}{2}$ number of speakers and then take average of all $S-1$ output pairs from a certain speaker through the following equation:

$$p(s, i) = \frac{1}{S-1} \sum_{q \in U - \{s\}} \sigma(\mathbf{z}_i^{s,(s,q)}) \quad (5)$$

where U is a complete set of the speakers, and $\sigma(\mathbf{z}_i^{s,(s,q)})$ is the sigmoid output of the i -th base scale segment of the s -th

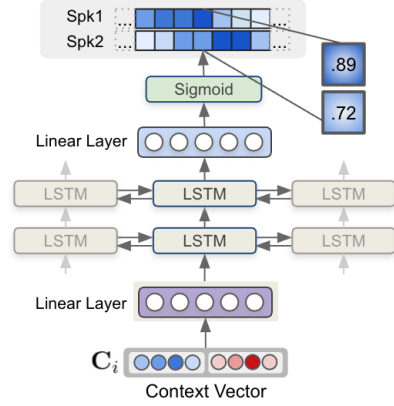


Figure 5: Sequence modeling using LSTM for generating speaker labels.

speaker, which is extracted from the pairing of the s -th and q -th speakers. We use a threshold value T on the average sigmoid value $p(s, i)$ to obtain the final speaker label output for each base scale segment. If both sigmoid values are below the threshold T for the given time step i , we then choose the label from the initial clustering result.

4. Experimental Results

4.1. Models and training

We train the meeting and telephonic models (3.4M parameters) separately using the same dataset but different parameters. Table 1 shows the length of scale and the speaker embedding extractor models. In this study, we use the TitaNet-L (25.3M parameters) model proposed in [13]. The shift lengths for each scale are half of each scale window length. The scale information in Table 1 is applied to both the neural network model and the clustering as initialization. The initial scale weights are set using the following equation:

$$w_k = r - \frac{(r-1)}{K-1}k, \quad (6)$$

where the weights are linearly increased or decreased from the parameter value r and the base scale (the highest index) is always assigned a value of 1.0. In this way, the initial clustering scale weight vector $\mathbf{w} = [w_0, w_1, \dots, w_{K-1}]$ is formed. Note that the overall scale of the scale weights w_k for clustering does not affect the result because the weighted sum of the cosine similarity values are min-max normalized during the clustering process [14]. The proposed system requires parameter tuning on r in Eq. (6) for clustering and threshold T described in Section 3.3 for the MSDD.

4.2. Evaluation setup

An evaluation of a speaker diarization system depends on numerous conditions regarding the use of the development set, whether information on the oracle number of speakers is used, and the diarization error rate (DER) evaluation conditions such as collar and overlap inclusion. We only report the results based on the estimated number of speakers and oracle voice activity detection (VAD). We report two different DER evaluation settings following the name of the setups in [12]: (1) *Forgiving*, DER_A , 0.25-s collar and ignored overlap speech; and (2) *Full*, DER_B , 0-s collar and included overlap speech.

4.3. Datasets

We use following datasets:

Table 1: Scale length and quantity for each model.

Model Domain	K	Scale (Segment) Length
Telephonic	5	[1.5, 1.25, 1.0, 0.75, 0.5]
Meeting	6	[3.0, 2.5, 2.0, 1.5, 1.0, 0.5]

- **CH(NS)**: NIST-SRE-2000 (LDC2001S97) is the most popular diarization evaluation dataset and referred to as CALL-HOME in the speaker diarization papers. To make a fair comparison with other studies and our previous research, a two-fold cross validation is used for tuning the parameters using the split appeared in [15, 16].
- **CH109**: Call Home American English Speech (CHAES, LDC97S42) is a corpus that contains only English telephonic speech data. We evaluate 109 sessions (CH109) which have two speakers per session. The parameters are optimized on the remaining 11 sessions in the CHAES dataset.
- **AMI-MH**: AMI meeting corpus [17] is a meeting speech dataset containing up to five speakers. For the evaluation, we use annotation files and the file lists of splits provided by [18] from the *MixHeadset* part. All parameters are optimized on the dev set.
- **The Fisher corpus**: The Fisher corpus contains 10 minute English conversations. We use Fisher corpus [19] to train the proposed model using the two-speaker setting. We use our own random split that has 10,000 sessions for training and 1,500 sessions for validation.

In terms of the hyperparameters, we use 256 nodes of hidden layer units for both hidden layers in the scale weight generation and the LSTM. We use 16 filters for the 1-D CNN described in 3.1. F1 score is used for stopping the training or select the model. The results reported herein can be reproduced using the NeMo¹ open-source toolkit [20].

4.4. Diarization performance

Table 2 shows the performance of the previously published speaker diarization performance on the same dataset. In Table 2, note that DER_A does not reflect the overlap detection capability. SS-Clus is our previous approach [13], where SS clustering-based diarization systems with TitaNet models [13] were used. The Equal-w-Clus system uses equal weights, which means $w_k=1$ in Eq. (6) for all elements in the scale weight vector \mathbf{w} . Accordingly, Equal-w-MSDD is the result obtained from the MSDD based on the initializing clustering with an equal scale weight. By contrast, Opt-w-Clus indicates a clustering-based diarization result with the optimized scale weight vector \mathbf{w} on each development set, and Opt-w-MSDD refers to the system based on the clustering result Opt-w-Clus. Whereas Opt-w-MSDD shows an overall better result, we want to emphasize the importance of the Equal-w-MSDD result because Equal-w-MSDD does not require parameter tuning for the initializing clustering.

5. Discussions

Although our proposed system shows a competitive performance, there is still room for improvement, which requires further investigation. First, although the trained model can adjust the scale weight on the fly, to obtain improved results, we need to have separate models for both the meeting and telephonic data. Nevertheless, the proposed MSDD model does not show a significantly improved performance on the AMI datasets. We believe this can be overcome by including datasets from many

¹<https://github.com/NVIDIA/NeMo>

Table 2: DER results from previous studies and the proposed methods. The number of speakers N_s is estimated within the range $1 \leq N_s \leq 8$.

Systems	Eval. Setup	Telephonic		Meeting	
		CH (NS)	CH 109	AMI-MH dev	AMI-MH test
Park <i>et al.</i> [5]	DER_A	6.46	2.04*	-	3.32
Aronow- <i>et al.</i> [21]	DER_A	5.1	-	-	-
Dawala- <i>et al.</i> [22]	DER_A	-	-	2.43	4.03
Landini <i>et al.</i> [12]	DER_A	4.42	-	2.14	2.17
	DER_B	21.77	-	22.98	22.85
SS-Clus [13]	DER_A	5.38	1.42	-	1.89
Equal-w-MS-Clus	DER_A	4.57	1.19	1.34	1.06
Equal-w-MSDD	DER_A	4.25	0.69	1.34	1.05
	DER_B	20.39	10.94	22.21	21.18
Opt-w-MS-Clus	DER_A	4.18	0.70	1.30	1.06
Opt-w-MSDD	DER_A	3.92	0.73	1.30	1.06
	DER_B	20.14	10.82	22.20	21.19

* This is *eval* set in the CHAES dataset.

different domains. Second, the temporal resolution is still limited to the shift length (0.25 s) of the finest scale. We observed that reducing the base scale to below 0.5 s only results in a larger error, which can also be an area of improvement.

Despite the aforementioned downsides, we found the following benefit of the MSDD. First, the proposed MSDD approach achieves similar or better results compared to clustering-based diarization results. Second, compared to the clustering based multi-scale diarization, the proposed system can achieve a more stable performance for unknown domains by using an equal scale weight. As investigated in previous multi-scale studies [5], the performance of clustering-based diarization results is heavily dependent on multi-scale weights. However, in terms of the proposed diarization decoder, the averaged speaker embeddings achieve a relatively minor effect from the scale weights used for the clustering. In particular, the smaller dependency on the tuned parameters is a benefit when the diarization system is deployed under real-life scenarios. Third, the proposed system does not rely on an iterative procedure, unlike the system in [7] where the speaker representation is extracted again after the first prediction. Thus, the proposed system can be applied to streaming diarization systems where we first apply clustering for a relatively short amount of time and then predict the diarization result on the incoming buffer in an incremental manner.

6. Conclusions

In this paper, we proposed a multi-scale diarization decoder with a scale weighting mechanism. The proposed system has the following benefits: First, this is the first study applying a multi-scale weighting concept with sequence model (LSTM) based speaker label estimation. Thus, the multi-scale diarization system enables overlap-aware diarization, which cannot be achieved with traditional clustering-based diarization systems. Moreover, because the decoder is based on a clustering-based initialization, the diarization system can deal with a flexible number of speakers. Finally, we showed a superior diarization performance compared to the previous published results. There are two future areas of research regarding the proposed system. First, we plan to implement a streaming version of the proposed system by implementing diarization decoder based on window-wise clustering. Second, the end-to-end optimization from speaker embedding extractor to diarization decoder can be investigated.

7. References

- [1] S. Chen, P. Gopalakrishnan *et al.*, “Speaker, environment and channel change detection and clustering via the bayesian information criterion,” in *Proc. DARPA broadcast news transcription and understanding workshop*, vol. 8. Virginia, USA, 1998, pp. 127–132.
- [2] M. Senoussaoui, P. Kenny, T. Stafylakis, and P. Dumouchel, “A study of the cosine distance-based mean shift for telephone speech diarization,” *IEEE/ACM TASLP*, vol. 22, no. 1, pp. 217–227, 2013.
- [3] G. Sell, D. Snyder, A. McCree, D. Garcia-Romero, J. Villalba, M. Maciejewski, V. Manohar, N. Dehak, D. Povey, S. Watanabe *et al.*, “Diarization is hard: Some experiences and lessons learned for the JHU team in the inaugural DIHARD challenge.” in *Proc. INTERSPEECH*, Sep. 2018, pp. 2808–2812.
- [4] P. Kenny, D. Reynolds, and F. Castaldo, “Diarization of telephone conversations using factor analysis,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 4, no. 6, pp. 1059–1070, 2010.
- [5] T. J. Park, M. Kumar, and S. Narayanan, “Multi-scale speaker diarization with neural affinity score fusion,” in *Proc. ICASSP*. IEEE, 2021, pp. 7173–7177.
- [6] Y. Fujita, N. Kanda, S. Horiguchi, Y. Xue, K. Nagamatsu, and S. Watanabe, “End-to-end neural speaker diarization with self-attention,” in *2019 ASRU*. IEEE, 2019, pp. 296–303.
- [7] I. Medennikov, M. Korenevsky, T. Prisyach, Y. Khokhlov, M. Korenevskaya, I. Sorokin, T. Timofeeva, A. Mitrofanov, A. Andrusenko, I. Podluzhny *et al.*, “Target-speaker voice activity detection: A novel approach for multi-speaker diarization in a dinner party scenario,” *Proc. INTERSPEECH*, pp. 274–278, 2020.
- [8] S. Horiguchi, Y. Fujita, S. Watanabe, Y. Xue, and K. Nagamatsu, “End-to-end speaker diarization for an unknown number of speakers with encoder-decoder based attractors,” *Proc. INTERSPEECH*, pp. 269–273, 2020.
- [9] S. Watanabe, M. Mandel, J. Barker, E. Vincent, A. Arora, X. Chang, S. Khudanpur, V. Manohar, D. Povey, D. Raj *et al.*, “Chime-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings,” in *Proc. CHiME*, 2020, pp. 1–7.
- [10] N. Ryant, K. Church, C. Cieri, J. Du, S. Ganapathy, and M. Liberman, “Third dihard challenge evaluation plan,” *arXiv preprint arXiv:2006.05815*, 2020.
- [11] K. Kinoshita, M. Delcroix, and N. Tawara, “Integrating end-to-end neural and clustering-based diarization: Getting the best of both worlds,” in *Proc. ICASSP*. IEEE, 2021, pp. 7198–7202.
- [12] F. Landini, J. Profant, M. Diez, and L. Burget, “Bayesian hmm clustering of x-vector sequences (vbv) in speaker diarization: theory, implementation and analysis on standard tasks,” *Computer Speech & Language*, vol. 71, p. 101254, 2022.
- [13] N. R. Koluguri, T. Park, and B. Ginsburg, “Titanet: Neural model for speaker representation with 1d depth-wise separable convolutions and global context,” *arXiv preprint arXiv:2110.04410*, 2021.
- [14] T. J. Park, K. J. Han, M. Kumar, and S. Narayanan, “Auto-tuning spectral clustering for speaker diarization using normalized maximum eigengap,” *IEEE Signal Processing Letters*, vol. 27, pp. 381–385, 2019.
- [15] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, “X-vectors: Robust DNN embeddings for speaker recognition,” in *Proc. ICASSP*, Apr. 2018, pp. 5329–5333.
- [16] D. Snyder, “Callhome diarization recipe using x-vectors,” Github, May 4, 2018. [Online]. Available: https://david-ryan-snyder.github.io/2018/05/04/model_callhome_diarization_v2.html, [Accessed Mar. 21, 2022].
- [17] I. McCowan, J. Carletta, W. Kraaij, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos *et al.*, “The ami meeting corpus,” in *Proceedings of the 5th international conference on methods and techniques in behavioral research*, vol. 88. Citeseer, 2005, p. 100.
- [18] H. Bredin, R. Yin, J. M. Coria, G. Gelly, P. Korshunov, M. Lavechin, D. Fustes, H. Titeux, W. Bouaziz, and M.-P. Gill, “Pyannote.audio: neural building blocks for speaker diarization,” in *Proc. ICASSP*. IEEE, 2020, pp. 7124–7128.
- [19] C. Cieri, D. Miller, and K. Walker, “The fisher corpus: A resource for the next generations of speech-to-text.” in *Proc. LREC*, vol. 4, 2004, pp. 69–71.
- [20] O. Kuchaiev, J. Li, H. Nguyen, O. Hrinchuk, R. Leary, B. Ginsburg, S. Krivan, S. Beliaev, V. Lavrukhin, J. Cook *et al.*, “Nvidia nemo,” <https://github.com/NVIDIA/NeMo>, [Accessed Mar. 21, 2022].
- [21] H. Aronowitz, W. Zhu, M. Suzuki, G. Kurata, and R. Hoory, “New advances in speaker diarization.” in *Proc. INTERSPEECH*, 2020, pp. 279–283.
- [22] N. Dawalatabad, M. Ravanelli, F. Grondin, J. Thienpondt, B. Desplanques, and H. Na, “Ecapa-tdnn embeddings for speaker diarization,” *arXiv preprint arXiv:2104.01466*, 2021.