

DIVE: End-to-end Speech Diarization via Iterative Speaker Embedding

Neil Zeghidour, Olivier Teboul and David Grangier

Abstract—We introduce DIVE, an end-to-end speaker diarization algorithm. Our neural algorithm presents the diarization task as an iterative process: it repeatedly builds a representation for each speaker before predicting the voice activity of each speaker conditioned on the extracted representations. This strategy intrinsically resolves the speaker ordering ambiguity without requiring the classical permutation invariant training loss. In contrast with prior work, our model does not rely on pretrained speaker representations and optimizes all parameters of the system with a multi-speaker voice activity loss. Importantly, our loss explicitly excludes unreliable speaker turn boundaries from training, which is adapted to the standard collar-based Diarization Error Rate (DER) evaluation. Overall, these contributions yield a system redefining the state-of-the-art on the standard CALLHOME benchmark, with 6.7% DER compared to 7.8% for the best alternative.

Index Terms—diarization, speech, end-to-end

I. INTRODUCTION

Speech diarization is the task of annotating speaker turns in a conversation [1], [2], [3]. It is both a crucial step for downstream tasks such as automatic transcription of conversational speech, as well as a challenge as it requires handling long-term dependencies. Traditional systems typically split the problem in three sub-problems. First, a model is trained to extract short-term speaker embeddings. Such embeddings can be i-vectors derived from a Gaussian Mixture Model [4], [5], [6], [7], [8], [9], or embeddings produced by a neural network [10], [11], [12], [13], [14]. Then, given a sequence to be diarized, a pre-trained speech activity detection algorithm [15], [16] extracts active timesteps from the sequence and removes silences. Eventually, a clustering algorithm runs on top of these embeddings to assign each timestep to a speaker. Such composite systems have two main limitations. First, the speaker representations are not optimized for diarization, and may not extract relevant features for disambiguating speakers in e.g. presence of overlap. Moreover, most clustering algorithms being unsupervised, they cannot benefit from the fine-grained annotations of speaker turns in diarization datasets.

This motivated recent advances towards end-to-end diarization systems [17], [18]. In particular, [17], [19] propose to cast the diarization task as a multi-label classification problem. By training a model to predict whether each speaker is active at each timestep, a single model jointly performs speech activity detection (silence vs speech), speaker modelling and clustering. This framework has been used to train various architectures including LSTMs [17] and self-attention [20]

models [21]. Since diarization is a permutation-invariant problem (any permutation of the predicted speakers is valid), these models use Permutation-Invariant Training (PIT) [22], [19], [23] to avoid penalizing the model for choosing a particular speaker ordering. [24] has shown that PIT can suffer from inconsistent assignments when applied to long sequences, and that it is preferable to explicitly learn long-term speaker representations. Moreover, fine-grained annotations can be unreliable around speaker turn boundaries, such that the standard is to remove the neighborhood of boundaries from evaluation [12]. As a consequence, inconsistent supervision around boundaries during training can adversely affect the final accuracy of the system.

In this work we introduce DIVE (Diarization by Iterative Voice Embedding), an end-to-end neural diarization system. DIVE combines three modules which are trained jointly: projection of the waveform to an embedding space, iterative selection of long-term speaker representations, and per-speaker per-timestep voice activity detection. The iterative speaker selection process addresses the problem of speaker order ambiguity and removes the need for training with PIT, similarly to attractor-based approaches [25], [26]. Moreover, we introduce collar-aware training, a modification to the standard multi-label classification loss which ignores errors in a defined radius around speaker turn boundaries to match the evaluation setting. DIVE obtains a state-of-the-art Diarization Error Rate (DER) of 6.7% on CALLHOME[27]. We also perform ablation studies that demonstrate the benefits of collar-aware training, and analyze the patterns of errors of our system.

II. METHOD

A. Setting and notations

We consider a single channel recording $x \in \mathbb{R}^L$ of N , partially overlapping, speakers, with L the length of the sequence. The goal of speech diarization is to produce per-speaker voice activity masks $y_i \in \{0, 1\}^T$ for $i = 1, \dots, N$, with $y_{i,t} = 1$ meaning that speaker i is active at time t , and conversely. Typically, $T < L$ as the model does not produce voice activity masks at the sampling rate of the audio but rather at a lower sampling rate, e.g. every millisecond. DIVE cascades three components. First, a *temporal encoder* projects the input waveform to a downsampled embedding space. Then, the *speaker selector* identifies one embedding (the speaker vector) that characterizes well each speaker, in an iterative fashion. Eventually, the *voice activity detector* consumes the embeddings produced by the temporal encoder as well as the selected speaker vectors and produces a binary voice activity

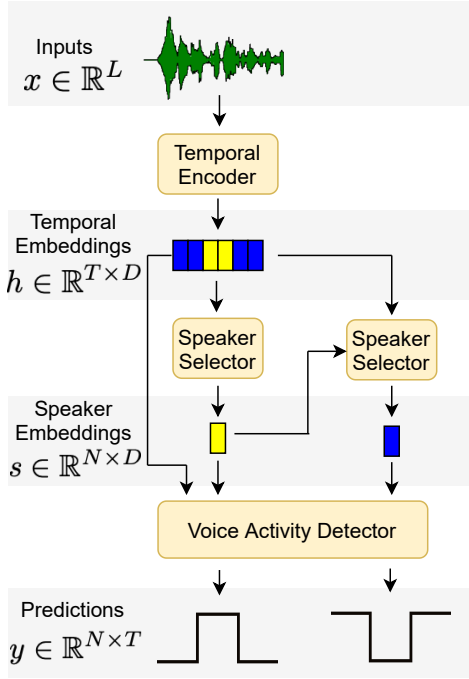


Fig. 1. DIVE for 2 speaker diarization. The temporal encoder first extracts a local speaker representation. The speaker selector iteratively selects the representation of a novel speaker when a single speaker is active. The voice activity module predicts speaker activity conditioned on the input signal and the selected representation.

mask for each speaker. We train these three modules jointly. In the following, we describe each component.

B. Temporal encoder

The *temporal encoder* projects the input waveform x to an embedding space, while performing downsampling. Precisely, the temporal encoder h produces T latent vectors of dimension D , i.e. $h(x) \in \mathbb{R}^{T \times D}$. We refer to these vectors as *temporal embeddings*. We use a temporal encoder similar to that of Wavesplit [24], which cascades residual blocks of dilated 1D-convolutions, with PReLU activations [28] and Layer Normalization [29]. Unlike Wavesplit [24], in which the temporal encoder maintains the original sampling rate of the signal, our temporal encoder performs downsampling by introducing 1D average pooling layers between residual blocks. As the length of the audio sequence L varies between examples, training on batches requires either truncating or padding sequences to a standard length. Given a batch of sequences, a typical training scheme is to randomly sample a fixed-length window from each sequence and to batch the resulting segments [30], [24]. As such windows are typically short (a few seconds), they are likely to only contain one to two speaker turns. This is not appropriate for training a diarization system that needs to model transitions between speaker turns and maintain long term consistency in speaker assignments. To address this issue, we instead sample W fixed-length windows per sequence, pass them through the temporal encoder, and then concatenate them along the temporal axis. This allows for more diversity and more speaker turns inside a single training example. Section III-D assesses the advantage of multi-window training.

C. Iterative speaker selector

The iterative speaker selector outputs a *speaker embedding* for each speaker detected in the signal. Iteratively, it takes as input the temporal embeddings h and a representation of the previously selected speakers μ and outputs a representation of a single non-selected speaker s along with a confidence score c . This process is repeated for two iterations in our two-speaker experiments but can be generalized to a variable number of iterations, stopping when the confidence drops below a given threshold. The representation of the previously selected speakers μ is defined recursively as the average of the embeddings of each previously selected speaker, starting with the zero vector for the first iteration. At each iteration i , this average embedding μ_i is mapped to a 4-by- D matrix $g_\mu(\mu_i)$ with a fully-connected network. This matrix represents a 4-class linear classifier to map each temporal embedding h_t to an event type e_t among four possibilities: a single novel speaker is active, a single already selected speaker is active, overlapped speech, and silence, i.e.

$$P(e_t|h_t, \mu_i) = \text{softmax}(g_\mu(\mu_i)g_h(h_t))$$

where g_h is a fully-connected network. At test time, the confidence of a selection iteration corresponds to the maximal confidence in the presence of a new speaker, i.e.

$$c_i(t_i^*) = \max_t c_i(t) = \max_t P(e_t = \text{novel speaker}|h_t, \mu_i)$$

and the speaker embedding corresponds to the temporal embedding where the maximal confidence is reached, i.e. $s_i = h_{t_i^*}$. During training, we do not output $h_{t_i^*}$ but instead a vector h_t sampled uniformly from times with a novel speaker marked as active in the labels. This allows visiting a larger set of speaker representations. The learning process is supervised and the parameters of the model minimize the negative log likelihood of the 4-way classifier,

$$\mathcal{L}_{\text{selector}}(h, \mu) = -\frac{1}{TN} \sum_{t=1}^T \sum_{i=1}^N \log P(e_t|h_t, \mu_i).$$

D. Voice activity detector

After selecting speaker embeddings, the last module of DIVE predicts the voice activity of each speaker $y_i \in \{0, 1\}^T$ for $i = 1, \dots, N$. The voice activity detector contains two parallel fully-connected neural networks f_h and f_s with PReLU [28] activations and Layer Normalization [29], except for the last layer which is a linear projection. To produce the voice activity $y_{i,t}$ of speaker i at timestep t , f_h and f_s project the current temporal embedding $h_t \in \mathbb{R}^D$ and the speaker vectors $[s_i; \bar{s}] \in \mathbb{R}^{2D}$ respectively:

$$\hat{y}_{i,t} = f_h(h_t)^\top f_s([s_i; \bar{s}]). \quad (1)$$

Here, $[s_i; \bar{s}]$ is the concatenation along the channel axis of s_i , the speaker vector of speaker i , and $\bar{s} = \frac{1}{N} \sum_{j=1}^N s_j$ the mean of all speaker vectors. Intuitively, this means that when predicting the voice activity of a speaker at a given time, we use three pieces of information: the temporal embedding that represents the current speech content, a speaker embedding that represents the identity of the speaker of interest, and

another embedding that represents all speakers. The latter allows the classifier to exploit contrasts between the current speaker of interest and other speakers in the sequence.

During training, we cast the problem of per-speaker, per-timestep voice activity detection as independent binary classification tasks and backpropagate the following loss:

$$\mathcal{L}_{\text{vad}}(\hat{y}, y) = -\frac{1}{TN} \sum_{t=1}^T \sum_{i=1}^N \log(\sigma(\hat{y}_{i,t}(2y_{i,t} - 1))). \quad (2)$$

E. Collar-aware training

When evaluating a diarization system in terms of DER, it is common to apply a *collar*, which is a tolerance around speaker boundaries such that the metric does not penalize the model for small annotation errors. A typical value for such a tolerance is 250ms on each side of a speaker turn boundary (500ms in total). Since we evaluate the model in these conditions, it would be beneficial to train it in a similar fashion i.e. to ignore errors within the collar tolerance. Thus, and as an additional contribution to the DIVE architecture, we propose a training scheme for supervised diarization systems. During training, when computing the loss of the voice activity detector, we remove the loss of frames that fall inside a collar from the total loss and backpropagate the resulting masked loss:

$$\mathcal{L}_{\text{vad}}^{\text{collar}}(\hat{y}, y) = -\frac{1}{TN} \sum_{\substack{t=1 \\ t \notin B_r}}^T \sum_{i=1}^N \log(\sigma(\hat{y}_{i,t}(2y_{i,t} - 1))), \quad (3)$$

with B_r the set of frames that lie within a radius r around speaker turn boundaries. The effect of *collar-aware training* is illustrated in Figure 3. In Section III-C, we show that training with the same collar as used for evaluation substantially improves the DER of the system. The total loss minimized by DIVE is therefore:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{selector}} + \mathcal{L}_{\text{vad}}^{\text{collar}}, \quad (4)$$

and is used to train jointly the temporal encoder, iterative speaker selector and voice activity detector.

III. EXPERIMENTS

We train our models on the "Fisher English Training Speech" Part 1 [31] and Part 2 [32], two datasets of conversational telephone speech. Since they contain clean sequences, we simulate noisy situations by adding background noise from the "noise" part of MUSAN [33]. More precisely, when sampling a training speech sequence, we also sample a random background noise. We renormalize the energy of both the speech and noise sequences, sample a gain uniformly in $[-20, 20]$ dB and apply it to the background noise before adding it to the speech sequence. We evaluate our models on the two-speaker evaluation of CALLHOME [27], a multilingual conversational speech dataset. Following the standard of [21], [25], [34] we report Diarization Error Rates averaged over the 148 test sequences. However, and unlike [21], we do not fine-tune our model on the 155 sequences of the "adaptation" set, but rather use it for hyperparameter selection. DERs are computed using the pyannote library [35].

TABLE I
DIARIZATION ERROR RATE (DER) IN % ON THE TEST SET OF CALLHOME. ALL MODELS ARE EVALUATED WITH A 250MS COLLAR. "NO OVERLAP" MEANS THAT THE EVALUATION EXCLUDES OVERLAPPED SPEECH.

Model	OVERLAP	NO OVERLAP
UIS-RNN V1 [37]	–	10.6
UIS-RNN V2 [37]	–	9.6
UIS-RNN V3 [37]	–	7.6
x-vector + LSTM [12]	–	6.6
BLSTM-EEND [17]	23.1	–
SA-EEND [21]	9.5	–
SA-EEND-EDA [25]	8.1	–
SA-EEND-EDA + Frame Selection [34]	7.8	–
DIVE	6.7	5.9

A. Hyperparameters

The temporal encoder first reduces the length T of temporal embeddings with a 1D-Convolution with a kernel of size 16 and a stride of 8. It then cascades 4 blocks of 10 dilated convolution layers with kernel size 3 and stride 1. The dilation factor δ_l at layer l follows the pattern of [30], [24], i.e. $\delta_l = 2^{l \bmod 10}$, which means that we reinitialize the dilation factor at the beginning of each block. Between the first two blocks, we perform average pooling with kernel size 3 and stride 2. Thus, the total downsampling factor of the model is 16 ($T = L/16$). All convolutional layers use 512 feature maps. The two branches g_μ and g_h of the iterative speaker selector, as well as those (f_h and f_s) of the voice activity detector have two hidden layers with 512 feature maps. We train our model with Adam [36] and a batch size of 512, using an initial learning rate of 0.0003, decayed by a factor of 0.7 every 50000 batches. We use multi-window training with 6 windows of length 32000 samples each.

B. CALLHOME

Table I reports the DER on the test set of CALLHOME. The UIS-RNN [37] is a hybrid system training an RNN on top of pre-trained speaker embeddings, with the V3 being trained on a proprietary dataset with 138k speakers. Similarly, [12] trains an LSTM to model the similarity between pre-trained speaker embeddings and performs diarization. Both models are evaluated without considering overlapped speech, and the latter uses oracle speech activity labels (removing silences). Table I shows that DIVE outperforms both systems in this condition, reaching 5.9% DER, even though DIVE is trained in an end-to-end fashion, without any speaker label and without oracle speech activity annotations. BLSTM-EEND [17] trains a bidirectional LSTM for per-speaker per-timestep voice activity detection, with an additional speaker clustering loss, similar in spirit to DIVE, with a DER of 23.1%. SA-EEND [21] replaces the LSTM by self-attention [20] and removes the deep clustering loss, with the best variation reaching 7.8% thanks to vast training data and fine-tuning on the CALLHOME validation set. DIVE outperforms these models, with a 6.7% DER, and despite not being fine-tuned on CALLHOME. In Table I, the results for DIVE are obtained with a 11-frame median filtering on top of the predictions of the model, as suggested in [17]. This avoids predicting

TABLE II
LABELS VS PREDICTIONS CONTINGENCY (%) FOR FRAME-WISE
DIARIZATION ON CALLHOME.

		Labels			
		Spkr. 1	Spkr. 2	Overlap	Silence
Prediction	Spkr. 1	49.6	0.9	1.8	3.5
	Spkr. 2	0.6	18.8	1.5	2.1
	Overlap	4.1	3.3	8.4	0.9
	Silence	0.7	0.4	0.0	3.3
Class prior		55.1	23.3	11.8	9.8

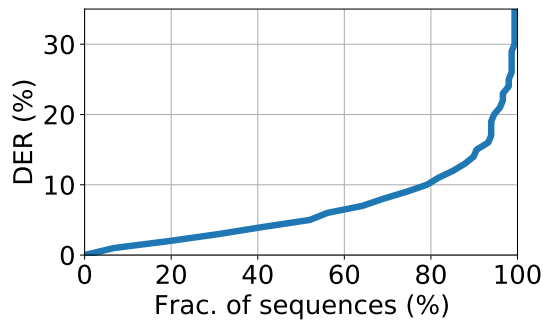


Fig. 2. Cumulative distribution of the Diarization Error Rate (DER) in % on CALLHOME, with the standard 250ms collar.

non-existing, extremely short segments. Without this median filtering, the DER of DIVE goes up from 6.7% to 6.8%, which shows that the model’s predictions are already reliable.

Table II analyzes error types. We observe few confusion errors where a speaker is mistaken for another (1.5%); most errors concentrate on mistaking single speaker activity for overlapped speech (7.4%), mistaking overlap for single speaker activity (3.3%) and mistaking silence for speaker activity (5.6%). Figure 2 plots the cumulative distribution of DER and shows that the median DER is below 5 while the average is higher due to a minority (5%) with DER over 20.

C. Impact of collar-aware training

Figure 3 shows the impact of collar-aware training. When using the standard loss function of Equation 2, the raw DER decreases steadily. On the other hand, when using the collar-aware loss defined in Equation 3, the raw DER plateaus early in training, but its DER with a 250ms collar converges faster and to a better score than its standard counterpart. This shows that when the target evaluation metric uses a collar, it is beneficial to integrate this tolerance into the training loss.

D. Impact of training on multiple windows

Diarization requires speaker representations that are reliable throughout long sequences of speech, e.g. several minutes. However, training a neural network over long speech sequences is slow since it prevents from training with a large batch size due to memory constraints. To solve this problem, our training process samples multiple short windows from the same sequence: this allows DIVE to observe the same speakers over snippets far apart in time while keeping memory usage low. Figure 4 illustrates the benefit of multi-window training:

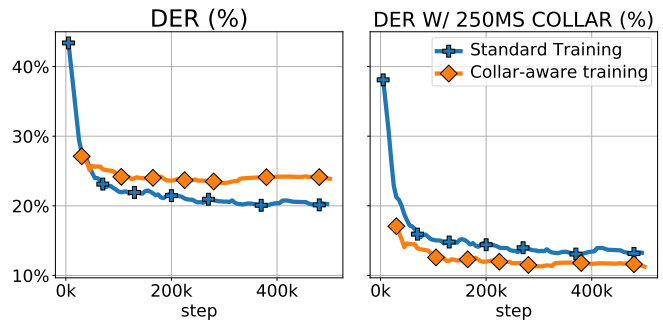


Fig. 3. Diarization Error Rate (DER) in % with and without collar-aware training, on the validation set of CALLHOME. On the left is the raw DER, that penalizes every error. On the right, the DER with the standard 250ms collar.

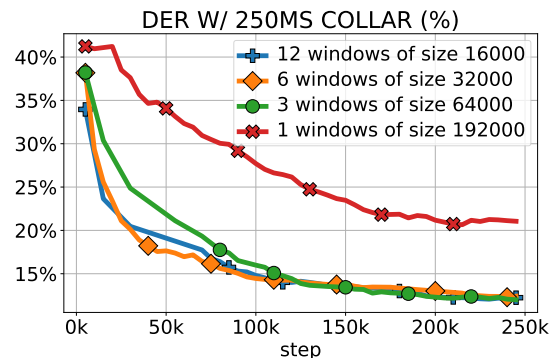


Fig. 4. Diarization Error Rate (DER) in % when varying the number and size of windows for a constant total of 192000 samples.

for a fixed budget of 192000 samples per training example, splitting it into several windows performs much better than using a single, contiguous window.

IV. CONCLUSIONS

This paper introduces DIVE, an end-to-end model for speaker diarization. DIVE decomposes the task into three stages: convolutional temporal encoding, iterative speaker selection and speaker-conditioned voice activity prediction. The iterative speaker selector repeatedly processes the whole sequence to select a representation of a speaker not selected during the previous iterations. The extracted representations condition voice activity prediction. This formulation resolves the ambiguity in speaker order and offers a generic formulation regardless of the number of speakers per sequence. The model does not rely on pretrained components and all parameters are trained to optimize the voice activity likelihood with a novel collar-aware loss function. This loss does not rely on supervision from unreliable speaker turn boundaries, and matches standard collar-aware evaluations. DIVE establishes a new state-of-the-art on the standard CALLHOME benchmark, with 6.7% DER compared to 7.8% for the best alternative. In the future, we aim to address experimental settings with variable number of speakers and noisier acoustic conditions [38], [39].

REFERENCES

- [1] S. E. Tranter and D. A. Reynolds, "An overview of automatic speaker diarization systems," *IEEE Trans. Speech Audio Process.*, vol. 14, no. 5, pp. 1557–1565, 2006. [Online]. Available: <https://doi.org/10.1109/TASL.2006.878256>
- [2] X. A. Miró, S. Bozonnet, N. W. D. Evans, C. Fredouille, G. Friedland, and O. Vinyals, "Speaker diarization: A review of recent research," *IEEE Trans. Speech Audio Process.*, vol. 20, no. 2, pp. 356–370, 2012. [Online]. Available: <https://doi.org/10.1109/TASL.2011.2125954>
- [3] T. J. Park, N. Kanda, D. Dimitriadis, K. J. Han, S. Watanabe, and S. Narayanan, "A review of speaker diarization: Recent advances with deep learning," *CoRR*, vol. abs/2101.09624, 2021. [Online]. Available: <https://arxiv.org/abs/2101.09624>
- [4] H. Ning, M. Liu, H. Tang, and T. S. Huang, "A spectral clustering approach to speaker diarization," in *INTERSPEECH 2006 - ICSLP, Ninth International Conference on Spoken Language Processing, Pittsburgh, PA, USA, September 17-21, 2006*. ISCA, 2006. [Online]. Available: http://www.isca-speech.org/archive/interspeech_2006/i06_1607.html
- [5] S. Shum, N. Dehak, R. Dehak, and J. R. Glass, "Unsupervised methods for speaker diarization: An integrated and iterative approach," *IEEE Trans. Speech Audio Process.*, vol. 21, no. 10, pp. 2015–2028, 2013. [Online]. Available: <https://doi.org/10.1109/TASL.2013.2264673>
- [6] S. H. Shum, N. Dehak, R. Dehak, and J. R. Glass, "Unsupervised methods for speaker diarization: An integrated and iterative approach," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 10, pp. 2015–2028, 2013.
- [7] G. Sell and D. Garcia-Romero, "Speaker diarization with plda i-vector scoring and unsupervised calibration," in *2014 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, 2014, pp. 413–417.
- [8] D. Dimitriadis and P. Fousek, "Developing on-line speaker diarization system," in *Interspeech 2017, 18th Annual Conference of the International Speech Communication Association, Stockholm, Sweden, August 20-24, 2017*, F. Lacerda, Ed. ISCA, 2017, pp. 2739–2743. [Online]. Available: http://www.isca-speech.org/archive/Interspeech_2017/abstracts/0166.html
- [9] D. Kounades-Bastian, L. Girin, X. Alameda-Pineda, S. Gannot, and R. Horaud, "An EM algorithm for joint source separation and diarisation of multichannel convolutive speech mixtures," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2017, New Orleans, LA, USA, March 5-9, 2017*. IEEE, 2017, pp. 16–20. [Online]. Available: <https://doi.org/10.1109/ICASSP.2017.7951789>
- [10] D. Garcia-Romero, D. Snyder, G. Sell, D. Povey, and A. McCree, "Speaker diarization using deep neural network embeddings," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 4930–4934.
- [11] Q. Wang, C. Downey, L. Wan, P. A. Mansfield, and I. L. Moreno, "Speaker diarization with lstm," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5239–5243.
- [12] Q. Lin, R. Yin, M. Li, H. Bredin, and C. Barras, "Lstm based similarity measurement with spectral clustering for speaker diarization," *arXiv preprint arXiv:1907.10393*, 2019.
- [13] K. Kinoshita, M. Delcroix, and N. Tawara, "Integrating end-to-end neural and clustering-based diarization: Getting the best of both worlds," *CoRR*, vol. abs/2010.13366, 2020. [Online]. Available: <https://arxiv.org/abs/2010.13366>
- [14] Q. Li, F. L. Kreyssig, C. Zhang, and P. C. Woodland, "Discriminative neural clustering for speaker diarisation," in *IEEE Spoken Language Technology Workshop, SLT 2021, Shenzhen, China, January 19-22, 2021*. IEEE, 2021, pp. 574–581. [Online]. Available: <https://doi.org/10.1109/SLT48900.2021.9383617>
- [15] T. Hughes and K. Mierle, "Recurrent neural networks for voice activity detection," in *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013, pp. 7378–7382.
- [16] S. Thomas, G. Saon, M. V. Segbroeck, and S. S. Narayanan, "Improvements to the IBM speech activity detection system for the DARPA RATS program," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2015, South Brisbane, Queensland, Australia, April 19-24, 2015*. IEEE, 2015, pp. 4500–4504. [Online]. Available: <https://doi.org/10.1109/ICASSP.2015.7178822>
- [17] Y. Fujita, N. Kanda, S. Horiguchi, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with permutation-free objectives," *arXiv preprint arXiv:1909.05952*, 2019.
- [18] T. von Neumann, K. Kinoshita, M. Delcroix, S. Araki, T. Nakatani, and R. Haeb-Umbach, "All-neural online source separation, counting, and diarization for meeting analysis," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 91–95.
- [19] Y. Fujita, S. Watanabe, S. Horiguchi, Y. Xue, and K. Nagamatsu, "End-to-end neural diarization: Reformulating speaker diarization as simple multi-label classification," *arXiv preprint arXiv:2003.02966*, 2020.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.
- [21] Y. Fujita, N. Kanda, S. Horiguchi, Y. Xue, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with self-attention," in *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*. IEEE, 2019, pp. 296–303.
- [22] D. Yu, M. Kolbæk, Z.-H. Tan, and J. Jensen, "Permutation invariant training of deep models for speaker-independent multi-talker speech separation," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 241–245.
- [23] Y. Xue, S. Horiguchi, Y. Fujita, S. Watanabe, P. García, and K. Nagamatsu, "Online end-to-end neural diarization with speaker-tracing buffer," in *IEEE Spoken Language Technology Workshop, SLT 2021, Shenzhen, China, January 19-22, 2021*. IEEE, 2021, pp. 841–848. [Online]. Available: <https://doi.org/10.1109/SLT48900.2021.9383523>
- [24] N. Zeghidour and D. Grangier, "Wavesplit: End-to-end speech separation by speaker clustering," *arXiv preprint arXiv:2002.08933*, 2020.
- [25] 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," *arXiv preprint arXiv:2005.09921*, 2020.
- [26] Y. Fujita, S. Watanabe, S. Horiguchi, Y. Xue, J. Shi, and K. Nagamatsu, "Neural speaker diarization with speaker-wise chain rule," *CoRR*, vol. abs/2006.01796, 2020. [Online]. Available: <https://arxiv.org/abs/2006.01796>
- [27] NIST, "2000 speaker recognition evaluation plan," 2000. [Online]. Available: <https://catalog.ldc.upenn.edu/LDC2001S97>
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*.
- [29] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [30] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. W. Senior, and K. Kavukcuoglu, "Wavenet: A generative model for raw audio," in *SSW*. ISCA, 2016, p. 125.
- [31] C. Cieri, D. Graff, O. Kimball, D. Miller, and K. Walker, "Fisher english training speech part 1," 2004. [Online]. Available: <https://catalog.ldc.upenn.edu/LDC2004S13>
- [32] —, "Fisher english training speech part 2," 2005. [Online]. Available: <https://catalog.ldc.upenn.edu/LDC2005S13>
- [33] D. Snyder, G. Chen, and D. Povey, "MUSAN: A Music, Speech, and Noise Corpus," *Tech. Rep.*, 2015. [Online]. Available: <http://www.itl.nist.gov/iad/mig/tests/sre/>
- [34] S. Horiguchi, P. Garcia, Y. Fujita, S. Watanabe, and K. Nagamatsu, "End-to-end speaker diarization as post-processing," *arXiv preprint arXiv:2012.10055*, 2020.
- [35] H. Bredin, "pyannote.metrics: a toolkit for reproducible evaluation, diagnostic, and error analysis of speaker diarization systems," in *Interspeech 2017, 18th Annual Conference of the International Speech Communication Association, Stockholm, Sweden, August 2017*. [Online]. Available: <http://pyannote.github.io/pyannote-metrics>
- [36] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR (Poster)*, 2015.
- [37] A. Zhang, Q. Wang, Z. Zhu, J. Paisley, and C. Wang, "Fully supervised speaker diarization," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6301–6305.
- [38] N. Ryant, P. Singh, V. Krishnamohan, R. Varma, K. Church, C. Cieri, J. Du, S. Ganapathy, and M. Liberman, "The third DIHARD diarization challenge," *CoRR*, vol. abs/2012.01477, 2020. [Online]. Available: <https://arxiv.org/abs/2012.01477>
- [39] S. Watanabe, M. I. Mandel, J. Barker, and E. Vincent, "Chime-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings," *CoRR*, vol. abs/2004.09249, 2020. [Online]. Available: <https://arxiv.org/abs/2004.09249>