# An Efficient Speaker Diarization Pipeline for Conversational Speech

**Wael A. Sultan[1], Mourad S. Semary[1], and Sherif M. Abdou[2]**
[1]Department of Basic Engineering Sciences, Benha Faculty of Engineering, Benha University, Benha, Egypt
[2]Information Technology Department, Faculty of Artificial Intelligence, Cairo University, Cairo, Egypt
**E-mail:** wael.ali@bhit.bu.edu.eg

**Abstract**
In the domain of audio signal processing, the accurate and efficient diarization of conversational speech is still a challenging task, particularly in environments with significant speaker overlap and diverse acoustic scenarios. This paper introduces a comprehensive speaker diarization pipeline that improves performance and efficiency in processing conversational speech. Our pipeline comprises several key components: Voice Activity Detection (VAD), Speaker Overlap Detection (SOD), Speaker Separation models, robust speaker embedding, clustering algorithms, and sophisticated post-processing techniques. Beginning with Voice Activity Detection (VAD), the pipeline efficiently discriminates between speech and non-speech segments, effectively reducing processing overhead. Following VAD, the Speaker Overlap Detection (SOD) component identifies segments featuring speaker overlap. Following this, a speaker separation model separates the overlapping speech into distinct streams. A pivotal enhancement in our pipeline is the integration of robust speaker embedding and clustering techniques, which capture and utilize speaker-specific characteristics to improve the grouping of speech segments. Finally, the post-processing stage refines these segments to ensure temporal consistency and improve the overall diarization accuracy. We evaluated our pipeline across multiple benchmark datasets, proving significant reductions up to 10% in Diarization Error Rate (DER) compared to existing methods.

**Keywords:** speaker diarization, speaker separation, voice activity detection, optimization, pipeline.

## 1. Introduction

Speaker diarization is about figuring out who is speaking and when in a conversation [1]. With more and more conversations being recorded, like in call centers or during meetings, there's a big need for systems that can do this well. However, the inherent challenges posed by conversational speech, including overlapping speech, varying acoustic conditions, and speaker characteristics, make speaker diarization a non-trivial task. Consequently, researchers have continuously sought innovative methodologies to address these challenges and enhance the performance of speaker diarization systems [2], [3].

Over the years, many methodologies have been proposed to tackle the complexities of speaker diarization. Early approaches primarily relied on traditional clustering algorithms such as Gaussian Mixture Models (GMMs) and Hidden Markov Models (HMMs) to model speech features and infer speaker identities. While these methods achieved moderate success, they often struggled with scalability and robustness, especially in real-world scenarios with diverse speaking styles and environmental conditions [4]. In recent years, advancements in machine learning and signal processing have revolutionized the field of speaker diarization, leading to the development of novel techniques that leverage deep neural networks (DNNs), convolutional neural networks (CNNs), and recurrent neural networks (RNNs). These deep learning-based approaches have demonstrated remarkable performance improvements in speaker diarization systems [2].

In this paper, we propose an efficient speaker diarization pipeline tailored to address the mentioned challenges while maintaining high accuracy and scalability. Our approach combines state-of-the-art deep learning techniques with efficient data processing strategies to achieve real-time performance and memory efficiency.

We evaluate our proposed pipeline on two different datasets and compare its performance against existing methods, showing its efficiency in conversational speech scenarios.

## 2. Speaker Diarization Pipeline

A typical speaker diarization pipeline involves several stages such as speech detection to identify speech from silence or noise, speech segmentation to break the speech into manageable chunks, speaker feature extraction to analyze each segment for speaker-related features like the Mel-frequency cepstral coefficients (MFCCs) or neural network embeddings, clustering to group segments by the speaker using algorithms like K-means or GMM, and re-segmentation and adjustment for fine-tuning the segment boundaries to enhance overall accuracy.

We'll outline our approach to constructing such an efficient pipeline in the upcoming sections. Before delving into details, however, we'll introduce our custom dataset and the chosen evaluation metric. This metric will enable us to assess and benchmark the effectiveness of our proposed system against existing solutions.

**Datasets**: To effectively investigate, compare, and fine-tune existing and proposed diarization pipelines, it is crucial to have access to a comprehensive dataset. This dataset should encompass a wide range of variables, including multiple speakers, speaker variability, diverse background noises, varying channel qualities, and differing rates of speaker turn-taking. Such a dataset will enable robust testing and optimization of diarization algorithms across multiple challenging scenarios.

In this study, we utilized two datasets. The first is the publicly available AMI [5] (Augmented Multi-party Interaction) dataset, a widely recognized resource in machine learning and natural language processing, particularly suited for analyzing meeting scenarios. This dataset comprises 100 hours of meeting recordings taken in both structured and natural environments. It includes recordings from real and scripted meetings designed around specific scenarios, encompassing audio, video, and textual data.

Additionally, we developed a private dataset consisting of 100 audio files that vary across calls, meetings, and recorded discussions in different settings. Most of the samples in this dataset feature dual-speaker recordings, with configurations including male/male, female/female, and male/female pairings. These audio samples mimic real-world communication scenarios found in typical calls and meetings, providing a rich resource for fine-tuning diarization systems to handle common communication setups effectively.

**Evaluation Metric:**

The most common evaluation metric for speaker diarization is the diarization error rate (DER), it can be calculated with this form:

$$DER = \frac{\text{False Alarm} + \text{Missed} + \text{Confusion}}{\text{Reference Length}}$$

Where,

• Reference Length: is the total length of the reference (ground truth).

• False Alarm is the length of segments considered speech in hypothesis but not in reference.

• Missed: is the Length of segments that are considered as speech in reference, but not in hypothesis.

• Confusion: is the length of segments that are assigned to different speakers in hypothesis and reference.

For our experiments, we utilized Pyannote.metrics [6] to calculate the Diarization Error Rate (DER).

**3. Proposed Methodology:**

There are a couple of popular open-source toolkits that present diarization pipelines for speaker diarization, namely pyannote.audio [7], and NeMo [8]. Both frameworks share a similar overall pipeline structure and pretrained models. But it's still very hard for both frameworks to generalize a speaker diarization due to many challenges so we propose the following pipeline.
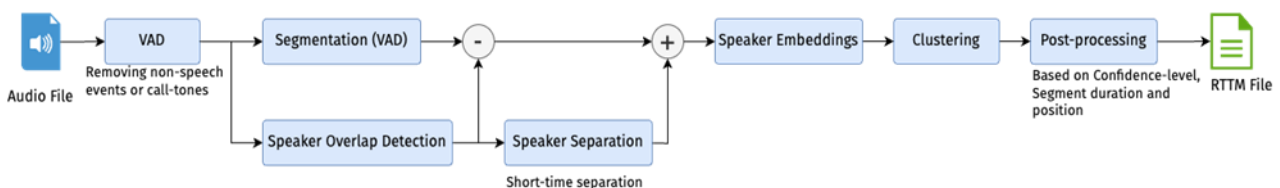
As **Fig (1)**, shows the pipeline takes audio files as input, applies a sophisticated series of processes, and outputs who spoke and when in RTTM formatted files. In the following sections, we'll go through each process individually.

**Speech detection**

The first step in our pipeline is to identify speech intervals within the audio file by removing non-speech segments. This is important because inaccurate filtering of non-speech events, like ringtones, call tones, and music, can introduce ambiguity that makes it harder to later categorize the speakers. We use voice activity detection (VAD) models to accomplish this task. Several pre-trained models have been explored, including Pyannote VAD [7], [9], MarbelNet [10], and SpeechBrain CRDNN [11]. Our experiments show that Pyannote VAD is particularly effective at filtering out non-speech segments, especially ringtones and call tones.

**Segmentation**

After speech detection is completed by the preceding step, the audio stream is segmented into smaller pieces, making it more manageable for feature extraction from each segment. At this stage, employing a VAD model is also beneficial, but unlike the previous step, it is advisable to use a more aggressive VAD that tends to segment the audio into shorter duration segments. We evaluated two potential models for this purpose: the WebRTC VAD [12], developed by Google for its WebRTC project, and Silero-VAD [13], an open-source initiative aimed at creating optimized and production-ready models using C++ and ONNX-runtime libraries. Based on its superior performance in terms of both accuracy and processing speed, Silero-VAD was selected as the preferred model.

**Speaker Overlap-Detection**

In everyday interactions such as casual conversations, calls, or meetings, it is common for individuals to speak simultaneously, resulting in overlapping speech. On the other hand, the speaker embedding process is designed to capture the unique acoustic characteristics of a single speaker, running under the assumption that each speech segment produced by the previous step contains speech from only one speaker. This assumption can lead to potential confusion during the clustering phase, as segments containing multiple speakers may yield poor-quality embeddings.



**Fig (1)** Proposed diarization pipeline

To address this challenge, we propose the implementation of a speaker overlap-detection model (SOD) to identify segments where speech overlaps, to later separate speakers in these segments, or to exclude these segments entirely when constructing embedding vectors. For this purpose, our pipeline has incorporated a pretrained pyannote.audio SOD [7], [9] model, which is specifically designed to handle such tasks.

**Speaker Separation**

In this step, segments identified as having overlapping speech by the preceding SOD model are processed through a speaker separation model to disentangle the overlapping voices. Recent advancements have seen end-to-end neural speech separation models, such as dual-path RNN (DPRNN) [14], Sepformer [15], and ConvTasNet [16], proving robust performance. However, a common drawback of these models is the potential for high latency when processing lengthy audio clips. Fortunately, in our pipeline, only the segments identified with overlapping speech are processed, mitigating this issue. Among these, the ConvTasNet model offers the best balance between effective separation and reduced latency.

Practical observations have shown that many segments flagged as overlapped speech in casual conversations are just one-word interruptions by another speaker. In scenarios where diarization is coupled with speech-to-text models for transcribing conversations, it may sometimes be preferable to disregard these overlapped segments altogether. So, in our pipeline, we opt to ignore segments with overlapping speech if they are shorter than 0.5 seconds.

**Speaker Feature Extraction and Embedding**

This step is critical in the entire process, as it focuses on representing each speech segment in a manner that helps the clear differentiation between segments from different speakers. Over the last decade, many techniques for speaker embedding have been developed, starting from i-vector and extending through various d-vector approaches [3]. However, ECAPA-TDNN-based speaker embedding models have recently appeared highly successful in this area [17], [18].

ECAPA-TDNN utilizes a Time Delay Neural Network (TDNN)-based architecture, enhanced with several key innovations. It incorporates a channel- and context-dependent attention mechanism within the pooling layer, uses 1-dimensional Squeeze-Excitation (SE) blocks, integrates 1-dimensional Res2Net blocks, and employs multi-layer feature aggregation. Moreover, the model leverages AAM-softmax loss for effective classification of speaker identities, enhancing its performance in distinguishing speakers [17], [18].

In the proposed pipeline, segments that appear from the segmentation process - after either separating or excluding segments with overlapping speech - are processed to extract embeddings using the ECAPA-TDNN model. This involves moving a window of 1.5 seconds across the audio with a shift of 0.75 seconds, resulting in the extraction of 192 speaker embedding vectors for each window.

We have implemented a recipe from SpeechBrain [11] for the ECAPA-TDNN and fine-tuned it on the AMI dataset [5]. Subsequently, it was benchmarked against our custom dataset of 100 samples to evaluate its performance.

To further validate and confirm the efficacy of the generated embedding vectors in speaker categorization, we conducted a comparison between vectors produced for all samples in our custom dataset. We use an implementation of a d-vector sourced from this repository[1] [19] and our ECAPA-TDNN model. These vectors are subsequently transformed into a two-dimensional array using the UMAP [20] technique for visualization purposes. The resulting figures (**Figure 2**, and **Figure 3**) clearly demonstrate the distinguishability of the ECAPA-TDNN-based vectors for each speaker.

**Clustering**

Following the generation of embedded vectors in the preceding step, the vectors are subjected to a clustering process, which involves grouping them based on the known number of speakers (referred to as the oracle) or, alternatively, when the number of speakers is unknown. While several straightforward clustering techniques such as K-means or K-nearest neighbors (KNN) can be applied, Spectral clustering (SC) [21] has emerged as particularly effective in the domain of speaker diarization. Spectral clustering leverages the spectral properties of the affinity matrix to partition the data into clusters. For our implementation, we have adopted a recipe provided by SpeechBrain for Spectral clustering[2], which provides robust and efficient clustering performance tailored to speaker diarization tasks.

**Post-processing**

In the concluding phase of the pipeline, a post-processing step is implemented to refine the clustering results by considering several factors. Mainly, the confidence level of the clustering assignment is assessed, and the duration of the current segment relative to the durations of the segments preceding and following it is considered. This analysis aims to mitigate the impact of short interruptions from other speakers, which may introduce ambiguities in the clustering process. By assessing the contextual continuity of speaker turns and the relative lengths of speech segments, this post-processing step helps to enhance the accuracy and coherence of the final speaker diarization output as shown in the following **Fig (4)**.
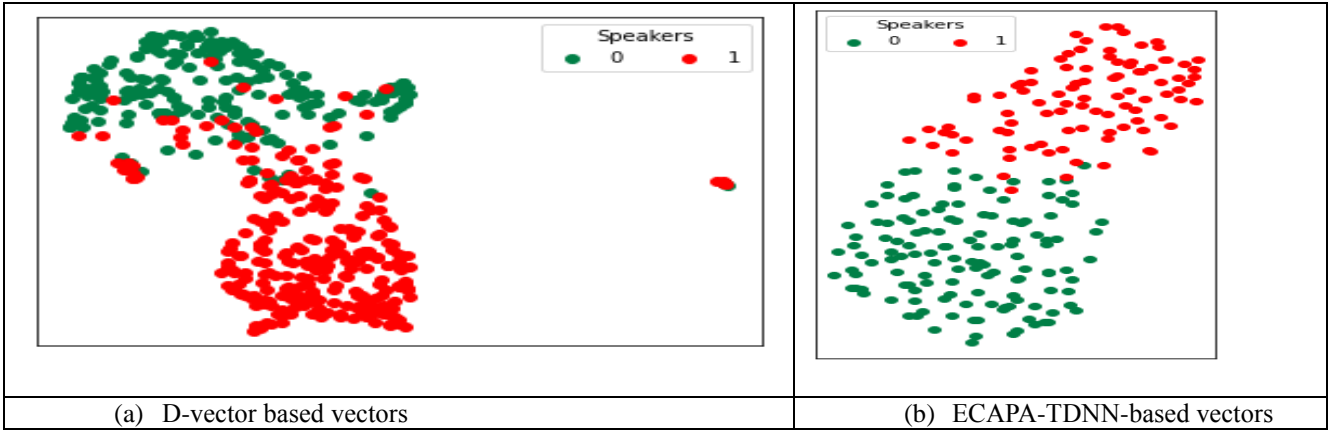
---

[1] https://github.com/hitachi-speech/EEND
[2] https://github.com/speechbrain/speechbrain/tree/develop/recipes/AMI/Diarization

| (a)   D-vector based vectors | (b)   ECAPA-TDNN-based vectors |

**Fig (2)** Visualization of UMP representation for extracted vector as a sample audio file



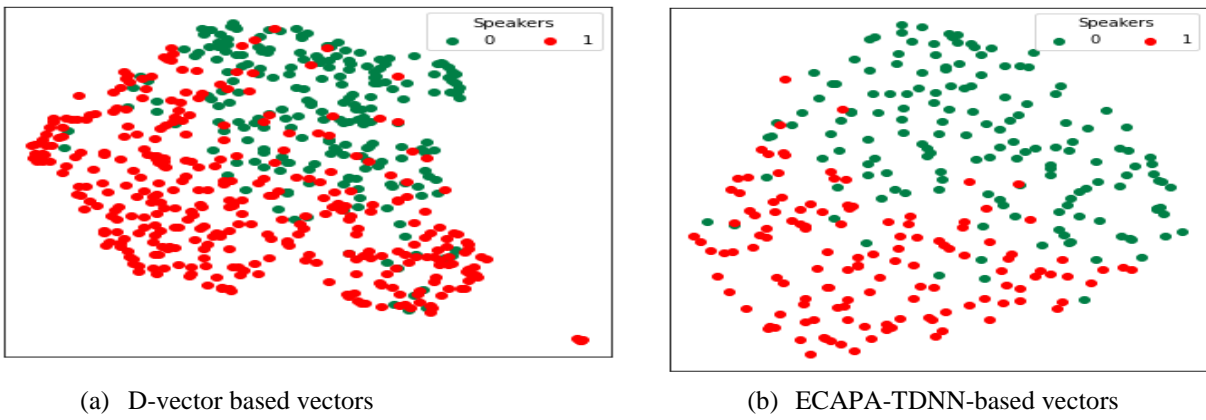(a)   D-vector based vectors          (b)   ECAPA-TDNN-based vectors

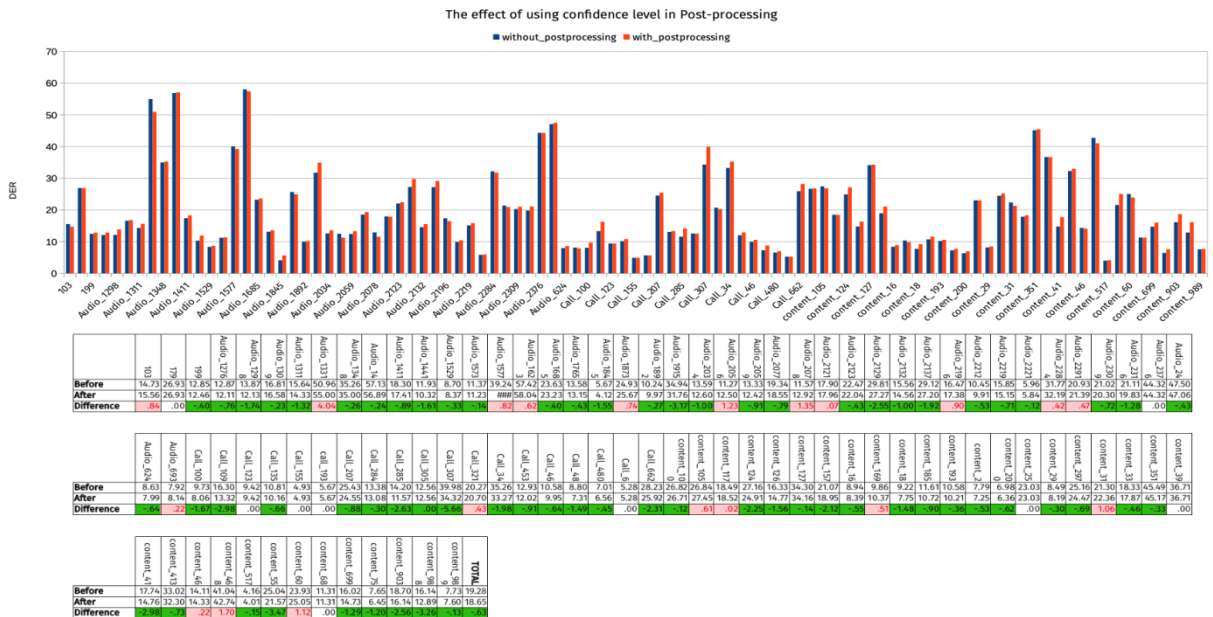**Fig (3)** Visualization of UMP representation for extracted vector as another sample audio file.



**Fig (4)** Post-processing DER analysis.

## 4. Resuts

We thoroughly assessed and evaluated each model employed at every stage of the pipeline, as discussed in earlier sections. To judge the overall performance, we benchmarked the entire system - using both the AMI and our custom datasets - against Pyannote.audio pipeline (version 3.1). The comparison results are summarized in the following **Table (1)**.

As demonstrated by the final table, our proposed pipeline significantly surpasses the existing Pyannote.audio pipeline in performance for both the public and the private datasets as follows:

- **AMI dataset (Mix-Headset)**:
  o **DER**: The proposed pipeline achieved a DER (Diarization Error Rate) of 1.91%, significantly lower than the 18.9% achieved by Pyannote.audio.
- **100-samples Custom Dataset**:
  o **DER**: The proposed pipeline achieved a DER of 16.0%, which is lower than the 26.53% for Pyannote.audio.

  o **Correct**: The proposed pipeline also had higher correctness, at 91.67%, compared to Pyannote.audio's 83.28%.
  o **False alarm**: The proposed pipeline had a lower false alarm rate of 7.66%, compared to Pyannote.audio's 9.81%.
  o **Missed**: The missed rates were similar between both systems, with the proposed pipeline at 1.85% and Pyannote.audio at 1.76%.
  o **Confusion**: The proposed pipeline had significantly lower confusion, at 6.47%, compared to Pyannote.audio's 14.95%.

Overall, the proposed pipeline outperforms Pyannote.audio in all metrics across both datasets. This indicates better diarization accuracy, fewer false alarms, fewer missed segments, and less confusion, making it a more reliable solution.

**Table (1)** Comparsion between the Pyannote.audio and the proposed pipeline

| System | AMI-dataset (Mix-Headset) | Our custom dataset | | | | |
|---|---|---|---|---|---|---|
| | DER | Correct | False alarm | Missed | Confusion | DER |
| **Pyannote.audio [22]** | 18.9% | 83.28% | 9.81 | 1.76% | 14.95% | 26.53% |
| **Proposed Pipeline** | **1.91%** | 91.67% | 7.66% | 1.85% | 6.47% | **16.0%** |

## 5. Conclusions

In this paper, we presented an efficient and comprehensive speaker diarization pipeline tailored to address the complex challenges associated with conversational speech processing. Our proposed pipeline incorporates a variety of state-of-the-art methods, including Voice Activity Detection (VAD), Speaker Overlap Detection (SOD), Speaker Separation, robust speaker embedding techniques, and Spectral clustering, which collectively contribute to its superior performance.

Through testing and benchmarking against the Pyannote.audio system, our pipeline demonstrated notable improvements across key metrics on both the AMI dataset and our custom dataset:

- AMI Dataset: The pipeline achieved a Diarization Error Rate (DER) of 1.91%, significantly lower than the 18.9% achieved by Pyannote.audio.
- Custom Dataset: The pipeline yielded a DER of 16.0%, outperforming the 26.53% recorded by Pyannote.audio.

Furthermore, in comparison with the Pyannote.audio system, our pipeline exhibited superior performance in terms of accuracy, false alarms, missed segments, and confusion across both datasets, showcasing its robustness and generalizability to various real-world audio scenarios.

## 6. Future Work

While the proposed pipeline offers considerable advancements, future research may focus on:

- Fine-tuning: Further refining the model by incorporating more varied datasets to enhance its adaptability and robustness across diverse environments.
- Real-time Processing: Optimizing the pipeline for real-time diarization tasks, especially for live transcription services and interactive voice-based applications.
- Integration: Exploring integration opportunities with other systems, such as speech-to-text or voice recognition.

Overall, our pipeline provides a significant step forward in speaker diarization, particularly in the realm of conversational speech, offering a more accurate, efficient, and reliable solution to the challenges posed by diverse audio processing scenarios.

## References

[1] N. Ryant et al., "The Third DIHARD Diarization Challenge." arXiv, Apr. 05, 2021. Accessed: Apr. 29, 2024. [Online]. Available: http://arxiv.org/abs/2012.01477

[2] 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." arXiv, Nov. 26, 2021. Accessed: Apr. 29, 2024. [Online]. Available: http://arxiv.org/abs/2101.09624

[3] Z. Bai and X.-L. Zhang, "Speaker Recognition Based on Deep Learning: An Overview." arXiv, Apr. 03, 2021. Accessed: Apr. 29, 2024. [Online]. Available: http://arxiv.org/abs/2012.00931

[4] G. Sell and D. Garcia-Romero, "Speaker diarization with plda i-vector scoring and unsupervised calibration," in 2014 IEEE Spoken Language Technology Workshop (SLT), Dec. 2014, pp. 413–417. doi: 10.1109/SLT.2014.7078610.

[5] "AMI Corpus." Accessed: Apr. 29, 2024. [Online]. Available: https://groups.inf.ed.ac.uk/ami/corpus/

[6] "pyannote/pyannote-metrics." pyannote, Apr. 26, 2024. Accessed: Apr. 29, 2024. [Online]. Available: https://github.com/pyannote/pyannote-metrics

[7] H. Bredin et al., "pyannote.audio: neural building blocks for speaker diarization." arXiv, Nov. 04, 2019. doi: 10.48550/arXiv.1911.01255.

[8] "NeMo/tutorials/speaker_tasks/Speaker_Diarization_Inference.ipynb at main · NVIDIA/NeMo," GitHub. Accessed: Apr. 28, 2024. [Online]. Available: https://github.com/NVIDIA/NeMo/blob/main/tutorials/speaker_tasks/Speaker_Diarization_Inference.ipynb

[9] H. Bredin and A. Laurent, "End-to-end speaker segmentation for overlap-aware resegmentation." arXiv, Jun. 10, 2021. Accessed: Apr. 28, 2024. [Online]. Available: http://arxiv.org/abs/2104.04045

[10] F. Jia, S. Majumdar, and B. Ginsburg, "MarbleNet: Deep 1D Time-Channel Separable Convolutional Neural Network for Voice Activity Detection," arXiv.org. Accessed: Apr. 28, 2024. [Online]. Available: https://arxiv.org/abs/2010.13886v2

[11] M. Ravanelli et al., "SpeechBrain: A General-Purpose Speech Toolkit." arXiv, Jun. 08, 2021. doi: 10.48550/arXiv.2106.04624.

[12] "WebRTC," WebRTC. Accessed: Apr. 28, 2024. [Online]. Available: https://webrtc.org/

[13] A. Veysov, "snakers4/silero-vad." Apr. 28, 2024. Accessed: Apr. 28, 2024. [Online]. Available: https://github.com/snakers4/silero-vad

[14] Y. Luo, Z. Chen, and T. Yoshioka, "Dual-path RNN: efficient long sequence modeling for time-domain single-channel speech separation." arXiv, Mar. 27, 2020. Accessed: Apr. 29, 2024. [Online]. Available: http://arxiv.org/abs/1910.06379

[15] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, "Attention is All You Need in Speech Separation," arXiv.org. Accessed: Apr. 29, 2024. [Online]. Available: https://arxiv.org/abs/2010.13154v2

[16] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing Ideal Time-Frequency Magnitude Masking for Speech Separation," IEEEACM Trans. Audio Speech Lang. Process., vol. 27, no. 8, pp. 1256–1266, Aug. 2019, doi: 10.1109/TASLP.2019.2915167.

[17] B. Desplanques, J. Thienpondt, and K. Demuynck, "ECAPA-TDNN: Emphasized Channel Attention, Propagation and Aggregation in TDNN Based Speaker Verification," in Interspeech 2020, Oct. 2020, pp. 3830–3834. doi: 10.21437/Interspeech.2020-2650.

[18] N. Dawalatabad, M. Ravanelli, F. Grondin, J. Thienpondt, B. Desplanques, and H. Na, "ECAPA-TDNN Embeddings for Speaker Diarization," in Interspeech 2021, Aug. 2021, pp. 3560–3564. doi: 10.21437/Interspeech.2021-941.

[19] 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, Oct. 05, 2020. Accessed: Apr. 29, 2024. [Online]. Available: http://arxiv.org/abs/2005.09921

[20] L. McInnes, J. Healy, and J. Melville, "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction." arXiv, Sep. 17, 2020. doi: 10.48550/arXiv.1802.03426.

[21] N. Raghav and M. Sahidullah, "Assessing the Robustness of Spectral Clustering for Deep Speaker Diarization." arXiv, Mar. 21, 2024. Accessed: Apr. 29, 2024. [Online]. Available: http://arxiv.org/abs/2403.14286

[22] H. Bredin, "pyannote.audio 2.1 speaker diarization pipeline: principle, benchmark, and recipe," in 24th INTERSPEECH Conference (INTERSPEECH 2023), Dublin, Ireland: ISCA, Aug. 2023, pp. 1983–1987. doi: 10.21437/Interspeech.2023-105.