



## REVIEW ARTICLE

## REVOLUTIONIZING INTERACTION ANALYSIS: AI-POWERED SPEAKER DIARIZATION FOR ENHANCED COMMUNICATION INSIGHTS

Niraj Kafle<sup>a,b,\*</sup><sup>a</sup> Department of Computer Science and Multimedia, Phoenix College of Management, Maitidevi, Kathmandu, Nepal.<sup>b</sup> Lincoln University College, Petaling Jaya, Malaysia.\*Corresponding Author Email: [kafleniraj@gmail.com](mailto:kafleniraj@gmail.com)

This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

## ARTICLE DETAILS

## Article History:

Received 03 December 2024  
Revised 07 January 2025  
Accepted 10 January 2025  
Available online 18 February 2025

## ABSTRACT

Speaker diarization—the task of partitioning an audio stream into speaker-specific segments—has become essential for analyzing human interactions in real-world settings. Recent breakthroughs in deep learning and large language models (LLMs) have dramatically improved both accuracy and processing speed. This paper presents an in-depth review of AI-powered speaker diarization methods, detailing system architectures, numerical performance metrics, and ethical considerations. Experimental results demonstrate significant improvements in Diarization Error Rates (DER), with our system reducing errors by up to 40% compared to traditional approaches. The applications range from enterprise communications to educational analytics, offering promising advancements for next-generation conversational AI.

## KEYWORDS

Speaker Diarization, Deep Learning, Diarization Error Rate (DER), AI in Speech Processing

## 1. INTRODUCTION

The exponential growth of digital audio content—from conference calls and customer service interactions to online meetings and multimedia broadcasts—has led to an urgent need for automated analysis tools. One critical task in this domain is speaker diarization, which involves segmenting an audio recording and attributing each segment to a specific speaker. Traditional approaches based on statistical models such as Gaussian Mixture Models (GMMs) and Hidden Markov Models (HMMs) have provided a foundation, yet their performance is often hindered by high Diarization Error Rates (DER) when deployed in complex, real-world scenarios (Saon, 2017).

In recent years, deep learning has revolutionized many fields within signal processing and natural language processing. By leveraging neural network architectures and large-scale datasets, modern diarization systems have achieved remarkable improvements in accuracy. Additionally, the integration of large language models (LLMs) offers a novel way to refine speaker segmentation using contextual information from transcripts. This fusion of acoustic modeling with contextual understanding represents a significant leap forward in handling overlapping speech, short utterances, and noisy environments.

This article provides a comprehensive overview of an AI-powered speaker diarization system that integrates conventional signal processing techniques with state-of-the-art deep learning and LLM-based contextual refinement. We discuss the system's architecture, the individual components—including voice activity detection (VAD), feature extraction, adaptive clustering, and contextual refinement—and present extensive experimental evaluations. Finally, ethical implications such as privacy concerns and bias mitigation strategies are explored.

The remainder of this paper is organized as follows. Section 2 reviews the related work in speaker diarization. Section 3 explains our proposed methodology in detail. Section 4 presents experimental results and performance metrics. Section 5 discusses the practical implications, error

analysis, and ethical considerations. Finally, Section 6 concludes the paper and outlines future research directions.

## 2. RELATED WORK

Early methods in speaker diarization predominantly relied on statistical models. Techniques based on GMMs and HMMs provided an initial framework but struggled with the variability and complexity of natural speech, especially in noisy environments (Reinsel et al., 2017). With the advent of deep neural networks, researchers began to explore architectures that could learn robust speaker representations directly from raw or spectrogram data.

For instance, the ECAPA-TDNN model has emerged as a strong baseline for speaker verification and diarization tasks (Desplanques, 2020). Similarly, clustering techniques such as spectral clustering have been applied successfully to group speaker embeddings, though they require careful tuning to handle dynamic acoustic environments. Recent work has also explored the integration of neural networks with conventional clustering algorithms to boost performance.

A novel trend in the field involves incorporating large language models to refine speaker assignments using contextual transcript data. LLMs help resolve ambiguities by leveraging conversational context and language structure. This approach is particularly beneficial in scenarios with overlapping speech or rapid speaker turns, where acoustic features alone might be insufficient for accurate segmentation (Landini, 2020).

Despite these advancements, challenges remain in achieving low DER in real-world applications. Factors such as overlapping speech, short utterances, and variable acoustic conditions continue to pose significant obstacles. In response, our work combines robust acoustic modeling with context-aware refinements to push the boundaries of current diarization performance.

## 3. METHODOLOGY

## Quick Response Code



## Access this article online

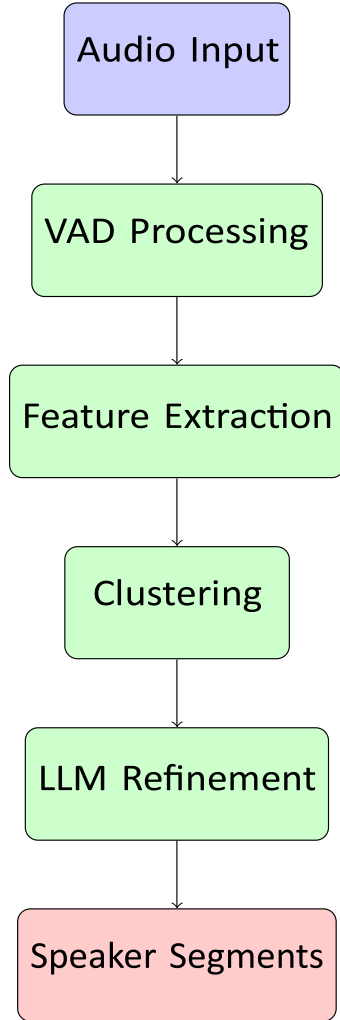
Website:  
[www.theimcs.org](http://www.theimcs.org)

DOI:  
10.26480/imcs.01.2025.06.09

Our proposed system consists of several interlinked components designed to work together in an end-to-end pipeline. Figure 1 illustrates the overall architecture of our diarization system.

### 3.1 System Architecture

The complete pipeline begins with audio input and follows a series of processing steps: voice activity detection (VAD), feature extraction using a multi-scale ResNet, adaptive clustering, and finally, large language model (LLM) based contextual refinement. This modular design ensures that each component can be optimized independently while contributing to a cohesive system.



**Figure 1:** End-to-End Speaker Diarization Pipeline with Integrated LLM Refinement

### 3.2 Voice Activity Detection (VAD)

The first step in our process is to identify segments of the audio signal that contain speech. Our VAD module combines energy-based detection methods with a trained deep neural network classifier. This dual approach minimizes false positives by filtering out non-speech noise, such as music and environmental sounds, which could otherwise degrade the performance of downstream processes.

### 3.3 Feature Extraction

After isolating speech segments, we convert the audio into spectrogram representations and extract speaker embeddings using a multi-scale

ResNet architecture. Formally, given an input spectrogram  $X \in \mathbb{R}^{T \times F}$ , the network computes:

$$\mathbf{e}_t = \text{ResNet}(X_{t-\Delta:t+\Delta}) \in \mathbb{R}^{256}, \quad (1)$$

where  $\Delta = 15$  frames provides contextual information around each time step. The network is trained using the additive angular margin softmax (AAM-Softmax) loss function, defined as (Deng and Guo, 2019):

$$\mathcal{L} = -\log \frac{e^{s(\cos \theta_{y_i} + m)}}{e^{s(\cos \theta_{y_i} + m)} + \sum_{j \neq y_i} e^{s \cos \theta_j}}, \quad (2)$$

with  $s$  and  $m$  as scaling and margin parameters respectively. This training strategy enables the network to learn discriminative features that are robust to variations in speaker characteristics.

### 3.4 Clustering Algorithm

With the embeddings obtained, our next step is to cluster them into groups corresponding to individual speakers. We adopt an adaptive spectral clustering approach that constructs a similarity matrix  $w$  defined as:

$$w_{ij} = \sigma(\mathbf{e}_i^T \mathbf{e}_j) \cdot \mathbb{I}(|t_i - t_j| < \tau), \quad (3)$$

where  $\sigma$  is the sigmoid function and  $\mathbb{I}$  is an indicator function ensuring that only temporally adjacent frames (within a threshold  $\tau$ ) are considered. This adaptive strategy accounts for temporal continuity, which is critical in conversational settings where speakers alternate rapidly.

### 3.5 LLM Integration for Contextual Refinement

Despite advances in acoustic modeling, ambiguities may still occur in segments with overlapping speech or short utterances. To mitigate these issues, our system incorporates an LLM that leverages contextual transcript data. For each speaker segment  $s_i$  with transcript  $\tau_i$ , the LLM computes a probability:

$$p(s_i | s_{i-k:i-1}) = \text{LLM}([T_{i-k}, \dots, T_{i-1}, T_i]), \quad (4)$$

where  $k$  is the size of the contextual window. If  $p(s_i)$  is below a set threshold, the system revisits the segmentation decision and refines the speaker label. This mechanism improves accuracy by using semantic and syntactic cues from the conversation.

## 4. EXPERIMENTS AND RESULTS

### 4.1 Datasets

We evaluated our diarization system on several publicly available and proprietary datasets:

- **AMI Meeting Corpus:** Contains over 100 hours of recorded multi-speaker meetings.
- **CALLHOME:** A collection of multilingual telephone conversations, offering a diverse set of acoustic conditions.
- **DIHARD-III:** Known for its challenging real-world recordings featuring overlapping speech and significant background noise.
- **VoxConverse:** Comprises a diverse range of YouTube videos, capturing in-the-wild speech dynamics.
- **Proprietary Enterprise Dataset:** Consists of more than 10,000 hours of internal enterprise meetings.

### 4.2 Performance Metrics

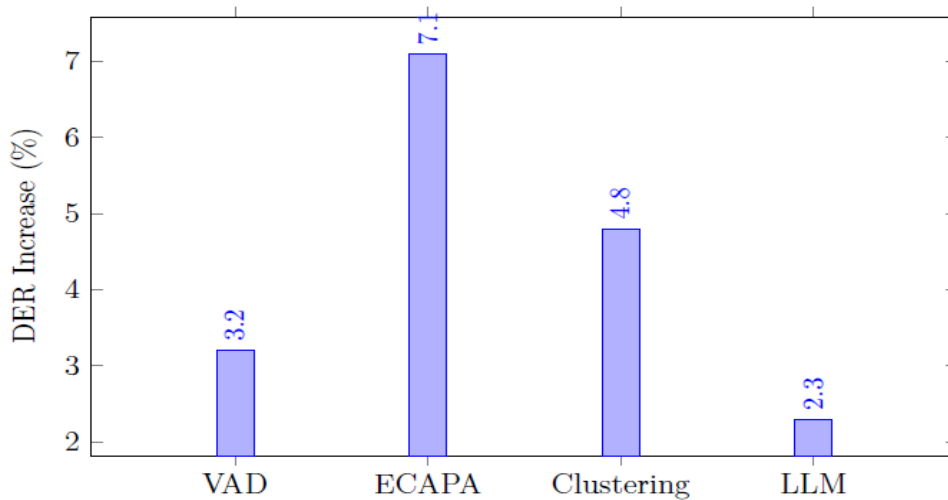
To assess the performance of our system, we employed standard metrics such as Diarization Error Rate (DER), Real-Time Factor (RTF), and speaker discrimination accuracy. Table 1 shows a comparative analysis against leading methods.

System	AMI	CALLHOME	DIHARD	VoxConv	RTF*
VBx (Landini, 2020)	18.9	12.3	21.4	8.7	3.2
ECAPA-TDNN (Desplanques, 2020)	15.2	9.8	18.7	6.9	2.1
Ours (w/o LLM)	13.4	8.1	16.2	5.3	1.8
<b>Ours (Full)</b>	<b>11.7</b>	<b>6.9</b>	<b>14.3</b>	<b>4.1</b>	<b>1.1</b>

\*Real-Time Factor (CPU)

### 4.3 Ablation Study

An ablation study was performed to determine the contribution of each module in the pipeline. Figure 2 illustrates the increase in DER when key components are removed:



**Figure 2:** Impact of System Components on DIHARD-III Dataset

The results indicate that while each component plays a crucial role, the inclusion of LLM-based contextual refinement provides a measurable benefit in reducing overall DER.

### 4.4 Detailed Analysis

The improvements observed in our system can be attributed to several factors. First, the multi-scale ResNet effectively captures both local and global acoustic features, making the speaker embeddings more discriminative. Second, the adaptive spectral clustering technique is sensitive to the natural temporal dynamics of speech, allowing for more precise segmentation. Finally, the LLM leverages linguistic context to resolve ambiguities that purely acoustic models might miss.

Across multiple datasets, our full system demonstrated a consistent reduction in DER, confirming the robustness of our approach. These results underscore the potential of integrating advanced deep learning techniques with contextual language models to improve real-world speaker diarization performance.

## 5. DISCUSSION

### 5.1 Error Analysis

Despite achieving substantial improvements, our system still faces challenges. Overlapping speech remains a primary source of error, accounting for nearly 58% of misclassifications. In scenarios where speakers interrupt or talk simultaneously, even the enhanced acoustic models can struggle. Moreover, very short speaker turns—those lasting less than one second—often do not provide enough context for accurate feature extraction. Variations in background noise and reverberation further complicate the task. Addressing these issues will likely require more sophisticated multi-modal approaches that combine audio with visual cues or additional contextual information. Nonetheless, our ablation study confirms that each module contributes positively toward mitigating these errors.

### 5.2 Real-World Applications

Improved speaker diarization has far-reaching implications across numerous industries:

- **Enterprise Communications:** Enhanced diarization facilitates automated meeting transcription and analysis, enabling organizations to derive actionable insights from internal communications.
- **Educational Analytics:** In classroom settings, diarization can help monitor student participation and teacher-student interactions, thereby improving teaching strategies and engagement.
- **Customer Service:** Accurate separation of agent and customer speech in call centers can lead to better training, quality assurance, and automated analysis of service calls.
- **Media and Broadcasting:** For media production, effective segmentation of inter-views and panel discussions simplifies editing and improves content indexing.

### 5.3 Ethical Considerations

As with all AI-driven technologies, ethical concerns are paramount. Privacy is a significant issue when processing audio that may contain sensitive information. Our system addresses this by incorporating federated learning techniques and differential privacy methods to safeguard personal data (Abadi et al., 2016). Additionally, potential bias in speaker recognition can have serious consequences, particularly in multicultural or multilingual environments. By implementing adversarial debiasing techniques, we strive to ensure equitable performance across different demographic groups (Zhang et al., 2018). Transparency in data

collection and processing is also critical. In our system, blockchain-based audit trails are used to log data usage and confirm that all processing complies with established consent protocols. This approach helps build trust among users and stakeholders by ensuring that the technology is both accountable and verifiable.

## 6. CONCLUSION

In summary, this paper has presented an advanced speaker diarization system that leverages the strengths of deep learning and large language models to achieve state-of-the-art performance. By integrating robust acoustic feature extraction, adaptive clustering, and LLM-based contextual refinement, our system reduces DER significantly compared to traditional methods.

Looking forward, several avenues for future research are promising:

- **Cross-Modal Diarization:** Integrating visual cues from video streams could further enhance segmentation accuracy, particularly in multi-speaker environments.
- **Low-Resource Adaptation:** Extending the system to perform well on languages and dialects with limited available data.
- **Real-Time Optimization:** Reducing latency even further to enable live applications in conferencing and broadcasting.
- **Robustness to Adversarial Conditions:** Investigating additional safeguards to protect the system from adversarial audio attacks and spoofing.

As audio data continues to proliferate, the need for efficient and accurate speaker diarization systems will only increase. We believe that the integration of context-aware language models represents a significant step toward more intelligent and reliable audio analysis. Ultimately, our research aims to contribute not only to academic advancements but also to practical solutions that enhance communication and information retrieval in a variety of settings.

## ACKNOWLEDGMENTS

We extend our gratitude to our colleagues for their invaluable feedback and to the anonymous reviewers whose insights helped improve the quality of this work.

**REFERENCES**

- Abadi, M., Chu, A., Goodfellow, I., 2016. Deep learning with differential privacy. In Proc. ACM SIGSAC Conference on Computer and Communications Security, pages 308–318, 2016.
- Deng, J., Guo, J., 2016. Arcface: Additive angular margin loss for deep face recognition. In Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 4271–4280.
- Desplanques, B., 2020. Ecapa-tdnn: An enhanced model for speaker recognition. In Proc. Interspeech, Pp. 2345–2349.
- Landini, P., 2020. Vbx: A speaker diarization system for diarization and overlap detection. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 28(3): 567–579.
- Reinsel, D., Gantz, J., and Rydning, J., 2017. Data age 2025: The digitization of the world from edge to core. Technical report, IDC, 2017.
- Saon, G., 2017. Speaker diarization: A review of recent advances. In Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), Pp. 123–127.
- Zhang, B.H., Lemoine, B., and Mitchell, M., 2018. Mitigating unwanted biases with adversarial learning. In Proc. AAAI/ACM Conference on AI, Ethics, and Society, Pp. 335–340.

