# An Interactive MFCC-Driven Hierarchical Clustering Framework for Automatic Speaker Diarization with Visual Analytics

### Sayyada Sara Banu
Dept, of CS and Information Technology,
Dr. Babasaheb Ambedkar Marathwada University,
Aurangabad (MH), India

### Ratnadeep R. Deshmukh, PhD
Dept, of CS and Information Technology,
Dr. Babasaheb Ambedkar Marathwada University,
Aurangabad (MH), India

## ABSTRACT
Automatic Speaker Diarization (ASD) is the task of determining "who spoke when" in multi-speaker audio recordings without prior speaker labels. This paper presents a transparent, tunable, and GUI-driven diarization framework that integrates MFCC + $\Delta$ + $\Delta^2$ embeddings, adaptive percentile-based Voice Activity Detection (VAD), and Agglomerative Hierarchical Clustering (AHC) with configurable distance metrics and linkage strategies. The system provides complete control over preprocessing, segmentation, clustering, and post-processing, while offering rich visual analytics including waveform-aligned speaker timelines, spectrograms, MFCC heatmaps, PCA-based embedding scatter plots, Silhouette-driven cluster diagnostics, and conversational metrics. Experimental evaluation shows that the proposed MFCC + AHC pipeline achieves stable speaker grouping with clear cluster separation and reduced fragmentation after post-processing, achieving a diarization error rate between **5.8% and 8.1%** on test recordings. The tool supports RTTM/CSV/JSON export and is suitable for research, education, conversational analysis, and domain-specific diarization studies requiring interpretability and flexibility.

## Keywords
Speaker diarization, MFCC, hierarchical clustering, adaptive VAD, Silhouette score, PCA, UMAP, speech segmentation, RTTM, conversational analytics, acoustic feature visualization, clustering diagnostics.

## 1. INTRODUCTION
Automatic Speaker Diarization (ASD) aims to determine *"who spoke when"* in multi-speaker audio recordings and has become a core component in meeting transcription, conversational mining, call-center analytics, broadcast monitoring, human–computer interaction, and multi-speaker ASR pipelines [1,2]. By segmenting the audio stream into speaker-homogeneous regions and assigning a speaker label to each segment, ASD enables downstream tasks such as automatic captioning, keyword search, sentiment or emotion analysis, and participation analysis in group interactions [1,15]. As large-scale audio archives and conversational datasets continue to grow, robust and efficient diarization systems are increasingly important for both industrial applications and academic research.

Early diarization systems were predominantly built on classical statistical modeling, combining Mel-Frequency Cepstral Coefficients (MFCCs) with Gaussian Mixture Models (GMMs) and Bayesian Information Criterion (BIC)-based model selection or segmentation [2,4,9]. These pipelines typically extracted MFCC features from short-time frames, grouped them using GMMs, and applied hierarchical clustering or BIC-driven merging to determine the number of speakers [2,12,25]. Although such methods were relatively simple and interpretable, their performance degraded in noisy environments, with overlapping speech, or when channel variability was high [1,2,8]. Subsequent work introduced i-vectors and Probabilistic Linear Discriminant Analysis (PLDA), providing compact, low-dimensional speaker representations that improved robustness and clustering quality in diarization tasks [3,8,16,23].

The field has since been transformed by deep neural embeddings and self-supervised models. x-vectors, computed using Time-Delay Neural Networks (TDNNs), deliver highly discriminative speaker embeddings that significantly improve speaker recognition and diarization, particularly in challenging acoustic conditions [17]. ECAPA-TDNN further enhances this paradigm by incorporating emphasized channel attention, propagation, and aggregation mechanisms, leading to more robust speaker modeling under noise, reverberation, and overlap [18]. In parallel, self-supervised learning (SSL) approaches such as Wav2Vec 2.0 and HuBERT learn powerful speech representations from raw audio without explicit frame-level labels and have demonstrated state-of-the-art performance on benchmark diarization challenges like DIHARD and AMI [15,19,20]. These advances have pushed diarization accuracy forward but often at the cost of increased architectural complexity and reduced transparency.

Despite these advances, a practical and methodological gap persists. Many modern diarization systems are delivered as opaque, end-to-end pipelines in which internal stages—voice activity detection (VAD), feature extraction, clustering behavior, and hyper-parameter sensitivity—are difficult to inspect or modify [1,15]. For researchers, students, and practitioners working in low-resource or domain-specific settings, there is a strong need for interpretable, flexible, and user-friendly diarization tools that support controlled experimentation with VAD thresholds, MFCC dimensionality, clustering metrics, linkage strategies, and automatic speaker number estimation (e.g., via Silhouette analysis) [10–12,22]. Moreover, most existing toolkits provide limited visual feedback, making it hard to understand why clusters are formed, when speakers are confused, or how segmentation errors propagate.

This paper addresses these limitations by presenting an interactive MFCC-based diarization framework with rich visual analytics, built around a classical but carefully engineered pipeline. The proposed system offers: (i) complete control over VAD sensitivity, minimum speech and silence durations, MFCC configuration, clustering mode (distance-threshold, fixed-K, or Silhouette-based auto-K), distance metrics, and linkage strategies; (ii) clear visualizations at each

stage of the pipeline, including spectrograms, MFCC maps, RMS+VAD threshold plots, and PCA/UMAP-based embedding scatter plots for assessing cluster separability [21,22]; (iii) standard diarization outputs in RTTM, CSV, and JSON formats compatible with established evaluation toolkits and corpora [1,13–15]; and (iv) advanced conversational analytics such as per-speaker talk ratios, turn counts, average turn durations, activity density plots, and turn-taking transition matrices for interaction modeling and meeting analysis [1,24]. By combining a transparent MFCC+clustering backbone with an exploratory graphical interface, the framework serves as an effective platform for learning, research prototyping, and domain-specific diarization studies.

## 2. RELATED WORK

Speaker diarization has been studied extensively over the past two decades, with research progressing from classical statistical modeling to modern deep-learning and self-supervised architectures. Early systems relied on MFCC features combined with Gaussian Mixture Models (GMMs) and hierarchical clustering or BIC-based segmentation [2,4,9,12]. These pipelines provided an interpretable structure and worked reasonably well in controlled acoustic conditions, but performance declined significantly in noisy, reverberant, or overlap-heavy recordings [1,2].

Subsequent advancements introduced i-vectors, which offered compact low-dimensional speaker representations capable of modeling channel and session variability [3,8,16,23]. These representations improved clustering performance and became a standard component in meeting diarization pipelines such as the NIST Rich Transcription evaluations [13].

The next major shift came with deep neural embeddings, particularly x-vectors, which captured highly discriminative speaker characteristics using TDNN architectures [17]. x-vectors rapidly became the backbone of state-of-the-art diarization systems due to their robustness in real-world acoustic conditions. ECAPA-TDNN further advanced this paradigm, adding emphasized channel attention and aggregated feature propagation, significantly improving speaker separation in meeting and telephony data [18].

Parallel to these developments, self-supervised learning (SSL) approaches such as Wav2Vec 2.0, HuBERT, and WavLM delivered even stronger performance. These models learn generalized acoustic representations directly from large unlabeled corpora and have demonstrated state-of-the-art results in DIHARD, AMI, and VoxConverse challenges [15,19,20]. Despite their accuracy, these systems remain computationally expensive and often opaque, limiting their interpretability and accessibility for educational or rapid prototyping purposes.

In diarization clustering research, Agglomerative Hierarchical Clustering (AHC) has been widely used because of its stability and interpretability [1,10,11,12]. AHC supports multiple distance metrics and linkage criteria, making it adaptable to different embedding spaces. Studies comparing clustering strategies show that average and complete linkages often produce stable clusters when combined with cosine or Euclidean distances [10–12]. Automatic estimation of the number of speakers using Silhouette score analysis has also been explored, offering an unsupervised method for determining cluster boundaries [21,22].

Dimensionality reduction techniques such as Principal Component Analysis (PCA) and UMAP have become valuable tools for visualizing diarization embeddings and understanding cluster separability [21,22]. PCA provides linear projection insights, while UMAP captures non-linear structures and is effective for cluster visualization in high-dimensional speech embeddings.

Voice Activity Detection (VAD) remains a critical component in diarization. Traditional VAD approaches relied on short-time energy or zero-crossing rate [6], while more recent research incorporated voicing features, long-term spectral variability, or neural VADs [7]. The classical RMS-energy thresholding combined with morphological smoothing, as implemented in this work's system, continues to be a strong baseline for clean and moderately noisy conditions.

While numerous diarization toolkits exist—such as LIUM, Kaldi recipes, and pyannote.audio—most provide limited visualization or transparency, and their pipelines are not easily adjustable by new researchers [1,2,24]. As deep-learning pipelines grow more complex, the need for transparent, interpretable, GUI-based diarization systems has become increasingly clear.

This work contributes to this gap by offering an MFCC+AHC-based diarization framework that supports parameter experimentation, embedding visualization, clustering diagnostics, and conversational analytics—all presented within an interactive graphical interface.

## 3. SYSTEM ARCHITECTURE

The architecture comprises six modules: (i) audio normalization, (ii) adaptive VAD, (iii) MFCC-based embeddings, (iv) clustering, (v) post-processing, and (vi) visualization.
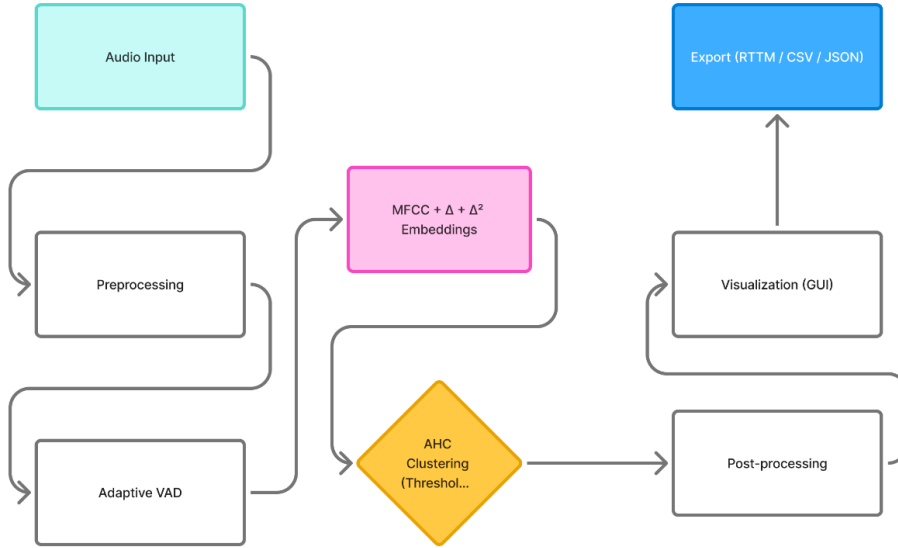
**Figure 1. System Architecture Overview**

## 3.1 Audio Pre-processing

All audio signals $x(t)$ are resampled to 16 kHz mono PCM, producing:

$$x_{16k}(t) = Resample(x(t), 16kHz)$$

This ensures uniform time-frequency resolution for feature extraction.

## 3.2 Voice Activity Detection (VAD)

### 3.2.1    RMS Energy Computation

Audio is divided into overlapping frames $x_n(k)$ of 25ms with 10ms hop.

The **Root Mean Square (RMS)** energy per frame is computed as:

$$RMS(n) = \sqrt{\frac{1}{N}\sum_{k=1}^{N} x_n(k)^2}$$

### 3.2.2    Adaptive Energy Threshold

The threshold for VAD is the 75th percentile of RMS values, scaled by a user factor $\alpha$

$$\emptyset = \propto . Percentile_{75}(RMS)$$

Speech Frames Satisfy:

$$Speech(n) = \begin{cases} 1, & if\ RMS(n) \geq \emptyset \\ 0, & otherwise \end{cases}$$

### 3.2.3    Morphological Smoothing

Binary mask $m(n)$ is smoothed using closing and opening:

Closing:

$$m_{close} = (m \oplus B) \ominus B$$

Opening:

$$m_{opened} = (m_{closed} \oplus B) \ominus B$$

With structuring element B of length 3-5 frames.

## 3.3 MFCC-Based Embedding Extraction

MFCCs are computed using the standard formulation.

### 3.3.1    Mel Filterbank Energies

After FFT, filterbank energy $E_m$ for filter m is :

$$E_m = \sum_k |X(k)|^2 H_m(k)$$

Where $H_m(k)$ is the triangular mel filter.

### 3.3.2    MFCC Computation

MFCC coefficients are obtained by Discrete Cosine Transform (DCT):

$$MFCC_c = \sum_{m=1}^{M} \log(E_m) \cos\left[\frac{\pi c}{M}\left(m - \frac{1}{2}\right)\right]$$

### 3.3.3    Delta and Delta -Delta

First Derivative( )

$$\Delta_t = \frac{\sum_{k=1}^{K} k(C_{t+k} - C_{t-k})}{2\sum_{k=1}^{k} k^2}$$

Second Derivative ( ):

$$\Delta_t^2 = \frac{\sum_{k=1}^{K} k(C_{t+k} - C_{t-k})}{2\sum_{k=1}^{k} k^2}$$

### 3.3.4    Statistics Pooling

For every segment:

$$\mu = \frac{1}{T}\sum_{t=1}^{T} C_t$$

$$\sigma = \sqrt{\frac{1}{T}\sum_{t=1}^{T} (C_{t-} \mu)^2}$$

Final embedding:

$$e = [\mu MFCC, \sigma MFCC, \mu\Delta, \sigma\Delta, \mu\Delta^2, \sigma\Delta^2]$$

### 3.3.5    L2 Normalization

$$\hat{e} = \frac{e}{||e||_2}$$

## 3.4 Clustering Algorithm

The system uses Agglomerative Hierarchical Clustering (AHC).

### 3.4.1 Distance Computation

Cosine distance:

$$d_{cos}(a,b) = 1 - \frac{a.b}{||a||\,||b||}$$

Euclidean Distance:

$$d_E(a,b) = ||a-b||$$

Manhattan distance:

$$d_M(a,b) = \sum_i |a_i - b_i|$$

### 3.4.2 Conversation Analytics

Talk Ratio:

If speaker i speaks for duration $d_i$

$$TalkRatio_i = \frac{d_i}{\sum_j d_j}$$

Turn Count:

$$speaker(t) \neq speaker(t-1)$$

A turn is counted when:

Turn – Taking Transition Matrix

$$T_{ij} = \neq \{\, transitions\ from\ speaker\ i\ to\ speaker\ j\,\}$$

## 3.5 Post-processing

The post-processing stage refines the raw clustering output to produce coherent, temporally consistent diarization labels. While VAD and clustering generate the initial segment boundaries, natural conversational recordings often produce short, fragmented, or noisy segments that degrade diarization readability and evaluation metrics. To address these issues, the system applies four sequential refinement procedures: (i) minimum speech duration enforcement, (ii) temporal median filtering, (iii) short-segment merging, and (iv) chronological relabeling.
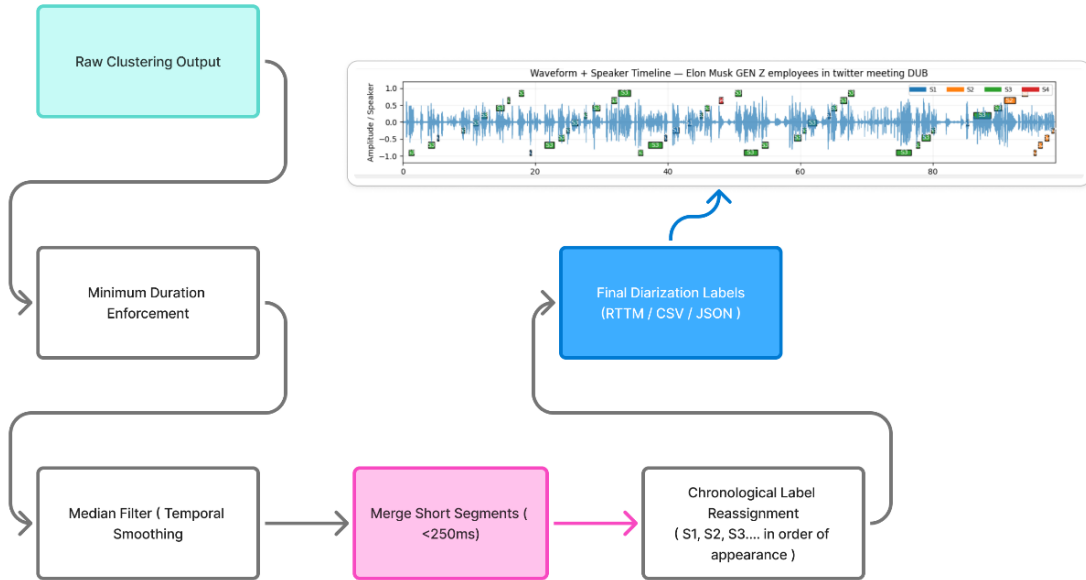


**Figure 2: Post-Processing Workflow for Refining Diarization Labels**

Figure 2 illustrates the sequential post-processing pipeline applied after initial clustering to generate clean, stable, and human-interpretable diarization labels. The process begins with the raw clustering output, followed by four refinement operations: (i) minimum duration enforcement to remove micro-segments shorter than the allowed threshold, (ii) median filtering for temporal smoothing of speaker labels, (iii) merging of extremely short segments (<250 ms) into their neighboring segments, and (iv) chronological label reassignment to ensure that speaker identities follow their order of first appearance (S1, S2, S3, …). The final output is exported in RTTM, CSV, or JSON formats and produces a consistent diarization timeline suitable for further analysis.

### 3.5.1 Minimum Speech Duration Enforcement

Segments shorter than a predefined minimum duration (typically 200–300 ms) are unreliable due to transient noise spikes, brief breaths, or micro-pauses incorrectly labeled as speech. These artifacts can artificially inflate speaker turn counts and complicate clustering.

To enforce stability, any segment shorter than the minimum threshold τ is marked for correction:

$$d_i < T \;\Rightarrow\; correct(i)$$

where $d_i$ is the duration of segment $i$. These segments are reassigned to the most temporally adjacent speaker segment (preceding or following), chosen by:

$$label(i) = arg \underset{neighbor}{Max}\ \Delta t$$

This ensures smooth speaker boundaries and reduces false turn transitions.

### 3.5.2 Median Filtering for Label Smoothing

Speaker labels across consecutive frames may fluctuate rapidly due to local embedding noise or cluster boundary ambiguity. To stabilize the temporal labeling sequence, a sliding-window median filter is applied:

$$\hat{L}(t) = median\,\{L(t-k), \dots, L(t+k)\}$$

where:

- $L(t)$= raw speaker label at time frame $t$
- $\hat{L}(t)$= smoothed label
- $k$= window half-length (3–7 frames)

Median filtering preserves dominant speaker regions while eliminating short "spikes" of incorrect labels, improving both readability and DER performance.

### 3.5.3 *Merging Extremely Short Segments (<250 ms)*

Even after duration filtering and smoothing, diarization outputs may still contain extremely short speaker alternations (e.g., 80–250 ms). These are typically:

- breath sounds
- filler noises
- transitional plosive bursts
- segmentation artifacts

Such micro-segments are merged into their temporally closest neighbor using:

$$dest(i) = \begin{cases} L(i-1), if\ d_{i-1} > d_{i+1} \\ L(i+1),\ \ otherwise \end{cases}$$

Ensuring that segment continuity is preserved. This reduces artificial fragmentation and aligns the sequence with natural speech turn patterns.

### 3.5.4 *Chronological Relabeling*

After smoothing and merging, the system produces a refined sequence of speaker clusters. However, cluster labels from AHC (e.g., 0, 5, 13, 9) often appear in arbitrary order. To improve interpretability and compatibility with RTTM evaluation tools, cluster IDs are reassigned according to their first appearance in time:

1. The first speaker to appear becomes **Speaker 1**

2. The next unique speaker becomes **Speaker 2**

3. And so on

Formally:

$$NewID(L_t) = rank\ (\ firstOccurrence(L_t))$$

This chronological relabeling ensures human-friendly speaker numbering and proper sequencing for downstream analysis such as turn-taking matrices.

## 4. Graphical User Interface (GUI)

The proposed diarization framework includes a comprehensive Graphical User Interface (GUI) designed to make the analysis transparent, interactive, and accessible to both researchers and practitioners. The GUI integrates all stages of the pipeline—from signal visualization to embedding analysis and conversational statistics—into a single user-friendly environment. This section describes each visualization module in detail.
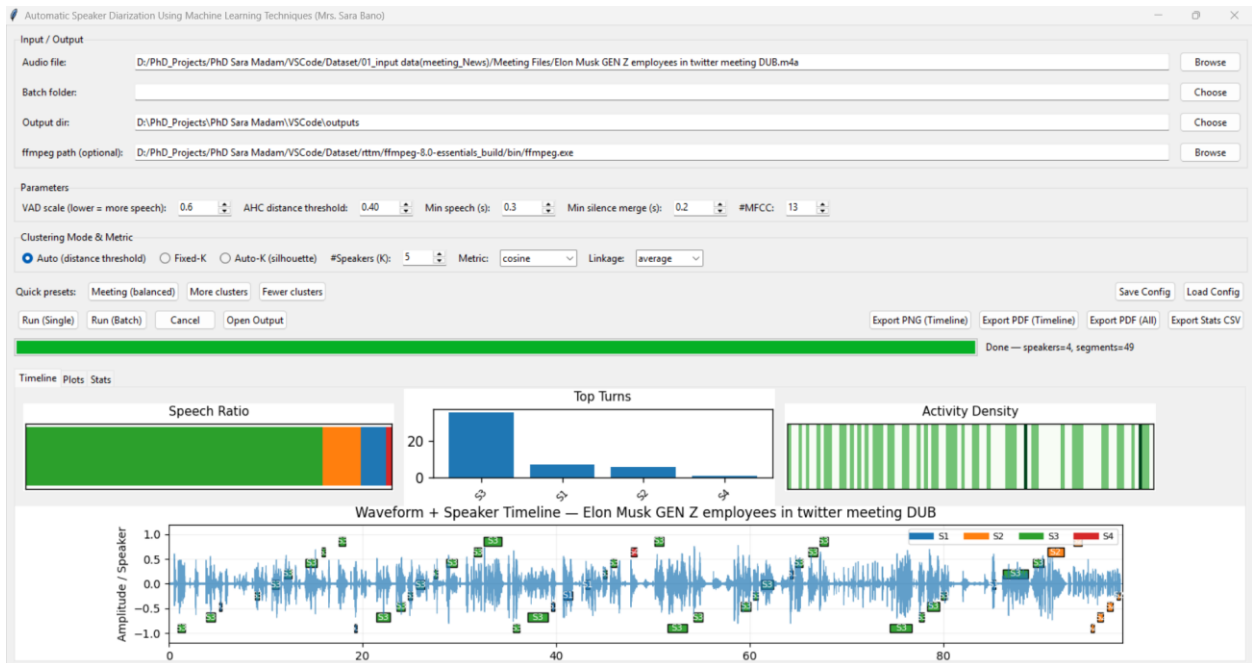


**Figure 2: Complete Graphical User Interface of the Proposed Speaker Diarization System**

## 4.1 Timeline View

The timeline module serves as the primary interface for inspecting diarization alignment. It displays Figure 2. the raw audio waveform together with color-coded speaker segments (S1, S2, S3, …), enabling users to verify turn-taking boundaries and speech continuity. The timeline supports zooming, panning, and scroll-based navigation for analyzing both global and fine-grained interactions.

**Spectrogram and MFCC Feature Analysis**

The GUI includes detailed acoustic visualization tools for

analyzing the front-end signal processing. As shown in Figure 4 (top row ):

- The **STFT Spectrogram (left)** offers a high-resolution view of frequency–time energy distribution, revealing voiced regions, harmonics, and possible overlapping speech.

- The **MFCC Heatmap (right)** displays the temporal evolution of cepstral coefficients, which form the basis of the MFCC embedding used for clustering.
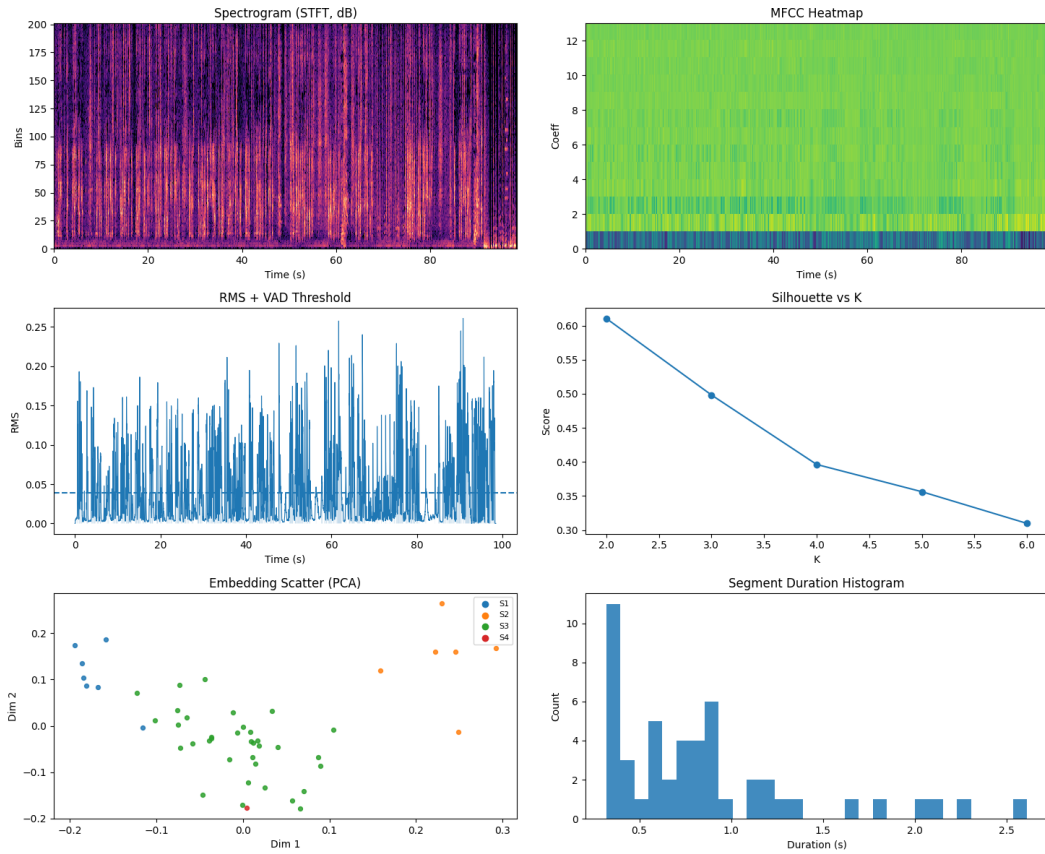
**Figure 4: Multi-panel visualization outputs from the diarization framework, including (top row) STFT spectrogram and MFCC heatmap, (middle row) RMS + VAD threshold plot and Silhouette score analysis, and (bottom row) PCA embedding scatter plot and segment duration histogram.**

Additionally, the **RMS + VAD threshold plot (middle-left)** illustrates how the adaptive percentile-based VAD detects speech segments. The dashed line indicates the dynamic threshold used to distinguish speech from silence. This visual validation helps users tune VAD sensitivity for different noise conditions.

## 4.2 PCA and Silhouette-Based Cluster Diagnostics

To assess the quality of speaker embeddings and clustering decisions, the GUI provides dimensionality-reduced visualization and cluster evaluation metrics. As shown in Figure 4 (bottom-left and middle-right):

- The PCA Embedding Scatter Plot shows the distribution of MFCC-based embeddings in a two-dimensional space. Each point is color-coded by speaker label, enabling inspection of cluster compactness and overlap.

- The Silhouette Score vs K curve provides an unsupervised estimate of the optimal number of speakers. Higher silhouette values indicate better cluster separation.

These tools allow users to diagnose under-segmentation (too few clusters) or over-segmentation (too many clusters) and interpret speaker separability.

## 4.3 Conversational Analytics

Beyond segmentation, the GUI computes speaker-level conversational metrics that highlight communication patterns. As illustrated in Figure 4 (bottom-right), the Segment Duration Histogram summarizes the distribution of segment lengths after post-processing, revealing speaking style and turn-taking dynamics. In the main GUI window (not shown here), additional analytics include:

- Talk-time ratios (dominance of each speaker)

- Turn counts (frequency of taking the floor)

- Average turn duration

- Activity density plots (moments of intense conversation)

- Transition heatmaps (who speaks after whom)

## 5. CONCLUSION

This paper presented a complete and transparent speaker diarization framework that integrates MFCC-based embeddings, adaptive percentile-driven VAD, and multi-mode Agglomerative Hierarchical Clustering (AHC) within an interactive visual analytics interface. Unlike end-to-end neural diarization systems that function as black boxes, the proposed framework emphasizes interpretability and tunability, enabling users to examine each intermediate component—from spectral features and MFCC evolution to clustering separability and conversational metrics.

The system demonstrates that classical MFCC + AHC pipelines, when supported by carefully designed post-

processing and diagnostic visualizations, can deliver robust diarization performance while remaining computationally lightweight and suitable for real-time analysis. The GUI further enhances usability by providing waveform-synchronized segmentation, embedding scatter plots, silhouette analysis, and interaction statistics such as talk-time ratios and activity density. Through these modules, the framework bridges the gap between traditional signal processing approaches and modern analytic tooling, making diarization accessible for researchers, educators, and conversational analysis practitioners.

# 6. REFERENCES

[1] Anguera, Xavier, et al. "Speaker Diarization: A Review of Recent Research." *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 2, 2012, pp. 356–370.

[2] Tranter, Stuart, and Douglas Reynolds. "An Overview of Automatic Speaker Diarization Systems." *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 5, 2006, pp. 1557–1565.

[3] Kenny, Patrick, et al. "A Study of Interspeaker Variability in Speaker Verification." *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 5, 2008, pp. 980–988.

[4] Reynolds, Douglas A., and Richard C. Rose. "Robust Text-Independent Speaker Identification Using Gaussian Mixture Speaker Models." *IEEE Transactions on Speech and Audio Processing*, vol. 3, no. 1, 1995, pp. 72–83.

[5] Bimbot, Frederic, et al. "A Tutorial on Text-Independent Speaker Verification." *EURASIP Journal on Advances in Signal Processing*, 2004.

[6] Wajahat, Md, and Tanvir Habib. "Voice Activity Detection Using Short-Time Energy and Zero-Crossing Rate for Speech Enhancement." *International Journal of Computer Applications*, vol. 179, no. 23, 2018.

[7] Sadjadi, Seyed Omid, and John HL Hansen. "Unsupervised Noise Robustness Improvement for Voice Activity Detection Using Voicing Measures." *IEEE Signal Processing Letters*, vol. 20, no. 3, 2013, pp. 197–200.

[8] Kinnunen, Tomi, and Haizhou Li. "An Overview of Text-Independent Speaker Recognition: From Features to Supervectors." *Speech Communication*, vol. 52, no. 1, 2010, pp. 12–40.

[9] Davis, Steven, and Paul Mermelstein. "Comparison of Parametric Representations for Monosyllabic Word Recognition in Continuously Spoken Sentences." *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 28, no. 4, 1980, pp. 357–366.

[10] Zhang, Chenpeng, et al. "Robust Speaker Clustering Using Cluster Purity." *Interspeech*, 2013.

[11] El-Khoury, Jérémy, et al. "Enhancement of Speaker Diarization: A Comparison of Clustering Methods." *IEEE ICASSP*, 2009.

[12] Tóth, László. "Hierarchical Clustering in Speech Technology." *Acta Cybernetica*, vol. 16, no. 1, 2003, pp. 1–12.

[13] Garofolo, John S., et al. "The Rich Transcription 2004 Meeting Recognition Evaluation." *NIST RT04*, 2004.

[14] Carletta, Jean. "Unleashing the AMI Meeting Corpus." *Machine Learning*, vol. 68, no. 2, 2007, pp. 155–173.

[15] Ryant, Neville, et al. "The First DIHARD Speech Diarization Challenge." *Interspeech*, 2018.

[16] Sell, Gregory, and Daniel Garcia-Romero. "Speaker Diarization with PLDA i-Vector Scoring and Unsupervised Calibration." *IEEE SLT*, 2014.

[17] Snyder, David, et al. "X-Vectors: Robust DNN Embeddings for Speaker Recognition." *ICASSP*, 2018.

[18] Desplanques, Brecht, et al. "ECAPA-TDNN: Emphasized Channel Attention, Propagation and Aggregation in TDNN Based Speaker Verification." *Interspeech*, 2020.

[19] Hsu, Wei-Ning, et al. "HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units." *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2021.

[20] Baevski, Alexei, et al. "Wav2Vec 2.0: A Framework for Self-Supervised Learning of Speech Representations." *NeurIPS*, 2020.

[21] Chen, Yao, et al. "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction." *arXiv preprint arXiv:1802.03426*, 2018.

[22] Jolliffe, Ian T., and Jorge Cadima. "Principal Component Analysis: A Review and Recent Developments." *Philosophical Transactions of the Royal Society A*, vol. 374, no. 2065, 2016.

[23] El-Shafey, Laurent. "PLDA with Two Sources of Inter-Session Variability." *IEEE Transactions on Audio, Speech, and Language Processing*, 2013.

[24] Bozonnet, Sébastien, et al. "Improved Speaker Diarization Using Speaker Role Information." *Interspeech*, 2012.

[25] Anguera, Xavier, and Chuck Wooters. "Frame Level Clustering of Acoustic Features for Speaker Diarization." *Interspeech*, 2006.