

# Integrating Audio, Visual, and Semantic Information for Enhanced Multimodal Speaker Diarization on Multi-party Conversation

Anonymous ACL submission

## Abstract

Speaker diarization aims to segment an audio stream into homogeneous partitions based on speaker identity, playing a crucial role in speech comprehension and analysis. Mainstream speaker diarization systems rely only on acoustic information, making the task particularly challenging in complex acoustic environments in real-world applications. Recently, significant efforts have been devoted to audio-visual or audio-semantic multimodal modeling to enhance speaker diarization performance; however, these approaches still struggle to address the complexities of speaker diarization on spontaneous and unstructured multi-party conversations. To fully exploit meaningful dialogue patterns, we propose a novel multimodal approach that jointly utilizes **audio**, **visual**, and **semantic** cues to enhance speaker diarization. Our approach structures visual cues among active speakers and semantic cues in spoken content into a cohesive format known as *pairwise constraints*, and employs a semi-supervised clustering technique based on **pairwise constrained propagation**. Extensive experiments conducted on multiple multimodal datasets demonstrate that our approach effectively integrates audio-visual-semantic information into the clustering process for acoustic speaker embeddings and consistently outperforms state-of-the-art speaker diarization methods, while largely preserving the overall system framework.

## 1 Introduction

Speaker diarization (SD) is the task of answering the question “who spoke when” by partitioning an audio stream into segments with timestamps and corresponding speaker labels. Speaker diarization is a crucial task in multi-party conversation scenarios, as it is important for speech comprehension and analysis to conduct automatic speech recognition (ASR) and also assign speaker labels to segments of audio or transcribed text. Many downstream

natural language processing (NLP) tasks (Ganesh et al., 2023a; Shen et al., 2023; Le et al., 2019; Ganesh et al., 2023b) have been proven to benefit from speaker diarization results.

Traditional speaker diarization systems rely solely on acoustic information and they can be generally categorized into two types: clustering-based approaches and end-to-end (E2E) approaches. Clustering-based approaches typically comprise three stages: voice activity detection (VAD) to filter out non-speech frames, speaker embedding extractor to obtain acoustic embeddings from each short speech segment, and an unsupervised speaker clustering to assign these embeddings into speaker classes (Anguera et al., 2012; Park et al., 2022). E2E approaches treat speaker diarization as a sequence labeling task, tagging each speech frame with its speaker identity, known as End-to-end neural diarization (EEND) (Fujita et al., 2020, 2019b). Although this modeling approach can unify the modeling of silence, single speaker speech, and speaker overlap, the absence of clustering often leads to a significant performance degradation in multi-party meeting scenarios with an uncertain number of participants, particularly when there are more than 3 speakers. The most popular acoustic-only speaker diarization systems are often rely on a clustering-based approach to determine the overall speaker results, while utilizing EEND as a sub-module to handle speaker changes and overlaps, such as Pyannotate (Bredin, 2023) and DiariZen (Han et al., 2024). Acoustic-only speaker diarization approaches often suffer significant performance degradation in challenging acoustic environments characterized by noise, reverberation, and speech overlapping between multiple speakers (Park et al., 2022). Recent studies have aimed to address this challenge by incorporating information from other modalities into the speaker diarization task. For instance, some works (Xu et al., 2022; Chung et al., 2020; Gebru et al., 2017) have inte-

085 grated visual cues, such as facial features and lip  
086 movement, with audio to determine active speak-  
087 ers. Other studies (Flemotomos and Narayanan,  
088 2022; Park and Georgiou, 2018; Zuluaga-Gomez  
089 et al., 2022) have utilized text data from automatic  
090 speech recognition (ASR) to identify speaker iden-  
091 tity and detect speaker change points. Although  
092 combining acoustic information with a single ad-  
093 ditional modality has shown some benefits, there  
094 is currently no effective approach to integrate in-  
095 formation from all three modalities—audio, visual,  
096 and textual—into the speaker diarization task.

097 In this paper, we propose a novel framework  
098 based on clustering-based speaker diarization, ca-  
099 pable of simultaneously modeling speaker-related  
100 information from multiple modalities. Specifi-  
101 cally, we incorporate visual information (e.g., face-  
102 tracking and lip movement) and textual information  
103 (e.g., dialogue and speaker-turn detection). These  
104 multimodal insights are integrated into *pairwise*  
105 *constraints* to enhance speaker clustering by replac-  
106 ing unsupervised clustering with a semi-supervised  
107 approach. This allows for effective multimodal  
108 fusion during the clustering stage. Our method  
109 is not limited by the absence of comprehensive  
110 multimodal datasets and maintains the structural  
111 integrity of traditional acoustic-only frameworks  
112 while benefiting from advancements in individual  
113 unimodal components. Experiments across multi-  
114 ple multimodal datasets have consistently demon-  
115 strated the effectiveness of our approach.

116 Our contributions can be summarized as follows:

- 117 • **We present a novel framework for speaker**  
118 **diarization, uniquely integrating audio, vi-**  
119 **sual, and semantic information.** This is  
120 the first framework to leverage all these three  
121 modalities, enhancing the robustness and ac-  
122 curacy of speaker diarization.
- 123 • **We introduce a joint pairwise constraint**  
124 **propagation method into the speaker clus-**  
125 **tering process,** effectively enhancing speaker  
126 clustering performance through multimodal  
127 information-derived constraints.
- 128 • To comprehensively evaluate the effectiveness  
129 of our method, **we contribute a 6.3-hour**  
130 **video evaluation set sourced from in-the-**  
131 **wild scenarios,** which has been annotated  
132 with speaker identity labels, corresponding  
133 speech activity timestamps, and speech con-  
134 tent.

## 2 Related Work 135

### 2.1 Multimodal Speaker Diarization 136

**Acoustic-Only speaker diarization** Audio-only  
137 speaker diarization has been studied exten-  
138 sively (Park et al., 2022). A typical speaker di-  
139 arization systems employ a multi-stage framework,  
140 including VAD (Gelly and Gauvain, 2018), speech  
141 segmentation (Xia et al., 2022), acoustic embed-  
142 ding extraction (Snyder et al., 2018; Zheng et al.,  
143 2020; Chen et al., 2023) and unsupervised cluster-  
144 ing such as agglomerative hierarchical clustering  
145 (AHC) (Day and Edelsbrunner, 1984) and spec-  
146 tral clustering(SC) (Wang et al., 2018). Recently,  
147 EEND where individual sub-modules in traditional  
148 systems can be replaced by one neural network  
149 has received more attention (Fujita et al., 2019a,c;  
150 Horiguchi et al., 2020) which treat speaker diari-  
151 zation as a frame-level sequence labeling task. Due  
152 to the absence of a clustering algorithm, the EEND  
153 method often experiences significant performance  
154 degradation in scenarios with a large number of  
155 speakers. Some approaches improve model per-  
156 formance by combining the global speaker pre-  
157 dictions from clustering with the local speaker  
158 change and overlap detection results from EEND,  
159 such as EEND-VC (Kinoshita et al., 2021) and Di-  
160 arizen (Han et al., 2024). Similar strategies have  
161 also been adopted by mainstream speaker diari-  
162 zation toolkits like Pyannote (Bredin, 2023). 163

**Audio-visual Speaker Diarization** Facial ac-  
164 tivities and lip motion are highly related to  
165 speech (Yehia et al., 1998). Visual information con-  
166 tains a strong clue for the identification of speak-  
167 ers and the location of speaker changes (Yoshioka  
168 et al., 2019), which can be used to significantly  
169 improve the accuracy of speaker diarization. Some  
170 methods leverage the audio and visual cues for di-  
171 arization using synchronization between talking  
172 faces and voice tracks (Chung et al., 2019). Other  
173 works (Xu et al., 2022; Wuerkaixi et al., 2022; Yin  
174 et al., 2024) utilized an attention-based network to  
175 perform middle-fusion and extract a unified rep-  
176 resentation of the two modalities. Recently, an  
177 interesting and promising direction is to use sep-  
178 arate neural networks to process data streams of  
179 two modalities and directly output speech probabil-  
180 ities for all speakers simultaneously (kui He et al.,  
181 2022), similar to audio-only EEND frameworks.  
182 All of these require expensive amounts of anno-  
183 tated audio-visual parallel data for training, which  
184 is expensive to acquire. 185

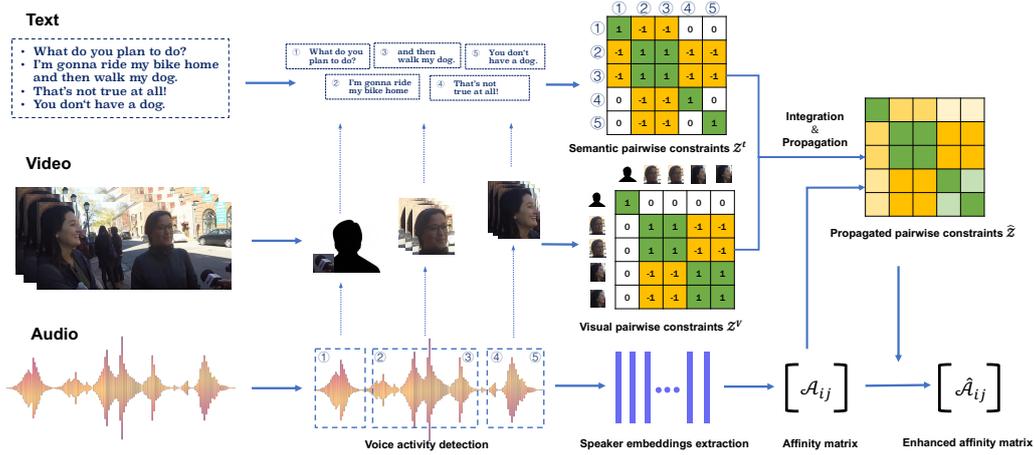


Figure 1: An overview of our proposed multimodal speaker diarization system.

**Audio-textual Speaker Diarization** Some previous works (Zuluaga-Gomez et al., 2022; Flemtomos and Narayanan, 2022; Park and Georgiou, 2018; Paturi et al., 2023) utilized semantic information derived from transcription to estimate the role profiles and detect speaker change point, demonstrating improvement in specific role-playing conversations, such as job interviews and doctor-patient medical consultations. Other works (Kanda et al., 2021; Xia et al., 2022; Khare et al., 2022) enhanced ASR models to capture speaker identity through joint training of paired audio and textual data, which typically require substantial annotated multi-speaker speech data. More recent works (Park et al., 2023; Wang et al., 2024; Cheng et al., 2023) employed large language models as post-processing to correct word speaker-related boundaries according to local semantic context.

## 2.2 Pairwise Constrained Clustering

Speaker diarization systems typically rely on unsupervised clustering to handle an unknown number of speakers. When integrating multimodal information, direct cross-modal similarity comparisons are not feasible. Thus, incorporating semi-supervised signals into the clustering process becomes essential, a technique known as constrained clustering (Bibi et al., 2023). Pairwise constrained clustering is a common approach within this framework, where supplementary information defines pairwise relationships among samples through **Must-link constraints** (indicating two samples belong to the same class) and **Cannot-link constraints** (indicating they do not) (Davidson and Ravi, 2007). The process of refining the affinity matrix using these pairwise constraints is referred to as pairwise con-

strained propagation. Initially confined to data mining domain (Hoi et al., 2007), the application of pairwise constrained clustering has expanded into multimodal areas such as vision and text (Yang et al., 2014; Yan et al., 2006). Advancing with theoretical progress, pairwise constraint propagation algorithms have increasingly integrated complex optimization techniques, including Lyapunov equation (Lu and Peng, 2011), Non-negative Matrix Factorization (NMF)(Fu, 2015), Inexact Augmented Lagrange Multiplier (IALM)(Liu et al., 2019), and deep learning outcomes (Zhang et al., 2021a,b). Among them, E2CP (Exhaustive and Efficient Constraint Propagation) (Lu and Peng, 2011) is widely adopted due to its simple and effective hyperparameter configuration. In this paper, we employ E2CP as the core pairwise constrained clustering method to integrate multimodal constraints into speaker clustering.

## 3 Methodology

Figure 1 provides an overview of how our approach leverages multimodal information. In addition to a clustering-based speaker diarization system, **video and text processing modules** are incorporated to independently extract visual and semantic information and derive pairwise constraints. Then a **joint propagation algorithm** will be employed to operate cross-modal pairwise constraints to enhance the affinity matrix constructed from acoustic speaker embeddings. The **enhanced affinity matrix** is subsequently integrated into the subsequent clustering procedure to assign speaker label for each speaker embedding. The following sections will present the joint propagation algorithm and the process of constructing visual and semantic constraints.

### 3.1 Joint Pairwise Constraint Propagation with multimodal Information

Considering that the audio contains comprehensive speaker-related information over time, we employ audio-based models, specifically a VAD model and a speaker embedding extractor, to obtain a sequence of acoustic speaker embeddings  $E = \{e_1, e_2, \dots, e_N | e_i \in \mathbb{R}^D\}$  by applying sliding windows to the audio data where  $D$  represents the dimension of the speaker embeddings and  $N$  denotes the number of speaker embeddings. Subsequently, we compute the affinity matrix  $\mathcal{A} = \{\mathcal{A}_{ij}\}_{N \times N}$ , where  $\mathcal{A}_{ij} = g(e_i, e_j)$  and  $g(\cdot)$  represents the measurement of similarity.

Assuming we have access to speaker-related cues from additional sources of information, we can derive various types of constraint pairs: must-link  $\mathcal{M}$  and cannot-link  $\mathcal{C}$ , defined as:

$$\begin{aligned} \mathcal{M}^k &= \{(e_i, e_j) | l(e_i) = l(e_j)\}, \\ \mathcal{C}^k &= \{(e_i, e_j) | l(e_i) \neq l(e_j)\}, \end{aligned} \quad (1)$$

where  $l(\cdot)$  denotes the speaker label associated with an acoustic speaker embedding, and  $k$  is the index of sources type. For different modality information, the criteria for establishing  $\mathcal{M}$  and  $\mathcal{C}$  are different, which will be described in Sec. 3.2 and Sec. 3.3 according to specific situation. Then each constraint is initially encoded into a matrix  $\mathcal{Z}^k$ :

$$\mathcal{Z}_{ij}^k = \begin{cases} +1 & \text{if } (e_i, e_j) \in \mathcal{M}^k, \\ -1 & \text{if } (e_i, e_j) \in \mathcal{C}^k, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

A series of constraint matrix  $\mathcal{Z}^k$  are integrated into a final constraint matrix  $\mathcal{Z}$ . During the integration process, some scenarios are relatively straightforward. For instance, if an embedding pair  $(e_i, e_j)$  belongs to  $\bigcap_k \mathcal{M}^k$ , then  $(e_i, e_j)$  is considered as a must-link constraint pair. Conversely, if  $(e_i, e_j)$  resides in  $\bigcap_k \mathcal{C}^k$ , it is a cannot-link constraint pair due to agreement between all modalities. However, there are evidently more complex scenarios, where the constraint matrices conflict with one another, such as  $(e_i, e_j) \in (\mathcal{M}^1 \cap \mathcal{C}^2)$  or  $(e_i, e_j) \in (\mathcal{M}^2 \cap \mathcal{C}^1)$ . To address these issues, we introduce acoustic information as the arbiter in the final determination. To summarize, we compute the integrated constraint scores following the given formula:

$$\mathcal{Z}' = \sum_k \alpha_k \mathcal{Z}^k + \beta \mathcal{A} - \theta \quad (3)$$

where  $\alpha_k, \beta$  represent the weight hyper-parameters for different modalities, and  $\theta$  is the bias. Then,  $\mathcal{Z}'$  is converted into a binarized constraint matrix  $\mathcal{Z}$  according to a threshold  $\delta$ .

$$\mathcal{Z}_{ij} = \begin{cases} +1 & \text{if } \mathcal{Z}'_{ij} > \delta, \\ -1 & \text{if } \mathcal{Z}'_{ij} < -\delta, \\ 0 & \text{else.} \end{cases} \quad (4)$$

The constraint matrix  $\mathcal{Z}$  may be sparse and constraints information is confined to discrete. It is essential to deploy a constraint propagation algorithm to efficiently broadcast the constraint information in  $\mathcal{Z}$  on a larger scale. Specifically, we employ E2CP (Lu and Peng, 2011) algorithm to obtain propagated constraints  $\hat{\mathcal{Z}}$ :

$$\hat{\mathcal{Z}} = (1 - \lambda)^2 (\mathbf{I} - \lambda \mathbf{L}_e)^{-1} \mathcal{Z} (\mathbf{I} - \lambda \mathbf{L}_e)^{-1}, \quad (5)$$

where  $\mathbf{L}_e = \mathbf{D}_e^{-1/2} \mathcal{A} \mathbf{D}_e^{-1/2}$  is the normalized Laplacian matrix, and  $\mathbf{D}_e$  is the degree matrix of  $\mathcal{A}$  and  $\mathbf{I}$  is a identity matrix. The parameter  $\lambda \in [0, 1]$  modulates the impact degree of the constraints. The refined affinity matrix  $\hat{\mathcal{A}} \in \mathbb{R}^{N \times N}$  is then updated to incorporate the influences of the propagated constraints  $\hat{\mathcal{Z}}$ :

$$\hat{\mathcal{A}}_{ij} = \begin{cases} 1 - (1 - \hat{\mathcal{Z}}_{ij})(1 - \mathcal{A}_{ij}) & \text{if } \hat{\mathcal{Z}}_{ij} \geq 0, \\ (1 + \hat{\mathcal{Z}}_{ij})\mathcal{A}_{ij} & \text{if } \hat{\mathcal{Z}}_{ij} < 0. \end{cases} \quad (6)$$

Upon calculating the affinity matrix  $\hat{\mathcal{A}}$ , it is then fed into the clustering process to derive the ultimate speaker diarization results. **It is worth noting that there is no limit to the number of constraint types  $k$ .** We can extract diverse constraint matrices related to different modal data. These constraint matrices can be considered as prior knowledge, guiding the clustering focus towards a specific perspective of the scenario. In this paper, we fix  $k$  at 2, thereby extracting two distinct constraint types: visual constraint  $\mathcal{Z}^v$  and semantic constraint  $\mathcal{Z}^t$ .

### 3.2 Visual constraints construction

The speaker-related visual constraints is constructed through the following steps, similar to (Chung et al., 2020; Xu et al., 2022): **(1) Face Tracking.** The first step involves detecting and tracking faces in video frames over time using a CNN-based face detector (Liu et al., 2018) and a position-based tracker. Only face tracks aligned with speech segments detected by VAD are retained for further processing. **(2) Active Speaker Detection.** This step determines whether tracked faces

correspond to active speakers at any given moment. A two-stream network (Tao et al., 2021), comprising temporal encoders and an attention-based decoder, analyzes audio-visual synchrony to identify speaker activity. Low-confidence frames are filtered using a predefined threshold. **(3) Face Clustering.** A face recognition CNN (Huang et al., 2020) extracts embeddings from face tracks at uniform intervals (e.g., every 200 ms). These embeddings are then clustered with AHC.

By integrating these steps, constraints based on visual information are obtained. Faces clustered to the same speaker are considered as must-link constraints, while those clustered to different speakers are cannot-link constraints. Each face is aligned with respective acoustic embeddings along the time axis. If an acoustic embedding corresponds to multiple faces, we will select the speaker associated with the majority of those faces.

### 3.3 Semantic constraints construction

To extract speaker-related information from the transcriptions, we construct two Spoken Language Processing (SLP) tasks: (1) **Dialogue Detection** discriminates between multi-speaker dialogues and monologues, conceptualized as a binary classification challenge. (2) **Speaker-Turn Detection** assesses each sentence in a sequence to estimate speaker change, functioning as a sequence labeling problem that identifies semantically significant speaker role transitions. Semantic constraints can be formulated based on the outputs of these two tasks. Specifically, must-link  $\mathcal{M}^t$  is formed between two embeddings if they are sourced from the same non-dialogue segment. Conversely, cannot-link  $\mathcal{C}^t$  is established between embeddings separated by a detected speaker-turn boundary, as illustrated in Figure 2.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets.** Our experiments are conducted on the AIShell-4 (Fu et al., 2021), Alimeeting (Yu et al., 2022), AVA-AVD (Xu et al., 2022), and our proposed datasets. The AVA-AVD dataset, which focuses on audio-visual diarization, provides diverse scenarios and face annotations but lacks ground truth transcripts. In contrast, AIShell-4 and Alimeeting are Mandarin datasets that include speaker-labeled transcripts, making them well-suited for audio-text-based tasks. Due to the

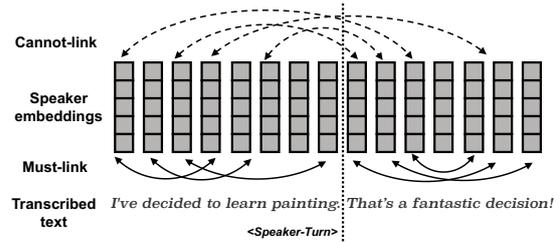


Figure 2: Semantic constraint construction is based on dialogue and speaker-turn detection. Text segments identified as non-dialogue imply that their embeddings are related through must-link constraints (solid connections). Conversely, detected transition points indicate that embeddings spanning these points should be connected via cannot-link constraints (dashed connections).

absence of publicly available evaluation datasets with annotations for visual, semantic, and acoustic modalities, we construct a new dataset comprising 6.3 hours of video, manually annotated with speaker timestamps and speech content. Further details about this dataset are provided in the Appendix A. The combination of these diverse datasets further demonstrates the effectiveness of our methods across different domains.

**Implementation Details.** In our system, the audio-based diarization modules follow the pipeline outlined in (Cheng et al., 2023). Our speaker embedding extractor is an adaptation of CAM++ (Wang et al., 2023), which has been trained on VoxCeleb dataset (Nagrani et al., 2020). The ASR we utilize is Paraformer (Gao et al., 2022), which has been trained with the aid of the FunASR (Gao et al., 2023) toolkits. For visual components, we employ a series of pre-trained models for different tasks: RFB-Net (Liu et al., 2018) for face detection, TalkNet (Tao et al., 2021) for active speaker detection, and CurricularFace model (Huang et al., 2020) for extracting face embeddings. For semantic tasks, we train models on open-sourced meeting datasets designed for various scenarios. Specifically, we use separate datasets for English and Mandarin to train corresponding semantic models, ensuring language-specific adaptations. All that training was conducted using a pre-trained BERT model (Devlin et al., 2019). We employ E2CP as the core algorithm for constraint propagation. In the post-clustering phase, our system adheres to the SC algorithm. Inspired by the work presented in (Park et al., 2020), our method incorporates refinement operations, such as row-wise thresholding and symmetrization, to enhance the performance of

Table 1: The results of speaker diarization on AIShell-4, Alimeeting and our proposed dataset.

Dataset	Methods	Modality	DER(%)↓	CpWER(%)↓
AIShell-4	Pyannote	Audio	12.2	-
	DiariZen	Audio	11.7	-
	Semantic-Aux SD	Audio + Semantic	-	15.23
	Proposed	Audio + Semantic	12.07	14.95
Alimeeting	Pyannote	Audio	24.4	-
	DiariZen	Audio	17.6	-
	Semantic-Aux SD	Audio + Semantic	-	36.15
	Proposed	Audio + Semantic	21.32	31.11
Proposed Dataset	Pyannote	Audio	15.57	-
	DiariZen	Audio	10.49	-
	CAM++ & VBx	Audio	10.31	18.03
	CAM++ & SC	Audio	9.37	17.04
	Proposed	Audio + Semantic	9.12	16.86
	Proposed	Audio + Visual	9.13	16.83
	Proposed	Audio + Semantic + Visual	9.01	16.36

spectral clustering. More details of the implementation and hyperparameter settings can be found in the appendix B.

**Evaluation Metrics.** To demonstrate the impact of the speaker diarization system, we report the Diarization Error Rate (DER) (Fiscus et al., 2006), which generally composed of three parts: Missed Speech (MS), False Alarms (FA) and Speaker Error (SPKE). As the ASR and forced-alignment module have been used in the pipeline, we also report the Concatenated Minimum-permutation Word Error Rate (Watanabe et al., 2020).

## 4.2 Main Results

**Multimodal Speaker Diarization with Audio-Text Modalities.** In Part 1 & 2 of Table 1, we compare our system with acoustic-only speaker diarization systems such as Pyannote and DiariZen, as well as with Semantic-Aux SD, an audio-text speaker diarization system, on the AIShell-4 and Alimeeting datasets. It can be observed that our system demonstrates a certain level of superiority over the classical speaker diarization toolkit, Pyannote, on both datasets. Specifically, our method achieves an absolute improvement of 0.13% in DER on the AIShell-4 dataset and 3.08% in DER on the Alimeeting dataset. However, our approach shows a slight disadvantage compared to DiariZen, which trains a frame-level EEND model using the training sets from AMI, Alimeeting, and AIShell-4 audio data, thus providing it with a noticeable advantage on these two homologous test sets.

Table 2: The results of audio-visual speaker diarization experiments on AVA-AVD datasets.

Models	SPKE(%)↓	DER(%)↓
AVR-Net	24.88	27.43
AFL-Net	21.10	23.65
AFL-Net + WavLM	19.57	22.12
DyViSE	20.86	23.46
Proposed	17.40	20.32

Compared with Semantic-Aux SD (Cheng et al., 2023), another speaker diarization method that combines audio and text modalities, our experiments maintain consistent ASR results. Our proposed solution shows clear improvements on both datasets, with an absolute gain of 5.04% in CpWER on the Alimeeting dataset. Unlike Semantic-Aux SD, which primarily uses semantic information for boundary refinement of speaker diarization results, our approach **integrates semantic information into the speaker clustering process, leveraging semantic cues to correct more errors** that arise from relying solely on acoustic-only information.

**Multimodal Speaker Diarization with Audio-Visual Modalities.** We have also compared our approach with several audio-visual joint training speaker diarization methods on the AVA-AVD dataset, such as AVR-Net (Xu et al., 2022), AFL-Net (Yin et al., 2024), and DyViSE (Wuerkaixi et al., 2022), to demonstrate the effectiveness of our method. Due to the lack of annotated transcripts in the AVA-AVD dataset, only the SPKE and DER

metrics are reported. Table 2 presents a comparison conducted on AVA-AVD, revealing the competitive performance of our method when utilizing only visual constraints. Compared to the baseline results of the AVA-AVD, AVR-Net, our model shows a 7.11% absolute improvement in DER. Additionally, when compared to audio-visual models such as AFL-Net and DyViSE, our model also exhibits a significant improvement in DER, with a relative 8.1% improvement over AFL-Net and a relative 13.3% improvement over DyViSE.

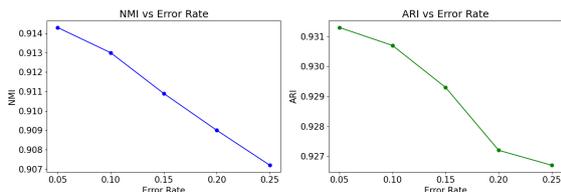


Figure 3: Simulated constraints with errors and the effect for constrained clustering

**Multimodal Speaker Diarization with Audio-Visual-Text Modalities.** In the last part of Table 1, we present the results of multiple speaker diarization systems on the proposed dataset. As previously mentioned, the acoustic-only speaker diarization SOTA(start-of-the-art) system, DiariZen, benefits from targeted training on the AIShell-4 and AI-meeting training sets, giving it a certain advantage over our method on these datasets. However, on the proposed dataset, our “CAM++ & SC” approach achieves a relative improvement of 10.7% in DER compared to DiariZen. Furthermore, when incorporating the semantic constraints proposed in this paper, the improvement over DiariZen reaches a relative 13.1% reduction in DER. In contrast, another well-known open-source speaker diarization toolkit, pyannote, exhibits a significantly higher DER of 15.57%, showing a larger gap compared to other approaches. When comparing the constrained propagation methods proposed in this paper, the results show that using only visual constraints or only semantic constraints result in very similar performance in terms of DER and CpWER. Compared to the acoustic-only baseline, incorporating visual constraints leads to a relative improvement of 2.56% in DER and 1.1% in CpWER, while incorporating semantic constraints yields a relative improvement of 2.66% in DER and 1.2% in CpWER. Further analysis reveals that combining **all three modalities provides even greater improvements over systems that combine only two**

**modalities.** Specifically, the DER is relatively reduced by 3.8%, and the CpWER is relatively reduced by 4.0%. This indicates that **semantic and visual constraints are complementary, and integrating multiple modalities can lead to further performance gains.** Some decoding cases and visualizations can be found in the appendix E.

### 4.3 Analysis and Discussion

**Constraint Construction.** It should be noted that constraints constructed based on multimodal data often contain some errors, and at the same time, constraints cannot cover all embedding pairs. In this section, we discuss the impact of varying quantities and qualities of pairwise constraints on the results of speaker clustering. We employ several simulation strategies to generate pairwise constraints, which allows for better control over both the quantity and quality of these constraints. All experiments are conducted on our proposed dataset, utilizing speaker embeddings extracted by CAM++ that remain fixed throughout the experiments; only the pairwise constraints used in each experiment are varied. The cluster metrics, Normalized Mutual Information (NMI) (Strehl and Ghosh, 2002) and Adjusted Rand Index (ARI) (Chac’ on and Rastrojo, 2020), will be reported in this section.

**(1) The Impact of constraint Quality.** In practice, the constraints we obtain often contain many errors. This is especially common in multi-party meeting or interview scenarios, such as when there is audio-visual asynchrony or errors from transcript text decoded by ASR due to complex acoustic environments. In order to investigate the impact of incorrect constraints on our method, we have established the following randomization strategy: First, we randomly generate a completely correct set of constraints, including must-links and cannot-links. We then randomly alter the status of a proportion  $p_{err}$  of these constraints—turning must-links into cannot-links and vice-versa—thereby introducing a certain level of constraint errors while keeping the total number of constraints constant. In our experiments,  $p_{err} \in \{5\%, 10\%, 15\%, 20\%, 25\%\}$ . The experimental results are illustrated in Figure 3. It can be observed that errors in the constraints do indeed lead to a decline in clustering performance. However, even when the error rate reaches 25%, the NMI experiences only a 0.7% relative decrease compared to the NMI at a 5% error rate. **This indicates that our method exhibits a certain degree of robustness to erroneous constraints.**

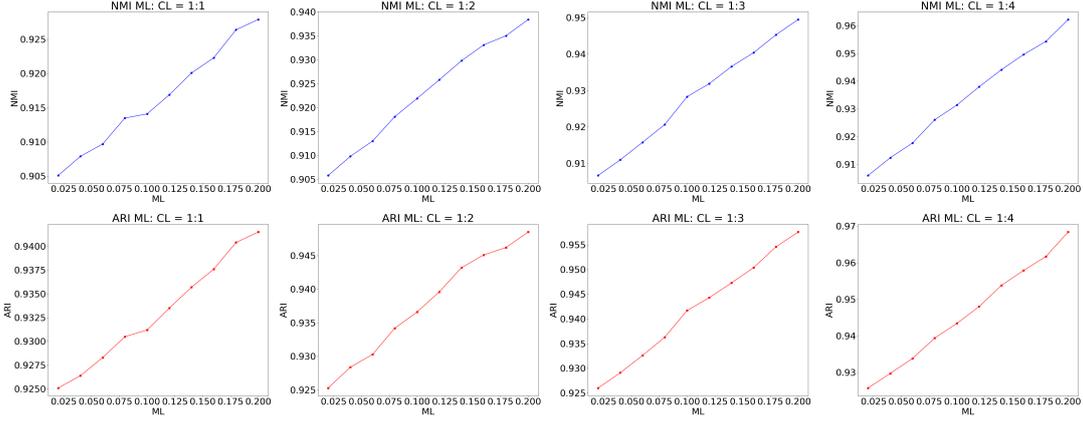


Figure 4: Results of constrained speaker cluster performance across various levels of constraints coverage, showcasing scenarios with imbalanced proportions of must-link and cannot-link constraints.

Table 3: The gain achieved by our proposed method, in comparison to the acoustic-only approach, varies across cases with different ASR levels in Alimeeting.

Test Subsets	DER(%)	DER(%)	DER relative gain(%)
	acoustic only	proposed models	
Easy subsets (WER <21%)	7.31	6.65	8.9
Hard subsets (WER >21%)	38.02	31.45	17.4

## (2) Impact of Constraint Quantity and Ratios

We investigate how the number of constraints and the ratio of must-link to cannot-link sets affect speaker clustering in our approach. We formulate a simulation strategy: for a sequence of speaker embeddings  $E = \{e_1, e_2, \dots, e_N | e_i \in \mathbb{R}^D\}$ , we vary the must-link coverage ( $p_{ml}$ ) and cannot-link coverage ( $p_{cl}$ ) proportions. Specifically,  $p_{ml} \in \{2\%, 4\%, 6\%, \dots, 20\%\}$ , and  $p_{cl} = k_{ratio} \times p_{ml}$  with  $k_{ratio} \in \{1, 2, 3, 4\}$ . We select  $p_{ml}\%$  of must-links and  $p_{cl}\%$  of cannot-links from all possible pairs. As shown in Figure 4, it can be observed that as the number of constraints increases, the clustering performance of the algorithm consistently improves. For instance, in the case of ML:CL = 1:1, the NMI increases from 0.905 to 0.925, and the ARI improves from 0.925 to 0.940. Additionally, our method demonstrates that **an imbalance between the quantities of ML and CL does not hinder performance gains.**

**Impact of ASR for Semantic Constraints.** In practice, visual and audio data are often collected independently using different devices, with semantic information extracted from audio via an ASR system. Complex acoustic environments can affect both speaker embedding extraction and ASR accu-

racy. In this section, we evaluate the performance of our method under varying ASR accuracy levels. Specifically, we first partition the Alimeeting test set into "Easy" and "hard" subsets based on whether the ASR Word Error Rate (WER) exceeds 21%. We then test both the acoustic-only solutions (CAM++ & SC) and our proposed method, which incorporates semantic constraints, on these two subsets. Table 3 presents our experimental results. It can be observed that on the "hard" subset, the DER is relatively higher, indicating that both speaker embedding extraction and the ASR system encounter certain errors in complex acoustic environments. Nevertheless, our method achieves notable improvements on both subsets. On the "Easy" subset, where the acoustic-only solution already performs well (DER = 7.31%), our approach achieves a relative improvement of 8.9%. On the "hard" subset, our method demonstrates a significant relative improvement of 17.4%, showcasing its **robustness in challenging acoustic conditions.**

## 5 Conclusions

In this study, we propose a novel multimodal approach that leverages audio, visual, and semantic information to enhance speaker diarization. By incorporating additional visual and textual processing modules, we generate complementary pairwise constraints that are integrated into the clustering process through a joint pairwise constraint propagation method. Experimental results demonstrate significant performance improvements. This research contributes to the advancement of more sophisticated systems for the speaker diarization task, providing potential directions for future exploration.

## 6 Limitations

The semantic information we utilize is derived from the Dialogue Detection and Speaker-Turn Detection models, which are trained based on the BERT architecture, rather than employing more advanced Large Language Models. Additionally, we observe that the embeddings corresponding to the constraints constructed for the current semantic tasks—Dialogue Detection and Speaker-Turn Detection—are relatively close to each other in the embedding space. While this proximity aids in identifying more precise speaker transition points, it also limits our ability to extract long-term semantic information from the text. In theory, our approach can be adapted to incorporate various multimodal sources of information. However, another modality that could significantly assist in speaker identity determination—speaker location information—has not been integrated into our experiments. We plan to explore this further in future work.

## References

Xavier Anguera, Simon Bozonnet, Nicholas Evans, Corinne Fredouille, Gerald Friedland, and Oriol Vinyals. 2012. Speaker diarization: A review of recent research. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(2):356–370.

Adel Bibi, Ali Alqahtani, and Bernard Ghanem. 2023. **Constrained clustering: General pairwise and cardinality constraints**. *IEEE Access*, 11:5824–5836.

Hervé Bredin. 2023. pyannote.audio 2.1 speaker diarization pipeline: principle, benchmark, and recipe. In *Proc. INTERSPEECH 2023*.

Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Maël Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, Guillaume Lathoud, Mike Lincoln, Agnes Lisowska, Iain McCowan, Wilfried Post, Dennis Reidsma, and Pierre Wellner. 2005. **The AMI meeting corpus: A pre-announcement**. In *Machine Learning for Multimodal Interaction, Second International Workshop, MLMI 2005, Edinburgh, UK, July 11-13, 2005, Revised Selected Papers*, volume 3869 of *Lecture Notes in Computer Science*, pages 28–39. Springer.

Jos’e E. Chac’on and Ana I. Rastrojo. 2020. Minimum adjusted rand index for two clusterings of a given size. *Advances in Data Analysis and Classification*, 17:125–133.

Yafeng Chen, Siqi Zheng, Hui Wang, Luyao Cheng, Qian Chen, and Jiajun Qi. 2023. An enhanced res2net with local and global feature fusion for speaker verification. *CoRR*, abs/2305.12838.

Luyao Cheng, Siqi Zheng, Zhang Qinglin, Hui Wang, Yafeng Chen, and Qian Chen. 2023. Exploring speaker-related information in spoken language understanding for better speaker diarization. In *Annual Meeting of the Association for Computational Linguistics*. 686–691.

Joon Son Chung, Jaesung Huh, Arsha Nagrani, Triantafyllos Afouras, and Andrew Senior. 2020. Spot the conversation: speaker diarisation in the wild. In *Interspeech*. 692–695.

Joon Son Chung, Bong-Jin Lee, and Icksang Han. 2019. Who said that?: Audio-visual speaker diarisation of real-world meetings. In *Interspeech*. 696–698.

Ian Davidson and S. S. Ravi. 2007. **Intractability and clustering with constraints**. In *Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007*, volume 227 of *ACM International Conference Proceeding Series*, pages 201–208. ACM. 699–704.

William H. E. Day and Herbert Edelsbrunner. 1984. **Efficient algorithms for agglomerative hierarchical clustering methods**. *Journal of Classification*, 1:7–24. 705–708.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*. 709–713.

Jonathan G. Fiscus, Jerome Ajot, Martial Michel, and John S. Garofolo. 2006. The rich transcription 2006 spring meeting recognition evaluation. In *Machine Learning for Multimodal Interaction, Third International Workshop, MLMI 2006, Bethesda, MD, USA, May 1-4, 2006, Revised Selected Papers*, volume 4299 of *Lecture Notes in Computer Science*, pages 309–322. Springer. 714–721.

Nikolaos Flemotomos and Shrikanth Narayanan. 2022. Multimodal clustering with role induced constraints for speaker diarization. 722–724.

Yihui Fu, Luyao Cheng, Shubo Lv, Yukai Jv, Yuxiang Kong, Zhuo Chen, Yanxin Hu, Lei Xie, Jian Wu, Hui Bu, Xin Xu, Jun Du, and Jingdong Chen. 2021. AISHELL-4: an open source dataset for speech enhancement, separation, recognition and speaker diarization in conference scenario. In *22nd Annual Conference of the International Speech Communication Association, Interspeech 2021, Brno, Czechia, August 30 - September 3, 2021*, pages 3665–3669. ISCA. 725–734.

Zhenyong Fu. 2015. **Pairwise constraint propagation via low-rank matrix recovery**. *Comput. Vis. Media*, 1(3):211–220. 735–737.

Yusuke Fujita, Naoyuki Kanda, Shota Horiguchi, Kenji Nagamatsu, and Shinji Watanabe. 2019a. End-to-End Neural Speaker Diarization with Permutation-free Objectives. In *Interspeech*, pages 4300–4304. 738–741.

742	Yusuke Fujita, Naoyuki Kanda, Shota Horiguchi, Yawen Xue, Kenji Nagamatsu, and Shinji Watanabe. 2019b. <a href="#">End-to-end neural speaker diarization with self-attention</a> . <i>2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 296–303.	798
743		799
744		800
745		801
746		802
747		803
748	Yusuke Fujita, Naoyuki Kanda, Shota Horiguchi, Yawen Xue, Kenji Nagamatsu, and Shinji Watanabe. 2019c. <a href="#">End-to-end neural speaker diarization with self-attention</a> . In <i>2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 296–303.	804
749		
750		805
751		806
752		807
753		808
754	Yusuke Fujita, Shinji Watanabe, Shota Horiguchi, Yawen Xue, and Kenji Nagamatsu. 2020. <a href="#">End-to-end neural diarization: Reformulating speaker diarization as simple multi-label classification</a> . <i>ArXiv</i> , abs/2003.02966.	809
755		810
756		811
757		812
758		
759	Ananya Ganesh, Martha Palmer, and Katharina Kann. 2023a. <a href="#">A survey of challenges and methods in the computational modeling of multi-party dialog</a> . <i>Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)</i> .	813
760		814
761		815
762		816
763		817
764	Ananya Ganesh, Martha Palmer, and Katharina Kann. 2023b. <a href="#">A survey of challenges and methods in the computational modeling of multi-party dialog</a> . In <i>Proceedings of the 5th Workshop on NLP for Conversational AI (NLP4ConvAI 2023)</i> , pages 140–154, Toronto, Canada. Association for Computational Linguistics.	818
765		
766		819
767		820
768		821
769		822
770		823
771	Zhifu Gao, Zerui Li, Jiaming Wang, Haoneng Luo, Xian Shi, Mengzhe Chen, Yabin Li, Lingyun Zuo, Zhihao Du, Zhangyu Xiao, and Shiliang Zhang. 2023. <a href="#">Funasr: A fundamental end-to-end speech recognition toolkit</a> . <i>ArXiv</i> , abs/2305.11013.	824
772		825
773		826
774		
775		827
776	Zhifu Gao, Shiliang Zhang, Ian Mcloughlin, and Zhijie Yan. 2022. <a href="#">Paraformer: Fast and accurate parallel transformer for non-autoregressive end-to-end speech recognition</a> . In <i>Interspeech</i> .	828
777		829
778		
779		830
780	Israel D Gebru, Sileye Ba, Xiaofei Li, and Radu Horaud. 2017. <a href="#">Audio-visual speaker diarization based on spatiotemporal bayesian fusion</a> . <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 40(5):1086–1099.	831
781		832
782		833
783		834
784		835
785	Gregory Gelly and Jean-Luc Gauvain. 2018. <a href="#">Optimization of rnn-based speech activity detection</a> . <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> , 26(3):646–656.	836
786		837
787		838
788		
789	Jiangyu Han, Federico Landini, Johan Rohdin, Anna Silnova, Mireia Diez, and Lukas Burget. 2024. <a href="#">Leveraging self-supervised learning for speaker diarization</a> . <i>arXiv preprint arXiv:2409.09408</i> .	839
790		840
791		841
792		842
793	Steven C. H. Hoi, Rong Jin, and Michael R. Lyu. 2007. <a href="#">Learning nonparametric kernel matrices from pairwise constraints</a> . In <i>Machine Learning, Proceedings of the Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007</i> , volume 227 of <i>ACM International Conference Proceeding Series</i> , pages 361–368. ACM.	843
794		844
795		845
796		846
797		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

855	Songtao Liu, Di Huang, and Yunhong Wang. 2018.	Laurens van der Maaten and Geoffrey Hinton. 2008.	909
856	Receptive field block net for accurate and fast object	<a href="#">Visualizing data using t-sne</a> . <i>Journal of Machine</i>	910
857	detection. In <i>The European Conference on Computer</i>	<i>Learning Research</i> , 9(86):2579–2605.	911
858	<i>Vision (ECCV)</i> .		
859	Zhiwu Lu and Yuxin Peng. 2011. Exhaustive and effi-	Ulrike Von Luxburg. 2007. A tutorial on spectral clus-	912
860	cient constraint propagation: A graph-based learning	tering. <i>Statistics and computing</i> , 17:395–416.	913
861	approach and its applications. <i>International Journal</i>	Hui Wang, Siqi Zheng, Yafeng Chen, Luyao Cheng,	914
862	<i>of Computer Vision</i> , 103:306–325.	and Qian Chen. 2023. Cam++: A fast and efficient	915
863	Arsha Nagrani, Joon Son Chung, Weidi Xie, and	network for speaker verification using context-aware	916
864	Andrew Zisserman. 2020. Voxceleb: Large-scale	masking.	917
865	speaker verification in the wild. <i>Comput. Speech</i>	Quan Wang, Carlton Downey, Li Wan, Philip Andrew	918
866	<i>Lang.</i> , 60.	Mansfield, and Ignacio Lopez Moreno. 2018. <a href="#">Speaker</a>	919
867	Tae Jin Park, Kunal Dhawan, Nithin Rao Koluguri, and	<a href="#">diarization with lstm</a> . In <i>2018 IEEE International</i>	920
868	Jagadeesh Balam. 2023. Enhancing speaker diariza-	<i>Conference on Acoustics, Speech and Signal Process-</i>	921
869	tion with large language models: A contextual beam	<i>ing (ICASSP)</i> , pages 5239–5243.	922
870	search approach. <i>ArXiv</i> , abs/2309.05248.	Quan Wang, Yiling Huang, Guanlong Zhao, Evan Clark,	923
871	Tae Jin Park and Panayiotis G. Georgiou. 2018. Mul-	Wei Xia, and Hank Liao. 2024. Diarizationlm:	924
872	timodal speaker segmentation and diarization using	Speaker diarization post-processing with large lan-	925
873	lexical and acoustic cues via sequence to sequence	guage models. <i>ArXiv</i> , abs/2401.03506.	926
874	neural networks. In <i>Interspeech</i> .	Shinji Watanabe, Michael Mandel, Jon Barker, and Em-	927
875	Tae Jin Park, Kyu J. Han, Manoj Kumar, and	manuel Vincent. 2020. Chime-6 challenge: Tackling	928
876	Shrikanth S. Narayanan. 2020. Auto-tuning spectral	multispeaker speech recognition for unsegmented	929
877	clustering for speaker diarization using normalized	recordings. <i>ArXiv</i> , abs/2004.09249.	930
878	maximum eigengap. <i>IEEE Signal Processing Letters</i> ,	Abudukelimu Wuerkaixi, Kunda Yan, You Zhang,	931
879	27:381–385.	Zhiyao Duan, and Changshui Zhang. 2022. Dyvise:	932
880	Tae Jin Park, Naoyuki Kanda, Dimitrios Dimitri-	Dynamic vision-guided speaker embedding for audio-	933
881	adis, Kyu J. Han, Shinji Watanabe, and Shrikanth	visual speaker diarization. In <i>24th IEEE Interna-</i>	934
882	Narayanan. 2022. A review of speaker diarization:	<i>tional Workshop on Multimedia Signal Processing,</i>	935
883	Recent advances with deep learning. <i>Comput. Speech</i>	<i>MMSP 2022, Shanghai, China, September 26-28,</i>	936
884	<i>Lang.</i> , 72(C).	2022, pages 1–6. IEEE.	937
885	Rohit Paturi, Sundararajan Srinivasan, and Xiang Li.	Wei Xia, Han Lu, Quan Wang, Anshuman Tripathi, Yil-	938
886	2023. Lexical speaker error correction: Leveraging	ing Huang, Ignacio Lopez Moreno, and Hasim Sak.	939
887	language models for speaker diarization error correc-	2022. Turn-to-diarize: Online speaker diarization	940
888	tion. <i>ArXiv</i> , abs/2306.09313.	constrained by transformer transducer speaker turn	941
889	Weizhou Shen, Xiaojun Quan, and Ke Yang. 2023.	detection. In <i>ICASSP 2022 - 2022 IEEE Interna-</i>	942
890	<a href="#">Generic dependency modeling for multi-party conver-</a>	<i>tional Conference on Acoustics, Speech and Signal</i>	943
891	<a href="#">sation</a> . In <i>ICASSP 2023 - 2023 IEEE International</i>	<i>Processing (ICASSP)</i> , pages 8077–8081.	944
892	<i>Conference on Acoustics, Speech and Signal Process-</i>	Eric Zhongcong Xu, Zeyang Song, Satoshi Tsutsui,	945
893	<i>ing (ICASSP)</i> , pages 1–5.	Chao Feng, Mang Ye, and Mike Zheng Shou. 2022.	946
894	David Snyder, Daniel Garcia-Romero, Gregory Sell,	Ava-avd: Audio-visual speaker diarization in the wild.	947
895	Daniel Povey, and Sanjeev Khudanpur. 2018. X-	MM '22, page 3838–3847.	948
896	vectors: Robust dnn embeddings for speaker recog-	Rong Yan, Jian Zhang, Jie Yang, and Alexander G.	949
897	nition. In <i>2018 IEEE International Conference on</i>	Hauptmann. 2006. <a href="#">A discriminative learning frame-</a>	950
898	<i>Acoustics, Speech and Signal Processing (ICASSP)</i> ,	<a href="#">work with pairwise constraints for video object clas-</a>	951
899	pages 5329–5333.	<a href="#">sification</a> . <i>IEEE Trans. Pattern Anal. Mach. Intell.</i> ,	952
900	Alexander Strehl and Joydeep Ghosh. 2002. Cluster en-	28(4):578–593.	953
901	sembles - a knowledge reuse framework for combin-	Zheng Yang, Yao Hu, Haifeng Liu, Huajun Chen, and	954
902	ing multiple partitions. <i>J. Mach. Learn. Res.</i> , 3:583–	Zhaohui Wu. 2014. <a href="#">Matrix completion for cross-view</a>	955
903	617.	<a href="#">pairwise constraint propagation</a> . In <i>Proceedings of</i>	956
904	Ruijie Tao, Zexu Pan, Rohan Kumar Das, and et al. 2021.	<i>the ACM International Conference on Multimedia,</i>	957
905	Is someone speaking? exploring long-term temporal	<i>MM '14, Orlando, FL, USA, November 03 - 07, 2014,</i>	958
906	features for audio-visual active speaker detection. In	pages 897–900. ACM.	959
907	<i>Proceedings of the 29th ACM International Confer-</i>	Hani C. Yehia, Philip Rubin, and Eric Vatikiotis-	960
908	<i>ence on Multimedia</i> , page 3927–3935.	Bateson. 1998. Quantitative association of vocal-	961
		tract and facial behavior. <i>Speech Commun.</i> , 26:23–	962
		43.	963

964	Yongkang Yin, Xu Li, Ying Shan, and Yuexian Zou.	to 29 minutes. The dataset covers a wide range of	1019
965	2024. Afl-net: Integrating audio, facial, and lip	scenarios, including interviews, talk shows, meet-	1020
966	modalities with cross-attention for robust speaker	ings, press conferences, round-table discussions,	1021
967	diarization in the wild.	and TV programs. It has been meticulously anno-	1022
968	Takuya Yoshioka, Igor Abramovski, Cem Aksoylar,	tated with speaker identity labels, corresponding	1023
969	Zhuo Chen, Moshe David, Dimitrios Dimitriadis, Yi-	speech activity timestamps, and transcribed speech	1024
970	fan Gong, Ilya Gurvich, Xuedong Huang, Yan Huang,	content.	1025
971	Aviv Hurvitz, Li Jiang, Sharon Koubi, Eyal Krupka,	These manual annotations come from annotators	1026
972	Ido Leichter, Changliang Liu, Partha Parthasarathy,	at external data company. Before annotation, we	1027
973	Alon Vinnikov, Lingfeng Wu, Xiong Xiao, Wayne	only ask for the speaker content of each speaker	1028
974	Xiong, Huaming Wang, Zhenghao Wang, Jun Zhang,	in the video, and the timestamp of each sentence.	1029
975	Yong Zhao, and Tianyan Zhou. 2019. Advances in	We provided annotators with a sample annotation	1030
976	online audio-visual meeting transcription. In <i>2019</i>	from Alimeeting, and the results they returned were	1031
977	<i>IEEE Automatic Speech Recognition and Understand-</i>	consistent with this batch of data. The speaker id	1032
978	<i>ing Workshop (ASRU)</i> , pages 276–283.	is completely anonymized in the annotation and	1033
979	Fan Yu, Shiliang Zhang, Yihui Fu, Lei Xie, Siqi Zheng,	labeled as $\{c1, c2, c3, \dots\}$	1034
980	Zhihao Du, Weilong Huang, Pengcheng Guo, Zhijie		
981	Yan, Bin Ma, Xin Xu, and Hui Bu. 2022. M2met:	<b>B Implementation details</b>	1035
982	The icassp 2022 multi-channel multi-party meeting	In this section, we provide the implementation de-	1036
983	transcription challenge. In <i>IEEE International Con-</i>	tails of our experiments.	1037
984	<i>ference on Acoustics, Speech and Signal Processing,</i>	In our system, the audio-based diarization mod-	1038
985	<i>ICASSP 2022, Virtual and Singapore, 23-27 May</i>	ules follow the pipeline outlined in (Cheng et al.,	1039
986	<i>2022</i> , pages 6167–6171. IEEE.	2023). Our speaker embedding extractor is an adap-	1040
987	Hongjing Zhang, Tianyang Zhan, Sugato Basu, and Ian	tation of CAM++ (Wang et al., 2023) <sup>1</sup> , which has	1041
988	Davidson. 2021a. A framework for deep constrained	been trained on VoxCeleb dataset (Nagrani et al.,	1042
989	clustering. <i>CoRR</i> , abs/2101.02792.	2020). To transcribe audio into text, we utilize the	1043
990	Hongjing Zhang, Tianyang Zhan, Sugato Basu, and Ian	ASR model, Paraformer (Gao et al., 2022), which	1044
991	Davidson. 2021b. A framework for deep constrained	has been trained with the aid of the FunASR (Gao	1045
992	clustering. <i>Data Min. Knowl. Discov.</i> , 35(2):593–	et al., 2023) toolkits <sup>2</sup> .	1046
993	620.	For the visual componets, we employ a se-	1047
994	Siqi Zheng, Yun Lei, and Hongbin Suo. 2020.	ries of pretrained models for different tasks:	1048
995	Phonetically-aware coupled network for short dura-	RFB-Net (Liu et al., 2018) <sup>3</sup> for face detection,	1049
996	tion text-independent speaker verification. In <i>Inter-</i>	TalkNet (Tao et al., 2021) <sup>4</sup> for active speaker de-	1050
997	<i>speech 2020, 21st Annual Conference of the Interna-</i>	tection, and CurricularFace model (Huang et al.,	1051
998	<i>tional Speech Communication Association, Virtual</i>	2020) <sup>5</sup> for extracting face embeddings.	1052
999	<i>Event, Shanghai, China, 25-29 October 2020</i> , pages	For semantic tasks, we train muliple models	1053
1000	926–930. ISCA.	with open-source meeting datasets for different	1054
1001	Juan Zuluaga-Gomez, Seyyed Saeed Sarfjoo, Amrutha	scenarios. Specifically, we employ AMI (Carletta	1055
1002	Prasad, Iuliia Nigmatulina, Petr Motlíček, Karel Ond-	et al., 2005), ICSI (Janin et al., 2003) and CHiME-	1056
1003	rej, Oliver Ohneiser, and Hartmut Helmke. 2022.	6 (Watanabe et al., 2020) to generate English se-	1057
1004	Bertraffic: Bert-based joint speaker role and speaker	semantic models, and used Alimeeting and AIShell-4	1058
1005	change detection for air traffic control communica-		
1006	tions. In <i>IEEE Spoken Language Technology Work-</i>		
1007	<i>shop, SLT 2022, Doha, Qatar, January 9-12, 2023</i> ,		
1008	pages 633–640. IEEE.		

## 1009 A Proposed Dataset

1010 In this section, we provide an overview of the pro-  
1011 posed dataset. The dataset includes a variety of  
1012 acoustic and visual scenarios, sourced from in-the-  
1013 wild videos. The dataset includes a total of 92  
1014 speakers, with each meeting involving 2 to 10 par-  
1015 ticipants. The data exhibits significant variability  
1016 in both content and environmental conditions. The  
1017 total duration of the dataset is approximately 6.3  
1018 hours, with individual video clips ranging from 7

<sup>1</sup>The pretrained CAM++ came from <https://github.com/modelscope/3D-Speaker>

<sup>2</sup>The ASR and forced-alignment models came from <https://github.com/modelscope/FunASR>

<sup>3</sup>The pretrained RFB-Net came from <https://github.com/Linzaer/Ultra-Light-Fast-Generic-Face-Detector-1MB>

<sup>4</sup>The pretrained TalkNet came from <https://github.com/TaoRuijie/TalkNet-ASD>

<sup>5</sup>The pretrained CurricularFace model came from [https://modelscope.cn/models/iic/cv\\_ir101\\_facerecognition\\_cfglint](https://modelscope.cn/models/iic/cv_ir101_facerecognition_cfglint)

Table 4: Constraints derived from various modalities. We separately evaluate the accuracy and coverage of these constraints.

Constraints	Accuracy(%)			Coverage(%)		
	Must-Link	Cannot-Link	Total	Must-Link	Cannot-Link	Total
Semantic Constraints	99.75	84.80	99.40	1.23	0.08	0.49
Visual Constraints	99.07	97.87	99.32	22.81	21.78	22.53
Semantic + Visual Constraints	99.11	97.83	99.34	23.65	21.84	22.87

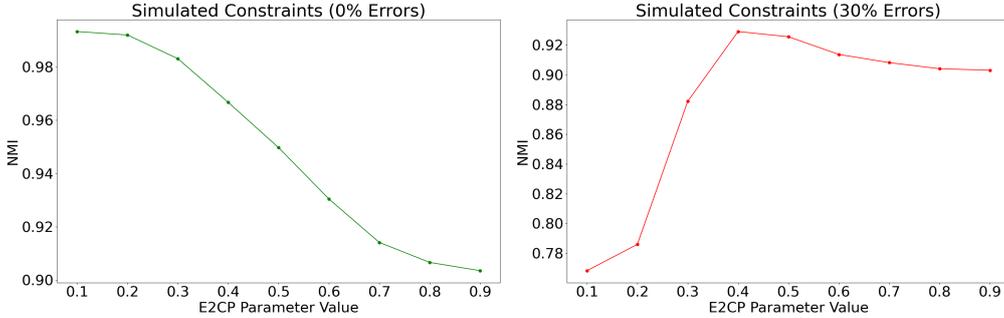


Figure 5: Analysis of constrained clustering outcomes with varying  $\lambda$  values. It is observed that when constructed constraints contain errors, the peak of the optimal  $\lambda$  shifts towards 1.0.

1059 training datasets to obtain Mandarin semantic mod- 1089  
 1060 els. In our experiments, a sliding window strategy 1090  
 1061 was employed, featuring a window size of 96 words 1091  
 1062 and a shift of 16 words, to construct training sets 1092  
 1063 for dialogue detection and speaker-turn detection 1093  
 1064 training from transcripts with speaker annotations 1094  
 1065 within these datasets. All that training was con- 1095  
 1066 ducted using a pre-trained BERT model (Devlin 1096  
 1067 et al., 2019). Subsequently, we employ the methods 1097  
 1068 described in Section 3.3 to construct the semantic 1098  
 1069 constraints. 1099

1070 The VBx approach (Landini et al., 2022) is a 1100  
 1071 canonical method for speaker diarization, where the 1101  
 1072 original paper utilizes speaker embeddings based 1102  
 1073 on the x-vector model. We replace this with the 1103  
 1074 more robust CAM++ model. Additionally, since 1104  
 1075 the post-processing step of the E2CP method in- 1105  
 1076 corporates SC (Von Luxburg, 2007), we also in- 1106  
 1077 vestigate the performance of a method that relies 1107  
 1078 solely on speaker embeddings and SC. These two 1108  
 1079 audio-only methods will serve as the baselines for 1109  
 1080 this study. 1110

1081 As introduced in Section 3.1, after obtaining mul- 1111  
 1082 timodal pairwise constraints, our clustering process 1112  
 1083 is divided into two submodules: constraint propa- 1113  
 1084 gation and post-clustering. When only visual 1114  
 1085 constraints are present, the parameter  $\lambda$  in E2CP 1115  
 1086 is set to 0.8, while it is set to 0.95 when semantic  
 1087 constraints are incorporated. In the post-clustering  
 1088 phase, we adhere to the Spectral Clustering (SC)

algorithm, consistent with the baseline. Inspired  
 by the work presented in (Park et al., 2020), our  
 method incorporates refinement operations such  
 as row-wise thresholding and symmetrization to  
 enhance the performance of spectral clustering.  
 For the row-wise thresholding step in SC, the p-  
 percentile parameter is set to 0.982.

### C Constraints Statistics

Table 4 illustrates the statistical information of the  
 visual constraints and semantic constraints con-  
 structed on our proposed dataset. Upon analysis, it  
 is evident that visual constraints significantly out-  
 perform semantic constraints in terms of coverage.  
 This disparity is attributed to the fact that the two  
 semantic tasks employed in our semantic model  
 are only capable of evaluating the relationships be-  
 tween embeddings within adjacent speaker turns,  
 whereas visual constraints are assessed across em-  
 bedding pairs with substantial temporal intervals.  
 Furthermore, the method we designed to combine  
 constraints from different modalities successfully  
 merges them.

### D Constrained Cluster Parameters Analysis

As mentioned in Section 3.1,  $\lambda$  is a critical param-  
 eter during the constraint propagation process. By  
 combining the analysis of Equations 5 and 6, it can

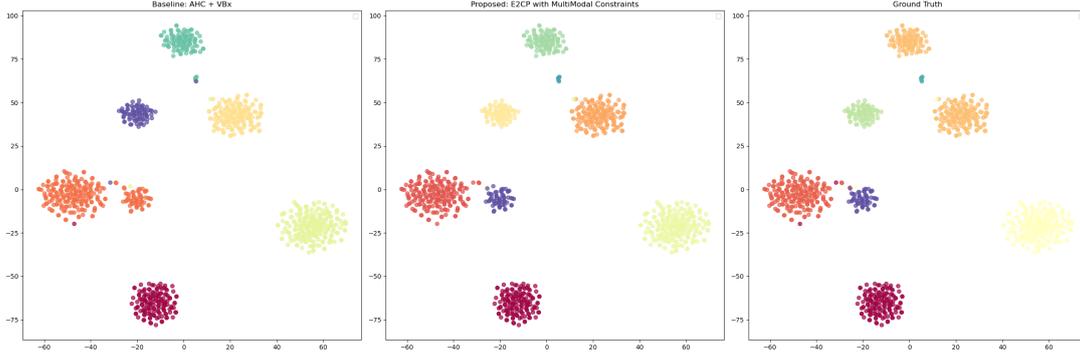


Figure 6: The t-SNE for cluster cases

1116 be found that when  $\lambda$  tends towards 0, the final  $\hat{\mathcal{Z}}$  1153  
 1117 will be closer to  $\mathcal{Z}$ , whereas when  $\lambda$  approaches 1, 1154  
 1118 the resulting  $\hat{\mathcal{A}}$  will be closer to  $\mathcal{A}$ . 1155

1119 Moreover, the parameter  $\lambda$  also signifies the level 1156  
 1120 of confidence that the model places in the con- 1157  
 1121 straints matrix. By adjusting the  $\lambda$  value, the model 1158  
 1122 can effectively handle different levels of error in the 1159  
 1123 constraints, enabling the constrained propagation 1160  
 1124 algorithm to adapt to models with varying perfor- 1161  
 1125 mance. This adaptability is essential for effectively 1162  
 1126 utilizing constraints in real-world scenarios. 1163

1127 We conducted simulations of constraints to compare 1164  
 1128 the optimal  $\lambda$  values when introducing errors 1165  
 1129 in the constraints. The Figure 5 illustrate that the 1166  
 1130 optimal E2CP parameter value  $\lambda$  for maximizing 1167  
 1131 NMI depends on the error rate within the con- 1168  
 1132 straints. With 0% errors, the best performance 1169  
 1133 is achieved at the lowest  $\lambda = 0.1$ , indicating that 1153  
 1134 with highly accurate constraints, the algorithm ben- 1154  
 1135 efits from a strong adherence to constraint guid- 1155  
 1136 ance. However, for constraints with a 30% error 1156  
 1137 rate, the peak NMI occurs at a higher  $\lambda = 0.4$ , 1157  
 1138 suggesting that with less reliable constraints, the 1158  
 1139 algorithm requires a more moderate constraint in- 1159  
 1140 fluence to balance error tolerance and performance. 1160  
 1141 These results highlight the importance of adjusting 1161  
 1142  $\lambda$  in accordance with the fidelity of constraints to 1162  
 1143 achieve optimal speaker diarization. 1163

## 1144 E Decoding cases and Cluster 1153 1145 Visualization 1154

1146 We utilized the t-SNE (van der Maaten and Hinton, 1153  
 1147 2008) algorithm to demonstrate the results of our 1154  
 1148 clustering method, as shown in Figure 6. We com- 1155  
 1149 pared the results of VBx, E2CP with semantic and 1156  
 1150 visual constraints, and ground-truth, and observed 1157  
 1151 that the E2CP with semantic and visual constraints 1158  
 1152 method effectively improved the clustering results 1159

1153 compared to VBx, especially in terms of clustering 1154  
 1155 the points at the edges of clusters, after introducing 1156  
 1157 constraints. 1158

1159 In Figure 7, we present a decoding case, where 1160  
 1161 each row follows the format: speaker-ID, start- 1161  
 1162 time, end-time and content. For the convenience 1162  
 1163 of aligning timestamps with textual information, 1163  
 1164 the decoding results presented here do not include 1164  
 1165 punctuation marks such as periods. In both the 1165  
 1166 acoustic-only and multimodal results, the same 1166  
 1167 force-alignment results was applied, resulting in 1167  
 1168 identical timestamp values. It can be observed that 1168  
 1169 there is a clear semantic transition point near 111.3 1169  
 1170 seconds. The acoustic-only result fails to correctly 1153  
 1171 segment this speaker change point; however, by 1154  
 1172 leveraging semantic constraints, our method suc- 1155  
 1173 cessfully separates the speakers. 1156

<b>Acoustic Results:</b>	<b>MultiModal Results:</b>	<b>Ground Truth Results:</b>
0 106.67 109.06 My baking experience is limited to eating raw cookie dough	3 106.67 109.06 My baking experience is limited to eating raw cookie dough	c1 105.894 111.193 My baking experience is limited to eating rock cookie dough,
0 109.07 110.96 and the last time we made one of these videos	3 109.07 110.96 and the last time we made one of these videos	and the last time we made one of these videos.
1 111.3 112.82 last time on without a recipe	5 111.3 112.82 last time on without a recipe	c2 111.193 120.328 Last time on without a recipe, my bread looked amazing and tasted
1 112.83 116.26 my bread looked amazing and tasted very plain	5 112.83 116.26 my bread looked amazing and tasted very plain	very plain, so this time I'm not gonna forget the flavor.
1 116.44 119.92 so this time I'm not going to forget the foot	5 116.44 119.92 so this time I'm not going to forget the foot	c4 120.328 130.223 I like to pride myself on the fact that I'm good at making food for others.
1 120.32 124.31 I like to pride myself on the fact that I'm good at making food for others	0 120.32 124.31 I like to pride myself on the fact that I'm good at making food for others	

Figure 7: Decoding case