

# Survey of End-to-End Multi-Speaker Automatic Speech Recognition for Monaural Audio

Xinlu He and Jacob Whitehill

**Abstract**—Monaural multi-speaker automatic speech recognition (ASR) remains challenging due to data scarcity and the intrinsic difficulty of recognizing and attributing words to individual speakers, particularly in overlapping speech. Recent advances have driven the shift from cascade systems to end-to-end (E2E) architectures, which reduce error propagation and better exploit the synergy between speech content and speaker identity. Despite rapid progress in E2E multi-speaker ASR, the field lacks a comprehensive review of recent developments. This survey provides a systematic taxonomy of E2E neural approaches for multi-speaker ASR, highlighting recent advances and comparative analysis. Specifically, we analyze: (1) architectural paradigms (SIMO vs. SISO) for pre-segmented audio, analyzing their distinct characteristics and trade-offs; (2) recent architectural and algorithmic improvements based on these two paradigms; (3) extensions to long-form speech, including segmentation strategy and speaker-consistent hypothesis stitching. Further, we (4) evaluate and compare methods across standard benchmarks. We conclude with a discussion of open challenges and future research directions towards building robust and scalable multi-speaker ASR.

**Index Terms**—Multi-speaker ASR, end-to-end ASR, monaural audio, speech overlap

## I. INTRODUCTION

MULTI-SPEAKER Automatic Speech Recognition (ASR) aims to transcribe speech from an audio containing multiple speakers whose speech may be overlapping. In contrast to single-speaker ASR, which focuses on *what was said*, multi-speaker ASR additionally determines *who says what* [1], [2], [3]. This task is closely related to the well-known *cocktail party problem* [4], where humans focus on one speaker in a noisy environment filled with competing talkers. Multi-speaker ASR extends this concept by transcribing all speakers in the mixture, and attributing words to individual speakers. This enables practical applications in real-world scenarios such as meetings, group discussions, and phone call recordings, and supports diverse downstream tasks such as meeting summaries, dialogue analytics, and intelligent conversational assistants.

Compared to single-speaker ASR, multi-speaker ASR poses unique challenges, primarily due to overlapping speech and the need for speaker distinction. These require advanced modeling and are further constrained by the scarcity of large, well-annotated datasets. Additionally, multi-speaker ASR is inherently multifaceted, involving not only recognition and diarization but also overlap detection, turn-taking detection, and target-speaker ASR. Although these tasks are interrelated

X. He is with the Department of Data Science, Worcester Polytechnic Institute, Worcester, MA, 01609 USA e-mail: xhe4@wpi.edu.

J. Whitehill is with Worcester Polytechnic Institute, Worcester, MA, 01609 USA e-mail: jrwhitehill@wpi.edu.

and can benefit from joint modeling, effectively leveraging their synergy remains challenging.

While there are multiple comprehensive literature surveys for single-speaker ASR [5], [6], [7], there has been no recent such survey for end-to-end multi-speaker ASR systems. Our paper seeks to fill this gap. Before exploring end-to-end solutions, we examine the limitations of traditional cascade architectures, which served as the initial attempts to address the multi-speaker ASR task.

### A. The Limits of Cascade Methods

Early multi-speaker ASR systems often adopted cascade (modular) methods [8], [9], as shown in Fig. 1. One approach is **diarization-segmented cascade system** (Fig. 1(a)): (1) Apply speaker diarization to determine *who spoke when* to obtain time-stamped boundaries for each speaker. (2) Split the audio into segments by detected speaker boundaries. (3) Use a single-speaker ASR model on each segment to obtain individual transcription. This approach leverages well-established single-speaker ASR and works well under minimal overlap. However, its accuracy degrades with overlapping speech. Moreover, traditional diarization [10], [11], [12], [13] assumes a single active speaker per frame. Even with later diarization methods that support multi-speaker labeling per frame [14], [15], [16], [17], the resulting segments may still contain overlapping speech, posing challenges for single-speaker ASR models.

Another **separation-based cascade system** (Fig. 1(b)) incorporates a speech separation module. (1) A separation model *enhances* and *denoises* the mixed audio into multiple single-speaker streams, addressing overlapping at the signal level. (2) Each stream is then transcribed using a single-speaker ASR model. While this method can separate overlapping speech in the initial stage, its overall accuracy is dependent on the separation models, which typically optimize signal-level objectives rather than directly targeting ASR performance. Consequently, errors introduced during the speech separation phase can propagate to the subsequent ASR process, compounding the overall error rate of the system [8]. These limitations have motivated end-to-end (E2E) multi-speaker ASR approaches that directly map mixed audio to speaker-attributed transcriptions.

### B. Preview of End-to-End Methods

Unlike cascade methods that explicitly separate speakers before transcription, E2E systems model *jointly* optimize the complementary problems “*who is speaking*” and “*what is being said*”. This can improve accuracy on both tasks. In E2E

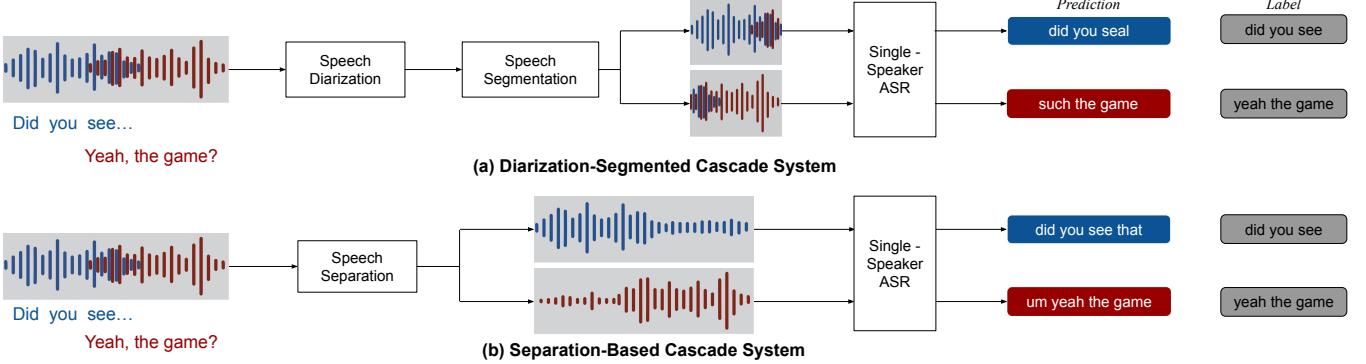


Fig. 1: **Two Types of Cascade Multi-Speaker ASR Systems.** (a) **Diarization-Segmented Cascade System:** The mixed audio is segmented by speaker and then processed by a single-speaker ASR model, potentially introducing errors in overlapping speech regions. (b) **Separation-Based Cascade System:** The mixture is first separated into individual speaker streams using speech separation, followed by single-speaker ASR processing. This method may propagate errors from the separation stage.

frameworks, the input is the raw mixed audio, and the output is the transcriptions partitioned by different speakers.

Initial explorations on E2E multi-speaker ASR [18], [19], [2], [20] primarily adopted a **single-input multiple-output** (SIMO) design, where mixed speech is processed through parallel branches to extract speaker-specific representations. These methods typically follow a separation-then-recognition process and are trained end-to-end, often using *permutation invariant training* (PIT) [21], [22]. A key limitation is that they assume a fixed number of speakers. To address this limitation, **single-input single-output** (SISO) methods were proposed [23], [24], [25], [26], notably using *serialized output training* (SOT) [27], which generates a unified sequence across speakers. Both SIMO and SISO have since been extended with various architectural and training improvements.

In light of recent advances, this review consolidates recent progress in end-to-end multi-speaker ASR by providing a structured taxonomy of representative models. We analyze core designs and improvements, and compare performance across benchmarks. Three key observations are summarized:

- 1) A key distinction in current multi-speaker ASR research lies in whether to separate the speech mixture explicitly. Explicit separation generates multiple outputs (SIMO), offering clearer modularity and easier integration with separation and ASR. In contrast, direct mixture processing produces a single output (SISO), preserving contextual information across stages and enabling multi-task learning with mixture-based tasks.
- 2) There is a growing trend to adapt foundation speech models [28], [29], [30] to multi-speaker scenarios via lightweight structure modifications and fine-tuning, aiming to mitigate data scarcity.
- 3) Model comparison is hindered by limited open-source implementations and inconsistent settings. We provide setting-wise comparisons on standard datasets (AMI, LibriSpeechMix, LibriMix), and analyze models on their design focus.

### C. Review Scope and Organization

In this paper, we review the recent progress on end-to-end multi-speaker ASR where multiple speakers may talk

simultaneously. In particular, we focus on monaural (single-channel) audio recordings and offline (not real-time/streaming) speech analysis. To structure the comparison of recent models, we identify three key dimensions: (1) **Model architecture:** whether the system follows a SIMO or SISO framework. (2) **Input granularity:** whether the system is designed and evaluated on pre-segmented clips or long-form continuous audio. (3) **Speaker enrollment:** whether the system incorporates speaker enrollment to improve the system.

The rest of the paper is organized as follows:

- Section II reviews background techniques for multi-speaker ASR: end-to-end ASR and speech separation techniques.
- Section III categorizes and reviews recent advances in SISO and SIMO models for pre-segmented audio.
- Section IV focuses on long-form audio, discussing segmentation techniques and hypothesis stitching.
- Section V reviews datasets and evaluation metrics for multi-speaker ASR, and presents performance comparisons across different experimental settings.

## II. BACKGROUND TECHNIQUES

This section reviews two background techniques for multi-speaker ASR. First, we discuss E2E ASR architectures, detailing prevalent model designs and objectives. Second, we outline the speech separation techniques that process mixed audio into isolated streams, which inspire architecture or serve as pre-processing modules for multi-speaker ASR.

### A. End-to-End ASR

Recent advances in deep learning have enabled E2E network approaches to achieve state-of-the-art performance in ASR. These systems utilize sequence-to-sequence models to directly map speech signals to text outputs. The three primary end-to-end ASR architectures are: (1) Connectionist Temporal Classification (CTC) which relies solely on the previous input; (2) Recurrent Neural Network Transducer (RNN-T), which depends on both previous input and previous output; and (3) Attention-based Encoder-Decoder (AED) architectures, which consider all inputs and previous outputs.

Among these, AED has gained wide appeal through the development of state-of-the-art systems such as Whisper [28]. In AED, the encoder processes the input acoustic features into a sequence of hidden states, and the decoder predicts output tokens conditioned on past outputs via an attention mechanism over encoder representations. Early AED implementations relied on RNNs for both the encoder and decoder, but recent models increasingly adopt Transformer [31] and Conformer [32] architectures for their ability to model long-range dependencies with global attention. Recently, large speech foundation models (e.g., Whisper [28], Wav2Vec [29], Hubert [30]) leverage self-supervised learning or multi-task training on massive datasets. After fine-tuning for ASR, these models achieve state-of-the-art performance.

The training objective typically employs a cross-entropy loss to maximize the probability of transcription. Additionally, a joint CTC loss is often applied to the encoder outputs to enforce monotonic alignment, leveraging its dependence only on previous inputs. This also enhances the encoder's acoustic modeling. However, due to CTC's strict monotonic constraint, it struggles with overlapping speech, requiring careful adaptation in multi-speaker ASR.

### B. Speech Separation Techniques

In end-to-end multi-speaker ASR, it sometimes includes a speech separation module and is trained together with the ASR part, inspired by the cascade model. Here we look back at the deep-learning-based speech separation techniques.

The deep-learning-based speech separation can be divided into frequency-domain methods and time-domain methods. The frequency-domain methods process the frequency feature of mixed speech and estimate time-frequency masks or spectral magnitudes for each speaker. The deep clustering framework (DPCL) [33] first maps each time-frequency spectral magnitude into a speaker-discriminative embedding, and then the clustering algorithm is used to get the speaker label. Alternatively, a mask output for speech separation can be directly estimated by a deep neural network without embedding. Chimera [34] combines the two, outputting both speaker embedding and mask. The time domain methods, such as TasNet [35], directly consume the waveforms using the encoder-separator-decoder framework. The encoder decomposes the mixture into learnable basis functions, the separator estimates speaker-specific weights, and the decoder reconstructs clean waveforms.

In training, label ambiguity arises when multiple outputs must be matched to multiple labels without a predefined order. Permutation Invariant Training (PIT) [36] addresses this by dynamically aligning model outputs with reference signals. It evaluates all possible permutations and selects the one that minimizes the loss. Consider a mixture of speech from  $S$  speakers, and a model produces a set of estimated signals  $\{\hat{Y}^s\}_{s=1}^S$ . The corresponding reference signals  $\{Y^s\}_{s=1}^S$  have no inherent correspondence to the outputs due to the unknown order of speakers. The objective function is defined as:

$$\mathcal{L}_{\text{PIT}} = \min_{\pi \in \mathcal{P}(S)} \sum_{s=1}^S \mathcal{L}(\hat{Y}^s, Y^{\pi(s)}), \quad (1)$$

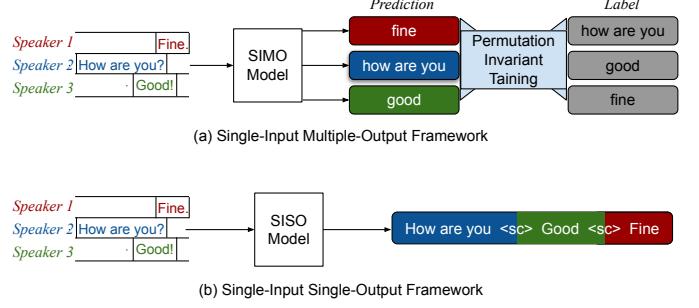


Fig. 2: Single-input multiple-output (SIMO) and single-input single-output (SISO) processes for multi-speaker ASR.

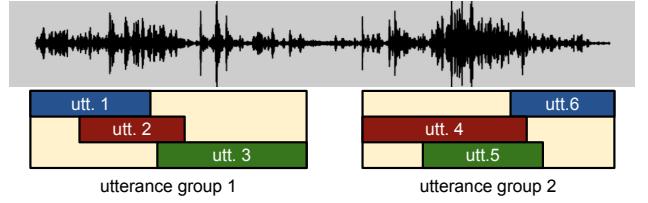


Fig. 3: Example of utterance groups consisting of overlapping speech segments from multiple speakers.

where  $\mathcal{P}(S)$  represents all permutations of  $\{1, \dots, S\}$ , and  $\mathcal{L}(\hat{Y}^s, Y^{\pi(s)})$  computes loss between  $\hat{Y}^s$  and the permuted reference  $Y^{\pi(s)}$ .

Despite its effectiveness, PIT faces computational challenges due to the need to evaluate all possible permutations. Additionally, alternative approaches like HEAT have been developed to reduce computational cost by using the Hungarian algorithm to efficiently find the optimal permutation.

## III. END-TO-END MULTISPEAKER ASR

In this section, we explore end-to-end multi-speaker ASR systems on two prominent frameworks: single-input single-output (SISO) and single-input multi-output (SIMO) (as shown in Fig. 2). While SIMO frameworks generate separate transcriptions for each speaker through parallel branches, SISO frameworks process mixed audio to produce a single transcription output. We review their architectural designs, training methodologies, and key improvements of both frameworks, highlighting their strengths and limitations.

This section focuses on processing predefined audio clips, typically segmented from continuous audio using silence points and heuristic rules (e.g., duration thresholds). In single-speaker ASR, such segments are commonly known as *utterances*. For multi-speaker ASR, the concept of *utterance group* [37] has been introduced as multiple utterances linked through overlapping regions (as shown in Fig. 3).

### A. Single-Input Multiple-Output (SIMO)

SIMO frameworks process mixture audio with exactly  $S$  speakers and produce  $S$  transcriptions, one per speaker, through parallel branches, as illustrated in Fig. 2(a). Here,  $S$  is a fixed parameter denoting the number of speakers. Starting

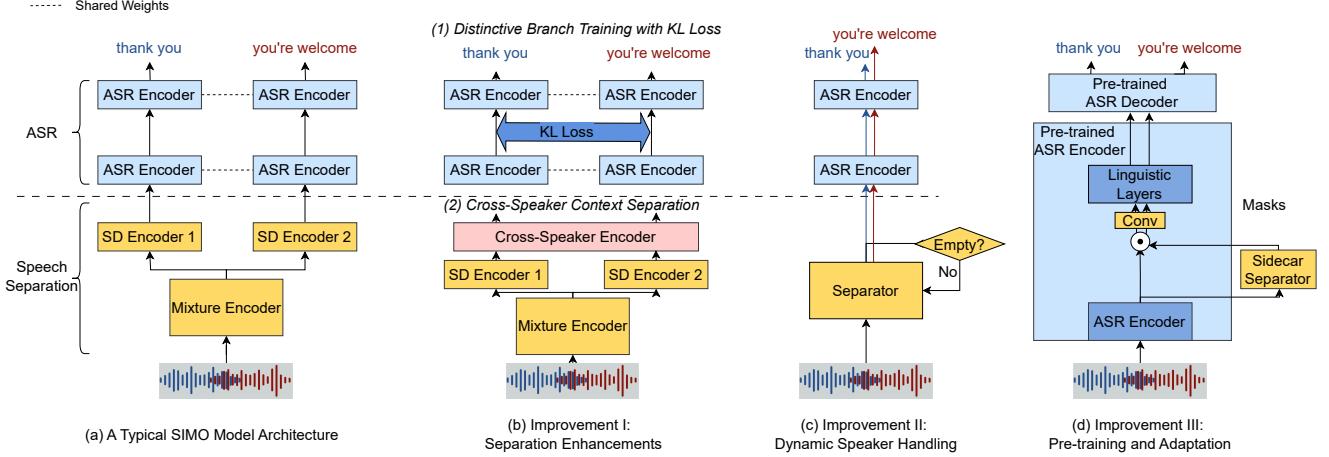


Fig. 4: **A typical SIMO architecture (a) and three types of key improvements (b-d).** (b) Enhance separation by (1) Introducing an auxiliary KL Loss to promote distinctive separation between speakers. (2) Incorporating a Cross-Speaker Encoder to provide cross-speaker contextual cues for compensating omission and reducing repetitions between branches; (c) Support dynamic speaker counts using iterative separation. (d) Adapt pre-trained large speech foundation model with a Sidecar separator.

from the standard architecture, we discuss advancements in separation enhancement, speaker scalability, and leveraging pre-trained models.

### 1) Model Architecture

The SIMO framework generates a separate transcription for each speaker through distinct output branches. A simple implementation involves a shared network followed by  $S$  individual networks, each dedicated to a speaker’s transcription. For instance, [18] employs a shared RNN with  $S$  linear heads to produce  $S$  transcriptions. Alternatively, inspired by separation-based cascade methods, a class of SIMO frameworks (e.g., [19], [38]) was proposed to integrate speech separation and ASR in a unified structure.

Fig. 4a illustrates a typical SIMO model architecture [19], [20], [39], [40] that integrates the speaker separation process and ASR into a single framework, enabling end-to-end training from scratch. The architecture adopts a stacked design comprising shared and unshared modules across branches. The process begins with a mixture encoder,  $\text{Encoder}_{\text{Mix}}$ , which extracts an intermediate feature sequence  $\mathbf{Z}$  from the mixed audio input  $\mathbf{X}$ , serving as the input for subsequent separation and ASR. This feature sequence is then fed into  $S$  parallel branches, each dedicated to one of the  $S$  speakers. Within branch  $s$ , a speaker differentiating encoder,  $\text{Encoder}_{\text{SD}^s}$ , disentangles the speech content  $\mathbf{Z}^s$  of the corresponding speaker from the mixture feature  $\mathbf{Z}$ . Finally, a shared ASR model, such as an attention-based encoder-decoder (AED) or transformer, generates the transcription hypothesis  $\mathbf{H}^s$  for speaker  $s$ . The process can be formally described as:

$$\mathbf{Z} = \text{Encoder}_{\text{Mix}}(\mathbf{X}), \quad (2)$$

$$\mathbf{Z}^s = \text{Encoder}_{\text{SD}^s}(\mathbf{Z}), \quad s = 1, \dots, S \quad (3)$$

$$\mathbf{H}^s = \text{ASR}(\mathbf{Z}^s), \quad s = 1, \dots, S \quad (4)$$

The training objective of each branch follows that of single-speaker ASR: cross-entropy loss  $\mathcal{L}_{\text{CE}}$  to maximize the transcription accuracy, and an optional CTC loss  $\mathcal{L}_{\text{CTC}}$  for

monotonic alignment. Similar to speech separation, SIMO models with multiple output branches introduce label ambiguity, where the correspondence between hypotheses and references is unclear. Permutation Invariant Training (PIT) is also utilized to align hypotheses  $\mathbf{H}^s = (h_1^s, \dots, h_{N_s}^s)$  with reference transcriptions  $\mathbf{R}^s = (r_1^s, \dots, r_{N_{ss}}^s)$  through permutations [18]. To reduce the permutation computational cost, label matching is often determined solely based on the CTC loss, while both CTC and cross-entropy losses are used ([19], [41], [39], [40]).

SIMO provides a natural framework for integrating traditional separation and ASR methods into an end-to-end system. This design facilitates the incorporation of state-of-the-art separation techniques to enhance performance. However, the fixed number of branches limits the model’s ability to handle a variable number of speakers, which remains a significant constraint in real-world scenarios. Also, limited separation performance still causes redundancy and omissions in the final transcription output.

### 2) SIMO Improvements

To address SIMO’s limitations, recent research has focused on three key goals: (1) enhancing separation performance, (2) enabling a flexible number of speakers, and (3) mitigating data scarcity. The following sections detail these categories. Fig. 5 summarizes the corresponding goals and methods.

#### a. Separation Enhancement

While end-to-end models jointly train separation and recognition modules, independent branches can propagate early separation errors, resulting in repeated or omitted transcriptions. Recent work addresses this by increasing inter-branch distinctiveness and leveraging context to refine separation.

**Distinctive Branch Training:** To encourage distinct transcriptions across output branches, auxiliary loss functions can be applied between ASR hidden state vectors. In AED-based

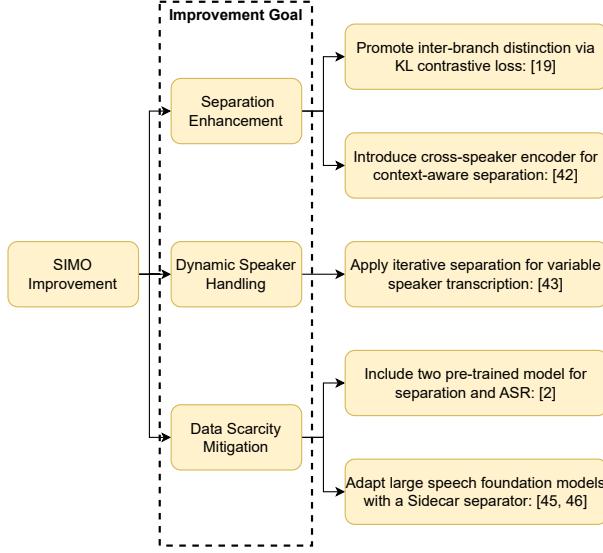


Fig. 5: **SIMO Improvements by Goal.** Targeting better separation performance, variable speaker handling, and data scarcity mitigation by module refinement, iterative processing, and pre-trained model integration.

models, Seki et al. [19] proposed a contrastive loss that maximizes Kullback-Leibler (KL) divergence between the ASR encoder states from different branches, as shown in Fig. 4(b.1).

This loss penalizes similarity between streams, reducing the likelihood of redundant transcriptions.

**Context-Aware Separation:** More recently, the separation has been further improved by introducing cross-speaker context to enhance the separation quality. Traditional SIMO systems process each speaker independently in parallel branches, restricting the model’s capacity to capture cross-speaker dependencies and perform mutual correction. To address this, Kang et al. [42] proposes a *Cross-Speaker Encoder* ( $\text{Encoder}_{\text{CSE}}$ ), positioned between the  $\text{Encoder}_{\text{SD}}$  and ASR model to enable information sharing across branches (Fig. 4(b.2)). It consumes the concatenation of mixture encoding from  $\text{Encoder}_{\text{Mix}}$  and different speaker’s individual encodings from  $\text{Encoder}_{\text{SD}}$ , and employs a conformer block to enable information sharing across branches. The conformer output is then partitioned into branch-specific representations and passed to the respective ASR encoders. This context-aware mechanism leverages the local context of the mixed speech to refine the single-speaker features, thereby improving separation accuracy.

### b. Dynamic Speaker Handling

A key limitation of traditional SIMO frameworks is their reliance on a fixed number of speaker streams, which constrains their applicability to real-world scenarios. Recent work addresses this challenge by introducing iterative and adaptive systems to determine dynamically the number of speakers. As illustrated in Fig. 5(c), Neumann et al. [43] propose an

iterative approach where the system separates one speaker’s audio at a time with a Dual-Path RNN TasNet [44] separator. The residual mixture is then passed to subsequent iterations until no speech remains. Each separated feature is then fed into a shared ASR model, which can be jointly fine-tuned with the separator. This design supports end-to-end training while supporting varying numbers of speakers.

### c. Pretraining and Adaptation

Leveraging pre-trained models mitigates the data scarcity challenge in multi-speaker ASR, as shown in Fig. 4(d). In the SIMO paradigm, a common approach is to adopt a dual-module framework consisting of a pre-trained speech separation module and a pre-trained ASR module, with joint fine-tuning to align their objectives. For example, Settle et al. [2] includes Chimera++, a pre-trained separation module, to generate speaker masks over mixture features, producing  $S$  separate streams. These streams are then fed into a shared ASR model for transcription, and the two modules are subsequently fine-tuned together to enable end-to-end optimization.

More recent advances introduce large speech foundation models like Whisper [28] and Wav2Vec [29] into the SIMO pipeline. Since these models are not inherently designed for multi-speaker inputs, a separation module must be inserted to enable SIMO-style processing. To this end, Meng et al. [45], [46] propose a lightweight convolutional network Sidecar separation module inserted between layers of a frozen Wav2Vec2.0 or Whisper model. This design splits the mixed audio into multiple streams, each processed independently by subsequent Wav2Vec layers to produce distinct transcriptions. Instead of building an external separation module, it operates directly on the internal feature representation of the pre-trained model, enabling an efficient and seamless adaptation of pre-trained single-speaker ASR models to multi-speaker scenarios.

## B. Single-Input Single-Output (SISO)

SISO is another multi-speaker ASR framework that only generates a single output sequence containing all speakers’ transcriptions. Unlike SIMO, which produces separate outputs per speaker, SISO relies on serialized output training (SOT) [27] to serialize transcriptions into a unified stream. The transcriptions of all speakers are concatenated into one output.

SISO offers several advantages. First, it naturally handles varying numbers of speakers, making it well-suited to real-world scenarios with unknown speaker counts. Second, by generating a single output sequence, SISO captures inter-speaker dependencies, improving both coherence and overall transcription accuracy. Third, it reduces computational cost by enforcing a fixed output order, thereby simplifying training.

This section first introduces two forms of SOT output: speaker-ordered sentence-based SOT and temporal-ordered token-based SOT. We then discuss improvements to the SISO architecture.

### 1) Serialized Output

SOT can be implemented in two primary forms: sentence-based SOT and token-based SOT (Fig. 6). In sentence-based

Speaker 1	HELLO HOW ARE YOU	GOOD THANK YOU
Speaker 2	NOT BAD AND YOU	
Speaker 3		FINE
<b>Arbitrary Speaker-Ordered SOT</b>		
NOT BAD AND YOU <sc> HELLO HOW ARE YOU GOOD THANK YOU<sc>FINE		
<b>FIFO Speaker-Ordered SOT</b>		
HELLO HOW ARE YOU GOOD THANK YOU <sc> NOT BAD AND YOU <sc>FINE		
<b>Temporal-Ordered SOT</b>		
HELLO HOW ARE <cc> NOT<cc> YOU <cc> BAD <sc>FINE <cc> AND YOU<cc> GOOD THANK YOU		

Fig. 6: An example of different SOT texts, containing speaker-change symbol  $\langle sc \rangle$  and  $\langle cc \rangle$ .

SOT, all sentences of each speaker are concatenated into a sequence. Special token  $\langle sc \rangle$  is inserted between the transcriptions to indicate changes in the speaker. Consecutive sentences from the same speaker are simply concatenated in the transcript, even if those sentences may be partially or even completely interrupted by another speaker’s sentence. The order of speakers can be arbitrary, with PIT used in loss. Alternatively, sentence-based SOT can follow a first-in, first-out (FIFO) order, where transcriptions are arranged by each speaker’s start time, from the earliest to the latest.

In contrast, token-based SOT serializes transcriptions based on token timestamps. The channel change symbol  $\langle cc \rangle$  is utilized in token-based SOT. This approach provides finer granularity in capturing overlapping speech by preserving the timing order, while remaining compatible with CTC loss, which is well-suited for mostly time-monotonic sequences [47] [48]. However, when the same speaker’s sentence is interleaved with others, a more advanced decoder is needed to effectively model long-range context and maintain coherence within that speaker’s speech.

## 2) SISO Model and Improvements

Since SISO models produce a single output stream, they typically build on single-speaker architectures without requiring explicit speaker separation. Fig. 7(a) shows an encoder-decoder model that has been successfully adapted to SISO-based multi-speaker ASR in [27]. In this setting, the ASR encoder generates hidden representations  $Z^{asr}$ , which are then passed to ASR decoder to obtain serialized transcriptions.

Although using a single processing path enables cross-speaker context modeling, the lack of explicit separation mechanisms in SISO systems makes accurate transcription of overlapping speech challenging. To address this and compensate for limited multi-speaker training data, researchers have proposed four main strategies: (1) auxiliary CTC-related losses, (2) integration of external speaker modules, (3) multi-task learning, and (4) pre-training and adaptation. Figure 8 shows the improvement type.

### a. Auxiliary CTC-Related Loss

As described in Section II-A, CTC loss commonly serves as an auxiliary objective to cross-entropy loss in AED to enhance the

encoder’s acoustic modeling. In the SISO framework, it can be applied similarly, through a parallel branch connected to the encoder output, independent of the attention-based decoder, as illustrated in Fig. 7. In addition, modified CTC variants can incorporate speaker information to improve speaker differentiation in multi-speaker transcription, as described below.

**CTC in SISO:** To leverage CTC for improving acoustic alignment in SISO systems, the CTC loss objective must maintain temporal consistency with the original mixed acoustic features. Typically, token-based SOT serves as the CTC target in this setting. If the model’s final output follows sentence-based SOT, the remaining components (not supervised by CTC) need to handle the utterance-level reordering. For instance, Liang et al. [25] proposed a two-stage CTC approach: the first CTC loss is applied after the initial encoder layers, optimized for token-based SOT to enhance acoustic modeling, while the second CTC loss operates on the full encoder output with sentence-based SOT supervision.

**Speaker-Enhanced CTC:** The second approach modifies CTC to explicitly integrate speaker information. A noticeable example is Speaker-Aware CTC (SACTC) [49] proposed by Kang et al., which formulates a Bayes-risk-based CTC that constrains the encoder to distinguish speaker-specific features at specific time frames, explicitly modeling speaker separation. Another advancement comes from Zheng et al. [50] with Weakening and Enhancing CTC (WECTC) loss, which enhances speaker change token ( $\langle sc \rangle$ ) prediction by adjusting pseudo-label posteriors.

Despite its widespread use in SISO frameworks and extensions with speaker-enhanced modules, CTC remains limited in overlapping speech scenarios, as shown in the results of [49]. It assumes conditional independence between output tokens and enforces a one-token-per-frame constraint, producing a single, linear output sequence. This makes it fundamentally incompatible with simultaneous speech from multiple speakers, often resulting in entangled or incomplete transcriptions. Future research may investigate approaches that relax this constraint while preserving CTC’s alignment ability.

### b. External Speaker Module

One direction for improving the SISO framework is the integration of explicit speaker information into the model architecture. By incorporating speaker-specific cues, the model can enhance its ability to differentiate between speakers and better handle overlapping speech. Current approaches can be broadly categorized into two methodologies: (1) Frame-level Speaker Conditioning (FSC) and (2) Token-level Speaker Conditioning (TSC).

**Frame-level Speaker Conditioning (FSC)** incorporates speaker information by aligning speaker and speech timestamps at the frame level. Modular FSC [9] directly aligns diarization results with SOT transcriptions in a post-processing stage, which is not an end-to-end approach and may lead to error propagation. Later end-to-end FSC methods (right side of Fig. 7(c)) integrate speaker information before the

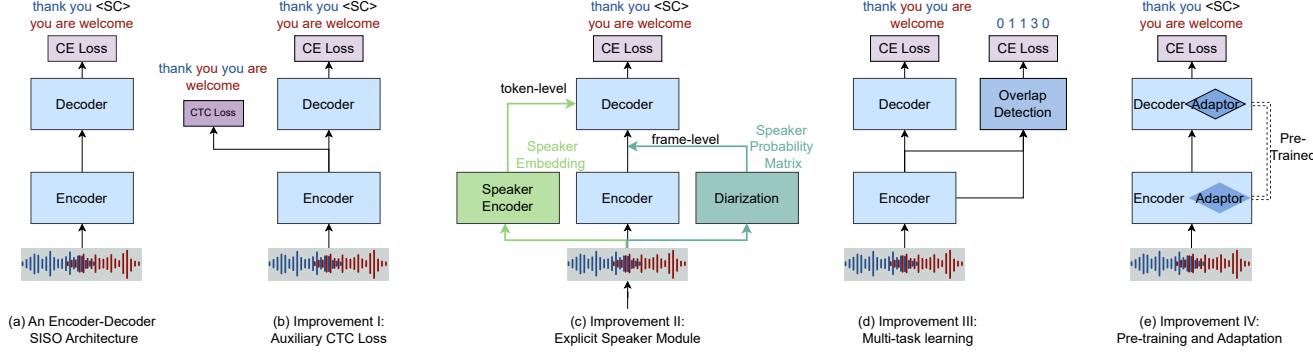


Fig. 7: An **encoder-decoder SISO architecture (a)** and four representative improvements (b-e). (b) Adding an auxiliary CTC loss to enhance acoustic modeling and speaker-awareness capability; (c) Incorporating explicit speaker modules to inject speaker information at the frame or token level; (d) Applying multi-task learning with auxiliary tasks such as overlap detection; (e) Integrating pre-trained models and fine-tuning them with mechanisms such as adapters.

decoder, aligning the diarization output with ASR encoding representation at the frame level. Specifically, the diarization module generates an  $S$ -speaker assignment probability matrix for  $T$  frames, denoted as  $\mathbf{P} \in \mathbb{R}^{S \times T}$ . This matrix is used to incorporate speaker information into the ASR encoder representations  $\mathbf{Z} \in \mathbb{R}^{M \times T}$  before passing the fused representation to the decoder through different integration strategies. One approach [51] incorporates a sinusoidal matrix, similar to the sinusoidal positional encoding used in Transformers [31]. Each speaker is associated with a unique sinusoidal pattern, which is then weighted by the probability matrix  $\mathbf{P}$  and added to the original encoder representation  $\mathbf{Z}$ . This mechanism can differentiate speakers in a structured, deterministic way, and can be disabled to fall back to a standard ASR pipeline. Alternatively, Wang et al. proposed Meta-Cat [52], which first computes a speaker-specific representation  $\mathbf{Z}_s$  by element-wise multiplying the encoder output  $\mathbf{Z}$  with the frame-level probability vector  $P_s$  for each speaker  $s$ . These representations are then concatenated across all speakers to form a supervector, which is passed to the decoder. This method expands the ASR embedding multiple times according to the number of speakers. Both approaches ensure the entire model remains differentiable, enabling joint training of the whole system.

**Token-level Speaker Conditioning (TSC)** introduces speaker embeddings as auxiliary input to the decoder during token generation. Unlike FSC, which aligns frame-level probabilities, TSC typically uses a speaker encoder to extract speaker embeddings for speaker-aware decoding (left side of Fig. 7(c)). A typical TSC speaker-attributed ASR system ([23], [53], [24]) consists of four modules: ASR encoder, speaker encoder, speaker decoder, and ASR decoder. The input mixture is processed in parallel by the ASR and speaker encoders to generate the speech representation  $\mathbf{Z}^{\text{asr}}$  and speaker embeddings  $\mathbf{Z}^{\text{spk}}$ , respectively. These are fed into the decoder, along with the previous token sequence, to produce the context-aware speaker representation. The ASR decoder then takes this speaker representation as extra input with ASR encoder states and previous output to predict the next token.

Here, the speaker representation is related to a pre-defined speaker inventory  $\mathcal{D} = \{\mathbf{d}_1, \dots, \mathbf{d}_K\}$ , where  $\mathbf{d}_k$  represents a speaker profile in the inventory. The output of the speaker decoder is used as a query to compute attention over this speaker inventory  $\mathcal{D}$ , and generates a weighted speaker profile  $\bar{\mathbf{d}}_n$  as final speaker representation given to the speaker decoder. Later approaches proposed by Shi et al. ([54]) considered contextual information for speaker representation generation. A separate contextual text encoder is deployed before the speaker decoder to aggregate the semantic information of the whole output utterance. In addition, when calculating  $\bar{\mathbf{d}}_n$ , an extra context-dependent scorer is employed to model the local speaker discriminability by contrasting with speakers in the context. Fan et al. [48] also enhance the speaker contextual relationship by replacing the weighted profile  $\bar{\mathbf{d}}_n$  as a similarity matrix and passing this matrix to the ASR decoder. In this way, the system can incorporate the speaker information without the speaker inventory  $\mathcal{D}$ .

### c. Multi-Task Learning

Multi-task learning enables the joint optimization of synergistic tasks via a fully or partially shared network, which can in turn improve the performance of individual tasks (e.g., multi-speaker ASR and overlapping detection). Compared to its application in SIMO models, multi-task learning in SISO can directly leverage the shared representation and jointly train the ASR model with mixture-related tasks such as diarization and overlap detection by sharing network components.

**Unified Labeling:** Task-specific labels can be inserted into the serialized output as special tokens. This unified labeling enables joint modeling of multiple tasks, such as speaker identification and timestamp predictions, without modifying the original auto-regressive multi-speaker ASR architecture. First, specific speaker tokens (e.g., `<spk0>` `<spk1>`) can replace general speaker changing tokens (`<sc>` and `<cc>`) to differentiate speakers, as in [55] [56]. This reduces the speaker ambiguity in SOT, especially token-level SOT. Second,

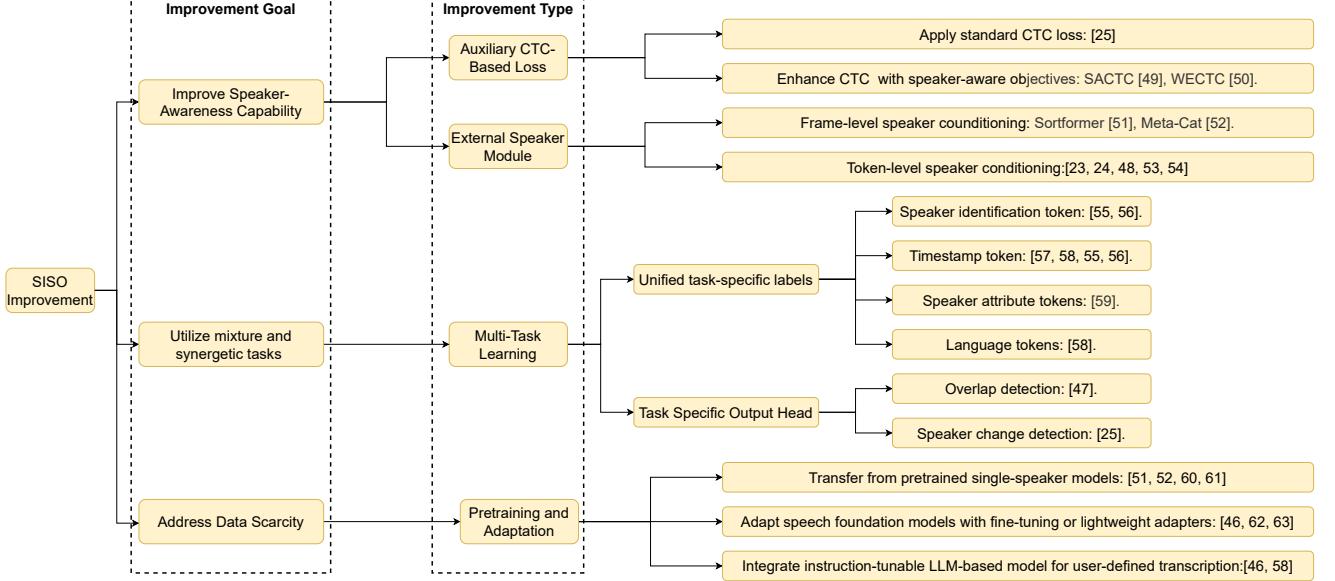


Fig. 8: **SISO Improvements by Goal and Type.** Aiming to enhance speaker-aware capabilities, utilize shared representation and alleviate data scarcity via auxiliary losses, external speaker modules, multi-task learning, and pre-trained models adaptation.

quantized timestamp tokens (e.g., 20ms resolution) can be inserted as special tokens to indicate the start and end time of an utterance [57], [58], [55], [56]. Timestamp tokens provide extra information about speaker turn-taking, temporal order, and overlapping speech. In some systems, additional timestamp tokens are inserted based on heuristic rules, such as when the gap between consecutive tokens exceeds two seconds, to indicate pauses or silence [58]. Additionally, Masumura et al. [59] and Li et al. [58] independently explored the use of speaker attribute tokens (e.g., gender, age) and language tokens, respectively, to enrich the output with contextual information. All tasks under this unified labeling scheme are handled by a single model with a shared output head.

**Task Specific Output Head:** Separate labels and output heads for each auxiliary task can serve as another implementation of multi-task learning in SISO. In this case, tasks share parts of the network, typically lower-layer acoustic or speaker-related representations, while maintaining task-specific output branches in the later stages. Li et al. [47] jointly optimize the overlapping prediction task and the ASR task to enhance multi-speaker ASR (Fig. 7(d)). Specifically, the overlap-aware task predicts two binary states for each token: whether it overlaps with other tokens and whether it is a boundary token. Both the ASR and overlapping prediction models adopt an encoder-decoder structure, sharing lower layers of the encoder to leverage acoustic features while diverging to predict overlap-aware labels. The total loss combines the ASR loss and the overlap-aware loss, enabling joint training to improve performance in overlapping speech scenarios. Speaker change detection has also been integrated as an auxiliary task in multi-task learning frameworks. For example, Liang et al. [25] proposed Boundary-Aware SOT (BA-SOT). Specifically, the

model adds a speaker change detection (SCD) module as the output head for the speaker change task after several decoder blocks. A binary speaker change label is used for the auxiliary task, while the ASR task adopts FIFO-style SOT labels without speaker change tokens. The total loss consists of ASR loss, SCD loss, and a newly introduced boundary constraint loss. Studies [25], [47] have shown that joint modeling of ASR with auxiliary tasks such as speaker change detection and overlap detection can lead to improved recognition accuracy.

#### d. Pretraining and Adaptation

Similar to SIMO, leveraging pre-trained modules in SISO architectures helps mitigate the data scarcity in multi-speaker ASR. Because SISO models share the same model structure and single-output format with single-speaker ASR, SISO models can be directly initialized from pre-trained single-speaker models without architectural modification [60], [61]. This facilitates transfer learning from large-scale simulated mixtures, which are artificially created multi-speaker audio samples (see section V). These simulated datasets can be used for pre-training and subsequently fine-tuned on real-world corpora to achieve competitive performance [37]. Moreover, speaker-specific modules in SISO models – such as speaker encoder and diarization model – can also be initialized from pre-trained modules [51], [52].

This structural compatibility also holds for speech foundation models: single-speaker models, such as Whisper and WavLM-based ASR, can be directly fine-tuned under the SISO framework for multi-speaker ASR, unlike SIMO models that require additional separation modules. These models can either be fully fine-tuned, or adapted using lightweight modules like LoRA, enabling resource-efficient training while keeping most

pre-trained parameters frozen [62], [63], [46]. In addition, Shi et al. [62] explored maintaining the multilingual property of foundation models by leveraging adapters.

Moreover, the unified output format of SISO models naturally aligns with multimodal large language models [64], which typically consist of a speech encoder and a large language model (LLM) performing instruction-driven decoding [58]. Such LLM-based systems offer strong potential for user-defined interaction scenarios. Recent work by Meng et al. [46] exemplifies this interactive paradigm: users can specify multi-speaker ASR behavior through natural language instructions, such as transcribing only the first speaker, a female speaker, or utterances containing specific keywords. The system integrates speaker- and context-aware speech representations from a pre-trained encoder with an LLM, enabling dynamic and controllable transcription tailored to diverse user-defined needs.

### C. SIMO vs. SISO: Summary

By performing speaker separation and speech recognition separately, SIMO provides more modularity than SISO, which performs both tasks in an integrated way. However, SIMO is less flexible than SISO in three key aspects: (1) SISO accommodates variable speaker counts, whereas SIMO requires knowing the number of speakers *a priori*; (2) SIMO’s branch-specific separation lacks cross-speaker validation for redundancy and complementarity and further underutilizes cross-speaker information in ASR. Very few (e.g., [42]) SIMO approaches attempt to address this shortcoming. (3) SISO’s retention of mixed features enables multi-task learning with mixture-based tasks, such as overlapping detection. Additionally, SISO’s fixed speaker order in transcription (e.g, FIFO) reduces computational overhead from label matching in SIMO training.

These structural differences lead to distinct improvement priorities. SIMO’s overall recognition performance is fundamentally limited by its separation module: poor separation introduces residual noise or truncated phonemes, degrading ASR accuracy. Thus, SIMO requires enhanced separation for optimal performance. In contrast, SISO lacks explicit separation, and thus harnessing speaker information to disambiguate overlapping speech is a focus for SISO improvement.

Both SIMO and SISO address data scarcity by utilizing pre-trained models, particularly recent large foundation models. For SIMO, this requires adding a separation process to the foundation model. SISO directly fine-tune the foundation ASR model on multi-speaker data. Both frameworks can achieve strong performance by tuning only 8-10% of all parameters.

### D. Future Directions

One promising direction is **SIMO-SISO hybridization**, which aims to combine SIMO’s separation precision with SISO’s comprehensive modeling. In such framework, the model can first separate the input audio into multiple speaker-specific branches, whose representations are then concatenated and passed to additional network layers for integrated decoding. Recent work has explored this hybridization at

different representation levels. For example, Cross-speaker encoder [42] combined acoustic features from different branches before ASR, to enhance separation by providing cross-speaker context. Also, Huang et al. [65] adapt WavLM to extract target-speaker transcriptions from mixture audio using speaker embeddings. The resulting utterances are concatenated and passed through a joint speaker module to generate the final serialized transcript. Future studies may further explore multi-level integration with advanced model architecture to enhance SIMO’s contextual modeling and improve SISO’s ability to handle overlapping speech.

Another direction is the **adapted foundation model enhancement**. While current work has explored such adaptation within SIMO and SISO frameworks, their combination with existing improvement strategies remains underexplored. For instance, can a foundation-adapted SIMO model also leverage cross-speaker context? How can foundation-adapted SISO models benefit from multi-task learning? Beyond architecture-specific adaptations, future work may further investigate the integration of multi-speaker ASR into multi-modal LLM-based unified frameworks. This can extend current instruction-based approaches [66] to more practical, real-world applications.

## IV. LONG-FORM MULTISPEAKER ASR

The previous section focuses on pre-segmented audio. We turn to continuous long-form audio in this section. Unlike pre-segmented data, long-form audio must be partitioned before being processed by either SIMO or SISO ASR models, and the results must be integrated into a coherent global transcription with consistent speaker identities. This section addresses these two key challenges. An ideal segmentation algorithm should be computationally efficient while preserving linguistic coherence. The segmentation methods are introduced in Section IV-A. Section IV-B discusses how to merge local hypotheses into a unified global transcription, with consistent speaker identities across all segments.

### A. Segmentation Methods

Segmenting long-form multi-speaker audio is crucial for enabling accurate and efficient ASR. Strategies need to balance computational efficiency and preserve linguistic coherence. Existing methods typically rely on acoustic or semantic cues.

Voice activity detection (VAD) segments speech by detecting silent intervals based on acoustic features such as energy levels and spectral patterns. Its simplicity leads to widespread use in both cascade and end-to-end systems [8], [53], [9]. However, VAD-based segments may still be excessively long, necessitating further clipping to meet the input requirements of downstream ASR models. Also, since VAD operates purely on acoustic cues, it may inadvertently split sentences at unnatural boundaries, particularly when speakers pause or hesitate mid-sentence, which can disrupt linguistic coherence and degrade the performance of subsequent ASR tasks.

Sliding window segmentation is another strategy, which processes audio using fixed-length windows and a pre-defined stride. Unlike VAD, the sliding window method generates uniformly sized segments ready for ASR models. However,

it still risks disrupting semantic coherence by splitting sentences at arbitrary boundaries, a problem exacerbated by rigid window constraints. Additionally, shorter strides can inflate computational costs due to extensive overlap.

To address these limitations, recent studies incorporate semantic information into segmentation. For example, Cornell et al. [55] introduce an adaptive sliding window strategy, which inserts the special token `<trunc>` to mark truncation, and resumes from the last silence point to preserve linguistic continuity and reduce redundancy. Huang et al. [67] predict segment boundaries in a streaming manner, based on both acoustic and text-level cues, enabling dynamic segmentation with minimal overhead.

### B. Hypothesis Stitching Methods

Generating the final global hypothesis requires concatenating and aligning local transcriptions from segmented audio. A straightforward approach is to concatenate local transcriptions directly with VAD-segmented clips [68], [53], [55]. When long audio is segmented with overlaps, the final global transcription cannot be obtained by simple concatenation due to redundancy and potential mismatch. High-confidence word selection can be employed in overlapping parts [69], while a neural-network-based hypothesis stitcher [70] fuses segment outputs into a coherent transcription without requiring alignment.

Maintaining globally consistent speaker labels is essential for long-form multi-speaker ASR. Existing approaches rely either on speaker profiles or on embeddings learned jointly within the ASR model. In profile-based methods, speaker identities can be resolved during speaker-attributed decoding when speaker enrollment is available. When enrollment is not possible, speaker profiles can be approximated through unsupervised clustering of embeddings from a pre-trained speaker encoder [68]. Even supplying dummy profiles, which do not appear in the input audio, can improve performance [53]. In the joint modeling approach, the multi-speaker ASR system outputs both transcriptions and speaker embeddings, typically by a multi-task learning framework. The global labels can be obtained by clustering the embeddings. In [55], each window of the E2E DAST model provides local speaker transcription and diarization with time-averaged speaker embeddings. Meanwhile, [71] demonstrates the advantages of jointly learning speaker embeddings and transcriptions for hour-long multi-speaker podcasts, leveraging lexical cues to improve speaker label assignment. After processing all windows, the final diarization and speaker embeddings can be obtained by clustering these time-averaged embeddings. This method eliminates the need for additional speaker embedding models and improves computational efficiency, but its performance depends heavily on the quality of the jointly learned embeddings.

## V. EVALUATION

Multi-speaker ASR research relies on both real-world and simulated datasets. Real-world datasets are collected from natural conversation scenarios such as meetings or phone calls; they provide authentic acoustic conditions, spontaneous speech patterns, and natural overlaps. However, they are often

small, noisy and domain-specific, posing challenges for early-stage model training and broader applicability. Simulated datasets are created by overlapping single-speaker recordings and introducing noise at the configurable overlapping rate and noise level, enabling scalable as well as controlled training and evaluation of overlap and noise robustness. To improve realism, Yang et al. [72] leverages statistical language models to guide overlap patterns. Additionally, Moriya et al. [73] introduce on-the-fly data generation to support dynamic parameter adjustment and improve memory efficiency. Table I highlights commonly used multi-speaker ASR datasets.

TABLE I: *Common real-world and simulated datasets. # Speakers refers to the number of speakers per recording in real datasets, and the exact mixed speaker number per sample in simulated datasets.*

Dataset	Scenario	Hours	Language	# Speakers
<i>Real</i>	AMI	100	English	3–5
	AliMeeting [74]	120	Chinese	2–4
	CallHome	60	Multilingual	2
	LibriCSS [75]	10	English	8
<i>Sim</i>	WSJ0-2mix [33]	45	English	2
	LibriMix [76]	500	English	2–3
	LibriHeavyMix [77]	20,000	English	2–4

Evaluating multi-speaker ASR systems requires measuring both transcription accuracy and the system’s ability to assign speaker labels. **Word error rate (WER)** evaluates overall transcription accuracy by measuring insertions, deletions, and substitutions. **Character Error Rate (CER)** is adopted for languages with non-alphabetic scripts, while **Sentence Error Rate (SER)** captures the proportion of sentences containing any error. In multi-speaker settings, these metrics can be extended with special tokens for speaker changes, overlapping speech, timestamps, and speaker identities. This extension enables evaluation not only for content accuracy but also of the system’s capability to handle speaker turns and overlaps.

In addition to standard WER, several evaluation metrics have been proposed to account for speaker distinctions in multi-speaker ASR. **Concatenated minimum-permutation WER (cpWER)** [78] concatenates utterances per speaker and computes the WER across all permutations of hypotheses and references, selecting the lowest:

$$\text{cpWER} = \min_{\pi \in \mathcal{P}(S)} \frac{\sum_{s=1}^S \text{WER}(H^s, R^{\pi(s)})}{S}, \quad (5)$$

where  $\mathcal{P}(S)$  denotes all permutations of  $\{1, \dots, S\}$ , and  $H^s$ ,  $R^{\pi(s)}$  denote the hypothesis and reference for speaker  $s$ , respectively. cpWER jointly reflects transcription and speaker assignment accuracy. **Speaker-attributed WER (SA-WER)** further enforces speaker correspondence by requiring that hypotheses be matched with the reference of the correct speaker label. It penalizes speaker assignment errors and provides a stricter evaluation of speaker-aware performance. Among these, WER, cpWER, and SA-WER are increasingly strict in evaluation criteria. To incorporate temporal alignment, Neumann et al. [79] propose **time-constrained cpWER (tcpWER)**, which restricts word matching to a fixed time window in addition to speaker permutation, which requires real or estimated token timestamps.

TABLE II: *cpWER of multi-speaker ASR models on AMI, LibrispeechMix, and LibriMix corpus. SDM denotes single distant microphone; IHM denotes mixture of independent headset microphones. Mix Pre. Hrs indicates the hours of simulated mixture data used for pretraining. Params Tr/To indicates trainable / total model parameters in millions. # Tr. Spk refers to the speaker configuration used during training. Methods are grouped by input granularity and listed chronologically. \* and † signify the paper reported results as standard WER or SA-WER, respectively.*

Dataset	Model	Input Gran.	Framework	Spk. Enrl.	Mix Pre. Hrs	Params Tr/To	SDM dev	SDM eval	IHM dev	IHM eval
AMI	Conformer AED [37]	Utterance Group	SISO	✗	900k	50 / 50	18.4	21.2	13.5	14.9
	WavLM/wTSE&JSM [65]	Utterance Group	Hybrid	✓	58	13 / 108	—	—	—	28.4†
	Adapted USM [58]	Utterance Group	SISO	✗	—	84 / 1630	—	21.4	—	—
	META-CAT [52]	Utterance Group	SISO	✗	—	600 / 723	—	—	—	22.8
	Transcribe-to-Diarize [68]	Long-Term	SISO	✗	—	146 / 146	22.6	24.9	15.9	16.4
Libri-speech-Mix	SLIDAR [55]	Long-Term	SISO	✗	5k	655 / 655	21.8	24.5	14.2	15.6
	LSTM SOT [27]	Utterance Group	SISO	✗	1, 2, 3	136 / 136	—	11.2*	—	24.0*
	Transformer SOT [24]	Utterance Group	SISO	✗	1, 2, 3	129 / 129	—	4.9*	—	6.2*
	LSTM SA-ASR [23]	Utterance Group	SISO	✓	1, 2, 3	146 / 146	—	9.9†	—	23.1†
	SA-MBR [80]	Utterance Group	SISO	✓	1, 2, 3	146 / 146	—	9.5†	—	20.7†
	Transformer SA-ASR [24]	Utterance Group	SISO	✓	1, 2, 3	142 / 142	—	6.4†	—	8.5†
	W2V-Sidecar [45]	Utterance Group	SIMO	✗	2	9 / 104	6.0	5.7	—	—
	CSE Network [42]	Utterance Group	SIMO	✗	1, 2	33 / 33	11.8	10.7	24.2	24.3
	CIF SA-SOT [48]	Utterance Group	SISO	✗	2	136 / 136	—	3.4	—	—
	SA-CTC SOT [49]	Utterance Group	SISO	✗	2	56 / 56	3.9	4.1	22.6	22.6
LibriMix	Whisper-SS-TTI [46]	Utterance Group	SIMO	✗	2, 3	9 / 250	—	5.2	—	8.6
	Whisper-SS-TTI [46]	Utterance Group	SIMO	✗	2, 3	13 / 779	—	4.0	—	7.5
	Whisper-SS-TTI [46]	Utterance Group	SIMO	✗	2, 3	18 / 1950	—	3.4	—	6.8
	MT-LLM [81]	Utterance Group	SISO	✗	2, 3	76.6 / 7550	—	5.2	—	10.2
	WavLM/wTSE&JSM [65]	Utterance Group	SIMO	✓	2	13 / 108	—	10.7†	—	—
LibriSpeechMix	Whisper-SS-TTI [46]	Utterance Group	SIMO	✗	2, 3	9 / 250	—	9.4	—	26.8
	Whisper-SS-TTI [46]	Utterance Group	SIMO	✗	2, 3	13 / 779	—	6.6	—	21.5
	Whisper-SS-TTI [46]	Utterance Group	SIMO	✗	2, 3	18 / 1950	—	4.7	—	16.8
	W2V-Sidecar [45]	Utterance Group	SIMO	✗	2	9 / 104	7.7	8.1	—	—
	GEncSep [82]	Utterance Group	SISO	✗	2, 3	—	6.4	6.6	13.3	13.1
	Hypothesis Stitcher [70]	Long-term	SISO	✗	1 - 6	—	—	11.5	—	13.4
	Hypothesis Clustering [83]	Long-term	SISO	✓	1, 2, 3	—	—	8.2†	—	21.5†

### A. Performance Comparison of E2E Approaches

This section compares various techniques and their accuracy across standard multi-speaker benchmarks. We report results on both real-world data (AMI) and simulated mixtures (LibriSpeechMix and LibriMix), as shown in Table II. LibriSpeechMix [27] provides the standard dev/test set for evaluation. To facilitate comparison between methods, we adopt cpWER as the primary evaluation metric to include speaker determination in WER. Note that some methods only report results on standard WER and SA-WER (see the results marked with \* and †); those accuracies cannot be directly compared with those of other methods. All results are reported directly from the original papers. Our analysis examines how different scenarios and training approaches impact results.

The results demonstrate no consistent superiority between SISO and SIMO approaches in multi-speaker ASR. For instance, SISO methods [37], [82] outperform SIMO [65], [45] on AMI and LibriMix respectively, while SIMO method [46] surpasses SISO [24] on LibriSpeechMix. Moreover, cpWER has not shown consistent improvement throughout the six-year development period of end-to-end multi-speaker ASR approaches. The currently best performance on AMI comes from a relatively small 50M model trained on extensive 900k hours of simulated mixture data [37] in 2021.

In general, contemporary research on multi-speaker ASR is less about pursuing marginal WER improvements on specific benchmark datasets. Instead, it tends to follow three trends: (1) **Understanding scenario-dependent factors:** For example, LibriMix methods demonstrate this evolution: while early works [27] used no speaker enrollment, subsequent studies [23], [24] introduced speaker-attributed (SA) methods with enrollment. This was further advanced by CIF SA-SOT [48], which achieved speaker attribution without enrollment. Meanwhile, Transcribe-to-Diarize [68] extends the approach from [23] to long-form scenarios. (2) **Exploring novel architectures and information fusion methods**, such as cross-speaker encoding for SIMO [42], enhancing CTC loss with speaker attribute [49]. (3) **Efficiently adapting single-speaker foundation ASR systems to multi-speaker settings with minimal training**, such as Adapted USM [58], WavLM/wTSE&JSM [65], Whisper-SS-TTI [46], W2V-Sidecar [45]. The # Parameters (Train / Total) in tables partially reflect the reduced training effort, though they don't fully capture cases like META-CAT [52], which achieves efficiency through fewer training epochs despite its larger trained size.

Currently, limited open-source availability makes fair comparisons difficult and slows research progress. Lots of methods show their improvements by comparing with their own

baseline systems that don't include the new architectures. This highlights the importance of developing standardized benchmarks and sharing reproducible models as a community.

## VI. CONCLUSION AND FUTURE DIRECTIONS

Recent research on multi-speaker ASR has increasingly focused on end-to-end (E2E) methods, aiming to produce more accurate speaker-distinguishable transcriptions in scenarios such as phone calls, meetings, and group discussions. E2E methods can overcome the limitations of traditional modular systems such as error propagation and failure to leverage cross-task synergies. Furthermore, progress in single-speaker ASR, speech separation, and diarization has provided the architectural foundation and pretrained model for initialization that facilitate the training of a multi-speaker ASR system, although with the scarcity of large-scale multi-speaker data.

This review provided an overview of end-to-end multi-speaker ASR approaches, from segment-level methods to processing continuous long-form recordings. Various architectural designs (Section III) were described for how to disentangle different speakers, and how to effectively leverage contextual cues from mixed speech, such as inter-speaker, overlapping, and temporal dependencies. Recent work on long-form processing (Section IV), which has further improved applicability to real-world cases, was also presented. Our accuracy comparisons indicate that, while end-to-end models often outperform modular approaches, no single architecture consistently outperforms others across end-to-end designs. Ongoing research continues to explore more sophisticated designs leveraging complex interaction patterns, while integrating advances in large-scale pretrained ASR models.

Overall, the future of the multi-speaker ASR lies in developing more robust, adaptive, and scalable systems for various scenarios. While recent designs have made significant progress, key challenges remain in robustly handling overlapping speech, and designing effective SISO/SIMO hybrids (Section III-D). Moreover, future research may also explore adaptive modeling, such as transforming single-speaker models into multi-speaker systems without performance degradation, and optional speaker enrollment. In addition, advancing multi-speaker understanding – through joint training with downstream objectives, cross-modal fusion, and injecting multi-speaker ASR into a unified multi-task foundation model.

**Acknowledgment:** This research was supported by the NSF National AI Institute for Student-AI Teaming (iSAT) under grant DRL #2019805, and also from grant #2046505. The opinions expressed are those of the authors and do not represent views of the NSF.

## REFERENCES

- [1] J. Hershey, S. Rennie, P. Olsen, and T. Kristjansson, "Superhuman multi-speaker speech recognition: A graphical modeling approach," *Computer Speech & Language*, 01 2010.
- [2] S. Settle, J. L. Roux, T. Hori, S. Watanabe, and J. R. Hershey, "End-to-end multi-speaker speech recognition," in *ICASSP*, 2018.
- [3] S. Watanabe, M. I. Mandel, J. Barker, and E. Vincent, "Chime-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings," 2020.
- [4] Y.-m. Qian, C. Weng, X.-k. Chang, S. Wang, and D. Yu, "Past review, current progress, and challenges ahead on the cocktail party problem," *Frontiers of Info. Technology and Electronic Engineering*, 01 2018.
- [5] S. Arora and R. Singh, "Automatic speech recognition: A review," *International Journal of Computer Applications*, 12 2012.
- [6] M. Malik, M. Malik, K. Mehmood, and I. Makhdoom, "Automatic speech recognition: a survey," *Multimedia Tools and Applications*, 03 2021.
- [7] R. Prabhavalkar, T. Hori, T. N. Sainath, R. Schlüter, and S. Watanabe, "End-to-end speech recognition: A survey," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [8] N. Kanda, X. Xiao, J. Wu, T. Zhou, Y. Gaur, X. Wang, Z. Meng, Z. Chen, and T. Yoshioka, "A comparative study of modular and joint approaches for speaker-attributed asr on monaural long-form audio," in *ASRU*, 2021.
- [9] F. Yu, Z. Du, S. Zhang, Y. Lin, and L. Xie, "A comparative study on speaker-attributed automatic speech recognition in multi-party meetings," in *Interspeech*, 2022.
- [10] X. Anguera, S. Bozonnet, N. Evans, C. Fredouille, G. Friedland, and O. Vinyals, "Speaker diarization: A review of recent research," *IEEE Transactions on Audio, Speech & Language Processing*, 02 2012.
- [11] G. Sell and D. Garcia-Romero, "Speaker diarization with plda i-vector scoring and unsupervised calibration," *SLT*, 2014.
- [12] M. Díez, L. Burget, and P. Matejka, "Speaker diarization based on bayesian hmm with eigenvoice priors," in *The Speaker and Language Recognition Workshop*, 2018.
- [13] F. Landini, J. Profant, M. Diez, and L. Burget, "Bayesian hmm clustering of x-vector sequences (vbx) in speaker diarization: theory, implementation and analysis on standard tasks," 12 2020.
- [14] Y. Fujita, N. Kanda, S. Horiguchi, K. Nagamatsu, and S. Watanabe, "End-to-end neural speaker diarization with permutation-free objectives," in *Interspeech*, 2019.
- [15] ———, "End-to-End Neural Speaker Diarization with Permutation-free Objectives," in *Interspeech*, 2019.
- [16] Z. Li and J. Whitehill, "Compositional embedding models for speaker identification and diarization with simultaneous speech from 2+ speakers," in *ICASSP*, 2021.
- [17] Z. Li, X. He, and J. Whitehill, "Compositional clustering: Applications to multi-label object recognition and speaker identification," *Pattern Recognition*, 07 2023.
- [18] D. Yu, X. Chang, and Y. Qian, "Recognizing multi-talker speech with permutation invariant training," in *Interspeech*, 2017.
- [19] H. Seki, T. Hori, S. Watanabe, J. Le Roux, and J. R. Hershey, "A purely end-to-end system for multi-speaker speech recognition," in *Proceedings of the Annual Meeting of the Assoc. for Comp. Linguistics*, Jul. 2018.
- [20] X. Chang, W. Zhang, Y. Qian, J. L. Roux, and S. Watanabe, "End-to-end multi-speaker speech recognition with transformer," in *ICASSP*, 2020.
- [21] D. Yu, M. Kolbæk, Z.-H. Tan, and J. Jensen, "Permutation invariant training of deep models for speaker-independent multi-talker speech separation," in *ICASSP*, 2017.
- [22] M. Kolbæk, D. Yu, Z.-H. Tan, and J. H. Jensen, "Multi-talker speech separation and tracing with permutation invariant training of deep recurrent neural networks," *ArXiv*, 2017.
- [23] N. Kanda, Y. Gaur, X. Wang, Z. Meng, Z. Chen, T. Zhou, and T. Yoshioka, "Joint speaker counting, speech recognition, and speaker identification for overlapped speech of any number of speakers," *ArXiv*, 2020.
- [24] N. Kanda, G. Ye, Y. Gaur, X. Wang, Z. Meng, Z. Chen, and T. Yoshioka, "End-to-end speaker-attributed asr with transformer," 08 2021.
- [25] Y. Liang, F. Yu, Y. Li, P. Guo, S. Zhang, Q. Chen, and L. Xie, "Bassot: Boundary-aware serialized output training for multi-talker asr," in *Interspeech*, 2023.
- [26] H. Shi, Y. Gao, Z. Ni, and T. Kawahara, "Serialized speech information guidance with overlapped encoding separation for multi-speaker automatic speech recognition," 09 2024.
- [27] N. Kanda, Y. Gaur, X. Wang, Z. Meng, and T. Yoshioka, "Serialized output training for end-to-end overlapped speech recognition," in *Interspeech*, 2020.
- [28] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *ICML*. JMLR.org, 2023.
- [29] A. Baevski, H. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," 2020.
- [30] W. Hsu, B. Bolte, Y. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," 2021.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," 2023.

- [32] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, and R. Pang, “Conformer: Convolution-augmented transformer for speech recognition,” 2020.
- [33] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *ICASSP*, 2016.
- [34] Z.-Q. Wang, J. L. Roux, and J. R. Hershey, “Alternative objective functions for deep clustering,” in *ICASSP*, 2018.
- [35] Y. Luo and N. Mesgarani, “Tasnet: time-domain audio separation network for real-time, single-channel speech separation,” 2017.
- [36] D. Yu, M. Kolbæk, Z.-H. Tan, and J. Jensen, “Permutation invariant training of deep models for speaker-independent multi-talker speech separation,” in *ICASSP*, 2017.
- [37] N. Kanda, G. Ye, Y. Wu, Y. Gaur, X. Wang, Z. Meng, Z. Chen, and T. Yoshioka, “Large-scale pre-training of end-to-end multi-talker asr for meeting transcription with single distant microphone,” 03 2021.
- [38] Y. Lin, Z. Du, S. Zhang, F. Yu, Z. Zhao, and F. Wu, “Separate-to-recognize: Joint multi-target speech separation and speech recognition for speaker-attributed asr,” in *ISCSLP*, 2022.
- [39] W. Zhang, X. Chang, and Y. Qian, “Knowledge distillation for end-to-end monaural multi-talker asr system,” in *Interspeech 2019*, 2019.
- [40] W. Zhang, X. Chang, Y. Qian, and S. Watanabe, “Improving end-to-end single-channel multi-talker speech recognition,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2020.
- [41] X. Chang, Y. Qian, K. Yu, and S. Watanabe, “End-to-end monaural multi-speaker asr system without pretraining,” in *ICASSP*, 2019.
- [42] J. Kang, L. Meng, M. Cui, H. Guo, X. Wu, X. Liu, and H. Meng, “Cross-speaker encoding network for multi-talker speech recognition,” *ICASSP*, 2024.
- [43] T. von Neumann, C. Boeddeker, L. Drude, K. Kinoshita, M. Delcroix, T. Nakatani, and R. Haeb-Umbach, “Multi-talker asr for an unknown number of sources: Joint training of source counting, separation and asr,” *ArXiv*, 2020.
- [44] Y. Luo, Z. Chen, and T. Yoshioka, “Dual-path rnn: efficient long sequence modeling for time-domain single-channel speech separation,” 2020.
- [45] L. Meng, J. Kang, M. Cui, Y. Wang, X. Wu, and H. Meng, “A sidecar separator can convert a single-talker speech recognition system to a multi-talker one,” in *ICASSP*, 2023.
- [46] L. Meng, J. Kang, Y. Wang, Z. Jin, X. Wu, X. Liu, and H. Meng, “Empowering whisper as a joint multi-talker and target-talker speech recognition system,” 09 2024.
- [47] T. Li, F. Wang, W. Guan, L. Huang, Q. Hong, and L. Li, “Improving multi-speaker asr with overlap-aware encoding and monotonic attention,” in *ICASSP*, 2024.
- [48] Z. Fan, L. Dong, J. Zhang, L. Lu, and Z. Ma, “Sa-sot: Speaker-aware serialized output training for multi-talker asr,” in *ICASSP*, 2024.
- [49] J. Kang, L. Meng, M. Cui, Y. Wang, X. Wu, X. Liu, and H. Meng, “Disentangling speakers in multi-talker speech recognition with speaker-aware ctc,” in *ICASSP*, 2025.
- [50] L. Zheng, H. Zhu, S. Tian, Q. Zhao, and T. Li, “Unsupervised domain adaptation on end-to-end multi-talker overlapped speech recognition,” *IEEE Signal Processing Letters*, 2024.
- [51] T. J. Park, I. Medennikov, K. Dhawan, W. Wang, H. Huang, N. R. Koluguri, K. C. Puvvada, J. Balam, and B. Ginsburg, “Sortformer: Seamless integration of speaker diarization and asr by bridging timestamps and tokens,” *ArXiv*, 2024.
- [52] J. Wang, W. Wang, K. Dhawan, T. Park, M. Kim, I. Medennikov, H. Huang, N. Koluguri, J. Balam, and B. Ginsburg, “Meta-cat: Speaker-informed speech embeddings via meta information concatenation for multi-talker asr,” in *ICASSP 2025*, 2025.
- [53] N. Kanda, X. Chang, Y. Gaur, X. Wang, Z. Meng, Z. Chen, and T. Yoshioka, “Investigation of end-to-end speaker-attributed asr for continuous multi-talker recordings,” in *SLT*, 2021.
- [54] M. Shi, Z. Du, Q. Chen, F. Yu, Y. Li, S. Zhang, J. Zhang, and L.-R. Dai, “Casa-asr: Context-aware speaker-attributed asr,” in *Interspeech*, 2023.
- [55] S. Cornell, J.-W. Jung, S. Watanabe, and S. Squartini, “One model to rule them all? towards end-to-end joint speaker diarization and speech recognition,” in *ICASSPs*, 2024.
- [56] N. Makishima, N. Kawata, M. Ihori, T. Tanaka, S. Orihashi, A. Ando, and R. Masumura, “Somsred: Sequential output modeling for joint multi-talker overlapped speech recognition and speaker diarization,” *Interspeech 2024*, 2024.
- [57] N. Makishima, K. Suzuki, S. Suzuki, A. Ando, and R. Masumura, “Joint autoregressive modeling of end-to-end multi-talker overlapped speech recognition and utterance-level timestamp prediction,” in *Interspeech*, 2023.
- [58] C. Li, Y. Qian, Z. Chen, N. Kanda, D. Wang, T. Yoshioka, Y. Qian, and M. Zeng, “Adapting multi-lingual asr models for handling multiple talkers,” *ArXiv*, 2023.
- [59] R. Masumura, D. Okamura, N. Makishima, M. Ihori, A. Takashima, T. Tanaka, and S. Orihashi, “Unified autoregressive modeling for joint end-to-end multi-talker overlapped speech recognition and speaker attribute estimation,” in *Interspeech*, 2021.
- [60] P. Denisov and N. T. Vu, “End-to-end multi-speaker speech recognition using speaker embeddings and transfer learning,” *ArXiv*, 2019.
- [61] R. Rose, O. Chang, and O. Siohan, “Cascaded encoders for fine-tuning asr models on overlapped speech,” 2023.
- [62] M. Shi, Z. Jin, Y. Xu, Y. Xu, S.-X. Zhang, K. Wei, Y. Shao, C. Zhang, and D. Yu, “Advancing multi-talker asr performance with large language models,” 2024.
- [63] W. Wang, K. Dhawan, T. Park, K. Puvvada, I. Medennikov, S. Majumdar, H. Huang, J. Balam, and B. Ginsburg, “Resource-efficient adaptation of speech foundation models for multi-speaker asr,” 12 2024.
- [64] S. Latif, M. Shoukat, F. Shamshad, M. Usama, H. Cuayáhuitl, and B. Schuller, “Sparks of large audio models: A survey and outlook,” 08 2023.
- [65] Z. Huang, D. Raj, P. García, and S. Khudanpur, “Adapting self-supervised models to multi-talker speech recognition using speaker embeddings,” in *ICASSP*, 2023.
- [66] M. Shi, Z. Jin, Y. Xu, Y. Xu, S.-X. Zhang, K. Wei, Y. Shao, C. Zhang, and D. Yu, “Advancing multi-talker asr performance with large language models,” *SLT*, 2024.
- [67] W. R. Huang, S. yiin Chang, D. Rybach, R. Prabhavalkar, T. N. Sainath, C. Allauzen, C. Peysier, and Z. Lu, “E2e segmenter: Joint segmenting and decoding for long-form asr,” in *Interspeech*, 2022.
- [68] N. Kanda, X. Xiao, Y. Gaur, X. Wang, Z. Meng, Z. Chen, and T. Yoshioka, “Transcribe-to-diarize: Neural speaker diarization for unlimited number of speakers using end-to-end speaker-attributed asr,” in *ICASSP*, 2022.
- [69] C.-C. Chiu, W. Han, Y. Zhang, R. Pang, S. Kishchenko, P. Nguyen, A. Narayanan, H. Liao, S. Zhang, A. Kannan, R. Prabhavalkar, Z. Chen, T. N. Sainath, and Y. Wu, “A comparison of end-to-end models for long-form speech recognition,” *ASRU*, 2019.
- [70] X. Chang, N. Kanda, Y. Gaur, X. Wang, Z. Meng, and T. Yoshioka, “Hypothesis sticher for end-to-end speaker-attributed asr on long-form multi-talker recordings,” in *ICASSP*, 2021.
- [71] H. H. Mao, S. Li, J. McAuley, and G. Cottrell, “Speech recognition and multi-speaker diarization of long conversations,” in *Interspeech*, 2020.
- [72] M. Yang, N. Kanda, X. Wang, J. Wu, S. Sivasankaran, Z. Chen, J. Li, and T. Yoshioka, “Simulating realistic speech overlaps improves multi-talker asr,” in *ICASSP*, 2023.
- [73] T. Moriya, S. Horiguchi, M. Delcroix, R. Masumura, T. Ashihara, H. Sato, K. Matsuura, and M. Mimura, “Alignment-free training for transducer-based multi-talker asr,” 09 2024.
- [74] F. Yu, S. Zhang, Y. Fu, L. Xie, S. Zheng, Z. Du, W. Huang, P. Guo, Z. Yan, B. Ma, X. Xu, and H. Bu, “M2met: The icassp 2022 multi-channel multi-party meeting transcription challenge,” in *ICASSP*, 2022.
- [75] Z. Chen, T. Yoshioka, L. Lu, T. Zhou, Z. Meng, Y. Luo, J. Wu, and J. Li, “Continuous speech separation: dataset and analysis,” 2020.
- [76] J. Cosentino, M. Pariente, S. Cornell, A. Deleforge, and E. Vincent, “Librimix: An open-source dataset for generalizable speech separation,” 2020.
- [77] Z. Jin, Y. Yang, M. Shi, W. Kang, X. Yang, Z. Yao, F. Kuang, L. Guo, L. Meng, L. Lin, Y. Xu, and S.-X. Zhang, “Libriheavymix: A 20,000-hour dataset for single-channel reverberant multi-talker speech separation, asr and speaker diarization,” 09 2024.
- [78] S. Watanabe, M. I. Mandel, J. Barker, and E. Vincent, “Chime-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings,” 2020.
- [79] T. von Neumann, C. Boeddeker, M. Delcroix, and R. Haeb-Umbach, “Meeteval: A toolkit for computation of word error rates for meeting transcription systems,” 07 2023.
- [80] N. Kanda, Z. Meng, L. Lu, Y. Gaur, X. Wang, Z. Chen, and T. Yoshioka, “Minimum bayes risk training for end-to-end speaker-attributed asr,” in *ICASSP*, 2021.
- [81] L. Meng, S. Hu, J. Kang, Z. Li, Y. Wang, W. Wu, X. Wu, X. Liu, and H. Meng, “Large language model can transcribe speech in multi-talker scenarios with versatile instructions,” 2024.
- [82] Y. Shi, L. Li, S. Yin, D. Wang, and J. Han, “Serialized output training by learned dominance,” *ArXiv*, 2024.
- [83] Y. Kashiwagi, H. Futami, E. Tsunoo, S. Arora, and S. Watanabe, “Hypothesis clustering and merging: Novel multitalker speech recognition with speaker tokens,” 09 2024.