

A Review of Speaker Diarization: Recent Advances with Deep Learning

Tae Jin Park^{a,*}, Naoyuki Kanda^{b,*}, Dimitrios Dimitriadis^{b,*}, Kyu J. Han^{c,*}, Shinji Watanabe^{d,*}, Shrikanth Narayanan^a

^aUniversity of Southern California, Los Angeles, USA

^bMicrosoft, Redmond, USA

^cASAPP, Mountain View, USA

^dJohns Hopkins University, Baltimore, USA

Abstract

Speaker diarization is a task to label audio or video recordings with classes corresponding to speaker identity, or in short, a task to identify “who spoke when”. In the early years, speaker diarization algorithms were developed for speech recognition on multi-speaker audio recordings to enable speaker adaptive processing, but also gained its own value as a stand-alone application over time to provide speaker-specific meta information for downstream tasks such as audio retrieval. More recently, with the rise of deep learning technology that has been a driving force to revolutionary changes in research and practices across speech application domains in the past decade, more rapid advancements have been made for speaker diarization. In this paper, we review not only the historical development of speaker diarization technology but also the recent advancements in neural speaker diarization approaches. We also discuss how speaker diarization systems have been integrated with speech recognition applications and how the recent surge of deep learning is leading the way of jointly modeling these two components to be complementary to each other. By considering such exciting technical trends, we believe that it is a valuable contribution to the community to provide a survey work by consolidating the recent developments with neural methods and thus facilitating further progress towards a more efficient speaker diarization.

Keywords: speaker diarization, automatic speech recognition, deep learning

1. Introduction

“Diarize” is a word that means making a note or keeping an event in a diary. Speaker diarization, like keeping a record of events in such a diary, addresses the “who spoke when” question [1, 2, 3] by logging speaker-specific salient events on multi-participant (or multi-speaker) audio data. Throughout the diarization process, the audio data would be divided and clustered into groups of speech segments with the same speaker identity/label. As a result, salient events, such as non-speech/speech transition, speaker turn changes, speaker classification or speaker role identification, are labeled in an automatic fashion. In general, this process does not require any prior knowledge of the speakers, such as their real identity or number of participating speakers in the audio data. Thanks to its innate feature of separating audio streams by these speaker-specific events, speaker diarization can be effectively employed for indexing or analyzing various types of audio data, e.g., audio/video broadcasts from media stations, conversations in conferences, personal videos from online social media or hand-held devices, court proceedings, business meetings, earnings reports in a financial sector, just to name a few.

Traditionally speaker diarization systems consist of multiple, independent sub-modules as shown in Fig. 1. In order to mitigate any artifacts in acoustic environments, various front-end

processing techniques, for example, speech enhancement, dereverberation, speech separation or target speaker extraction, are utilized. Voice or speech activity detection is then applied to separate speech from non-speech events. Raw speech signals in the selected speech portions are transformed to acoustic features or embedding vectors. In the clustering stage, the speech portion represented by the embedding vectors are grouped and labeled by speaker classes and in the post-processing stage, the clustering results are further refined. Each of these sub-modules is optimized individually in general.

1.1. Historical development of speaker diarization

During the early years of diarization technology (in the 1990s), the research focus was on unsupervised speech segmentation and clustering of acoustic events including not only speaker-specific ones but also those related to environmental or background changes [4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14]. In this period some of the fundamental approaches to speaker change detection and clustering, such as leveraging Generalized Likelihood Ratio (GLR) and Bayesian Information Criterion (BIC), were developed and quickly became the golden standard. Most of the works benefited Automatic Speech Recognition (ASR) on broadcast news recordings, by enabling speaker adaptive training of acoustic models [10, 15, 16, 17, 18]. All these efforts collectively laid out a path to consolidate activities across research groups around the world, leading to several research consortia and challenges in the early 2000s,

*Authors contributed equally

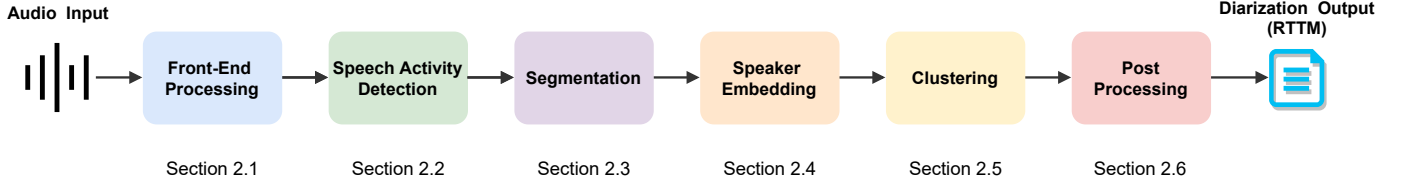


Fig. 1: Traditional Speaker Diarization Systems.

among which there were the Augmented Multi-party Interaction (AMI) Consortium [19] supported by the European Commission and the Rich Transcription Evaluation [20] hosted by the National Institute of Standards and Technology (NIST). These organizations, spanning over from a few years to a decade, had fostered further advancements on speaker diarization technologies across different data domains from broadcast news [21, 22, 23, 24, 25, 26, 27, 28, 29] and Conversational Telephone Speech (CTS) [24, 30, 31, 32, 33, 34] to meeting conversations [35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45]. The new approaches resulting from these advancements include, but not limited to, Beamforming [42], Information Bottleneck Clustering (IBC) [44], Variational Bayesian (VB) approaches [33, 45], Joint Factor Analysis (JFA) [46, 34].

Since the advent of deep learning in the 2010s, there has been a considerable amount of research to take the advantage of powerful modeling capabilities of the neural networks for speaker diarization. One representative example is extracting the speaker embeddings using neural networks, such as the d-vectors [47, 48, 49] or the x-vectors [50], which most often are embedding vector representations based on the bottleneck layer output of a “Deep Neural Network” (DNN) trained for speaker recognition. The shift from i-vector [51, 52, 53, 54] to these neural embeddings contributed to enhanced performance, easier training with more data [55], and robustness against speaker variability and acoustic conditions. More recently, End-to-End Neural Diarization (EEND) where individual sub-modules in the traditional speaker diarization systems (c.f., Fig. 1) can be replaced by one neural network gets more attention with promising results [56, 57]. This research direction, although not fully matured yet, could open up unprecedented opportunities to address challenges in the field of speaker diarization, such as, the joint optimization with other speech applications, with overlapping speech, if large-scale data is available for training such powerful network-based models.

1.2. Motivation

Till now, there are two well-rounded overview papers in the area of speaker diarization surveying the development of speaker diarization technology with different focuses. In [2], various speaker diarization systems and their subtasks in the context of broadcast news and CTS data are reviewed up to the point of mid 2000s. As such, the historical progress of speaker diarization technology development in the 1990s and early 2000s are hence covered. In contrast, the focus of [3] was put more on speaker diarization for meeting speech and its respective challenges. This paper thus weighs more in the

corresponding technologies to mitigate problems from the perspective of meeting environments, where there are usually more participants than broadcast news or CTS data and multi-modal data is frequently available. Since these two papers, especially thanks to leap-frog advancements in deep learning approaches addressing technical challenges across multiple machine learning domains, speaker diarization systems have gone through a lot of notable changes. We believe that this survey work is a valuable contribution to the community to consolidate the recent developments with neural methods and thus facilitate further progress towards a more efficient diarization.

1.3. Overview and Taxonomy of speaker diarization

Attempting to categorize the existing, most-diverse speaker diarization technologies, both on the space of modularized speaker diarization systems before the deep learning era and those based on neural networks of the recent years, a proper grouping would be helpful. The main categorization we adopt in this paper is based on two criteria, resulting in the total four categories, as shown in Table 1. The first criterion is whether the model is trained based on speaker diarization-oriented objective function or not. Any trainable approaches to optimize models in a multi-speaker situation and learn relations between speakers are categorized into the “Diarization Objective” class. The second criterion is whether multiple modules are jointly optimized towards some objective function. If a single sub-module is replaced into a trainable one, such method is categorized into the “Single-module Optimization” class. On the other hand, for example, joint modeling of segmentation and clustering [55], joint modeling of speech separation and speaker diarization [76] or fully end-to-end neural diarization [56, 57] is categorized into the “Joint Optimization” class.

Note that our intention of this categorization is to help readers to quickly overview the broad development in the field, and it is not our intention to divide the categories into superior-inferior. Also, while we are aware of many techniques that fall into the category “Non-Diarization Objective” and “Joint Optimization” (e.g., joint front-end and ASR [67, 68, 69, 70, 71, 72], joint speaker identification and speech separation [73, 74], etc.), we exclude them in the paper to focus on the review of speaker diarization techniques.

1.4. Paper Organization

The rest of the paper is organized as follows.

- In Section 2, we overview techniques belonging to the “Non-Diarization Objective” and “Single-module Optimization” class in the proposed taxonomy, mostly those

Table 1: Table of Taxonomy

	Non-Diarization Objective	Diarization Objective
Single-module Optimization	Section 2 Front-end [58, 59, 60], speaker embedding [61, 62, 50], speech activity detection [63], etc.	Section 3.1 IDEC [64], affinity matrix refinement [65], TS-VAD [66], etc.
Joint Optimization	Out of scope Joint front-end & ASR [67, 68, 69, 70, 71, 72], joint speaker identification & speech separation [73, 74], etc.	Section 3.2 UIS-RNN [55], RPN [75], online RSAN [76], EEND [56, 57], etc. Section 4 Joint ASR & speaker diarization. [77, 78, 79, 80], etc.

used in the traditional, modular speaker diarization systems. While there are some overlaps with the counterpart sections of the aforementioned two survey papers [2, 3] in terms of reviewing notable developments in the past, this section would add more latest schemes as well in the corresponding components of the speaker diarization systems.

- In Section 3, we discuss advancements mostly leveraging DNNs trained with the diarization objective where single sub-modules are independently optimized (subsection 3.1) or jointly optimized (subsection 3.2) toward fully end-to-end speaker diarization.
- In Section 4, we present a perspective of how speaker diarization has been investigated in the context of ASR, reviewing historical interactions between these two domains to peek the past, present and future of speaker diarization applications.
- Section 5 provides information of speaker diarization challenges and corpora to facilitate research activities and anchor technology advances. We also discuss evaluation metrics such as Diarization Error Rate (DER), Jaccard Error Rate (JER) and Word-level DER (WDER) in the section.
- We share a few examples of how speaker diarization systems are employed in both research and industry practices in Section 6 and conclude this work in Section 7 with providing summary and future challenges in speaker diarization.

2. Modular Speaker Diarization Systems

This section provides an overview of algorithms for speaker diarization belonging to the “Single-module Optimization, Non-Diarization Objective” class mostly modular speaker, as shown in Figure 1. Each subsection in this section corresponds to the explanation of each module in the traditional speaker diarization system. In addition to the introductory explanation of

each module, this section also summarizes the latest schemes within the module.

2.1. Front-end processing

This section describes mostly front-end techniques, used for speech enhancement, dereverberation, speech separation, and speech extraction as part of the speaker diarization pipeline. Let $s_{i,f,t} \in \mathbb{C}$ be the STFT representation of source speaker i on frequency bin f at frame t . The observed noisy signal $x_{t,f}$ can be represented by a mixture of the source signals, a room impulse response $h_{i,f,t} \in \mathbb{C}$, and additive noise $n_{t,f} \in \mathbb{C}$,

$$x_{t,f} = \sum_{i=1}^K \sum_{\tau} h_{i,f,\tau} s_{i,f,t-\tau} + n_{t,f}, \quad (1)$$

where K denotes the number of speakers present in the audio signal.

The front-end techniques described in this section is to estimate the original source signal $\hat{\mathbf{x}}_{i,t}$ given the observation $\mathbf{X} = (\{x_{t,f}\}_f)_t$ for the downstream diarization task,

$$\hat{\mathbf{x}}_{i,t} = \text{FrontEnd}(\mathbf{X}), \quad i = 1, \dots, K, \quad (2)$$

where $\hat{\mathbf{x}}_{i,t} \in \mathbb{C}^D$ is the i -th speaker’s estimated STFT spectrum with D frequency bins at frame t .

Although there are numerous speech enhancement, dereverberation, and separation algorithms, e.g., [81, 82, 83], herein most of the recent techniques used in the DIHARD challenge series [84, 85, 86], LibriCSS meeting recognition task [87, 88], and CHiME-6 challenge track 2 [89, 90, 91] are covered.

2.1.1. Speech enhancement (Denoising)

Speech enhancement techniques focus mainly on suppressing the noise component of the noisy speech. Single-channel speech enhancement has shown a significant improvement in denoising performance [92, 93, 94] thanks to deep learning, when compared with classical signal processing based speech enhancement [95]. For example, LSTM-based speech enhancement [96, 94] is used as a front-end technique in the DIHARD

II baseline [85], i.e.,

$$\hat{\mathbf{x}}_t = \text{LSTM}(\mathbf{X}), \quad (3)$$

where we only consider the single source example (i.e., $K = 1$) and omit the source index i . This is a regression-based approach by minimizing the objective function,

$$\mathcal{L}_{\text{MSE}} = \|\mathbf{s}_t - \hat{\mathbf{x}}_t\|^2. \quad (4)$$

The log power spectrum or ideal ratio mask is often used as the target domain of the output \mathbf{s}_t . Also, the speech enhancement used in [95] applies this objective function in each layer based on a progressive manner.

The effectiveness of the speech enhancement techniques can be boosted multi-channel processing, including minimum variance distortionless response (MVDR) beamforming [81]. [88] shows the significant improvement of the DER from 18.3% to 13.9% in the LibriCSS meeting task based on mask-based MVDR beamforming [97, 98].

2.1.2. Dereverberation

Compared with other front-end techniques, the major dereverberation techniques used in various tasks is based on statistical signal processing methods. One of the most widely used techniques is Weighted Prediction Error (WPE) based dereverberation [99, 100, 101].

The basic idea of WPE, for the case of single source, i.e. $K = 1$, without noise, is to decompose the original signal model Eq. (1) into the early reflection $x_{t,f}^{\text{early}}$ and late reverberation $x_{t,f}^{\text{late}}$ as follows:

$$x_{t,f} = \sum_{\tau} h_{f,\tau} s_{f,t-\tau} = x_{t,f}^{\text{early}} + x_{t,f}^{\text{late}}. \quad (5)$$

WPE tries to estimate filter coefficients $\hat{h}_{f,t}^{\text{wpe}} \in \mathbb{C}$, which maintain the early reflection while suppress the late reverberation based on the maximum likelihood estimation.

$$\hat{x}_{t,f}^{\text{early}} = x_{t,f} - \sum_{\tau=\Delta}^L \hat{h}_{f,\tau}^{\text{wpe}} x_{f,t-\tau}, \quad (6)$$

where Δ is the number of frames to split the early reflection and late reverberation, and L is the filter size.

WPE is widely used as one of the golden standard front-end processing methods, e.g., it is part of the DIHARD and CHiME both the baseline and the top-performing systems [84, 85, 86, 89, 90]. Although the performance improvement of WPE-based dereverberation is not significant, it provides solid performance improvement across almost all tasks. Also, WPE is based on the linear filtering and since it does not introduce signal distortions, it can be safely combined with downstream front-end and back-end processing steps. Similar to the speech enhancement techniques, WPE-based dereberberation shows additional performance improvements when applied on multi-channel signals.

2.1.3. Speech separation and target speaker extraction

Speech separation is a promising family of techniques when the overlapping speech regions are significant. Similarly to other research areas, DL-based speech separation has become popular, e.g., “Deep Clustering” [58], “Permutation Invariant Training” (PIT) [59], and Conv-TasNet [60]. The effectiveness of multi-channel speech separation based on beamforming has been widely confirmed [102, 103], as well. For example, in the CHiME-6 challenge [89], “Guided Source Separation” (GSS) [103] based multi-channel speech extraction techniques have been used to achieve the top result. On the other hand, single-channel speech separation techniques do not often show any significant effectiveness in realistic multi-speaker scenarios like the LibriCSS [87] or the CHiME-6 tasks [89], where speech signals are continuous and contain both overlapping and overlap-free speech regions. The single-channel speech separation systems often produce a redundant non-speech or even a duplicated speech signal for the non-overlap regions, and as such the “leakage” of audio causes many false alarms of speech activity. A leakage filtering method was proposed in [104] tackling the problem, where a significant improvement of speaker diarization performance was shown after including this processing step in the top-ranked system on the VoxCeleb Speaker Recognition Challenge 2020 [105].

2.2. Speech activity detection (SAD)

SAD distinguishes speech segments from non-speech segments such as background noise. A SAD system is mostly comprised of two parts. The first one is a feature extraction frontend, where acoustic features such as Mel-Frequency Cepstral Coefficients (MFCCs) are extracted. The other part is a classifier, where a model predicts whether the input frame is speech or not. These models may include Gaussian Mixture Models (GMMs) [106], Hidden Markov Models (HMMs) [107] or DNNs [63].

The performance of SAD largely affects the overall performance of the speaker diarization system because it can create a significant amount of false positive salient events or miss speech segments [108]. A common practice in speaker diarization tasks is to report DER with “oracle SAD” setup which indicates that the system output is using speech activity detection output that is identical to the ground truth. On the other hand, the system output with an actual speech activity detector is referred to as “system SAD” output.

2.3. Segmentation

Speech segmentation breaks the input audio stream into multiple segments so that the each segment can be assigned to a speaker label. Before re-segmentation phase, the unit of the output of speaker diarization system is determined by segmentation process. There are two ways of performing speech segmentation for speaker diarization tasks: either with speaker change point detection or uniform segmentation. The segmentation by detecting the speaker change point was the gold standard of the earlier speaker diarization systems, where speaker change

points are detected by comparing two hypotheses: Hypothesis H_0 assumes both left and right samples are from the same speaker and hypothesis H_1 assumes the two samples are from the different speakers. Many algorithms for the hypothesis testing, such as Kullback Leibler 2 (KL2) [10], “Generalized Likelihood Ratio” (GLR) [109] and BIC [110, 111] were proposed with the BIC method been the most widely used method. The BIC approach can be applied to segmentation process as follows: assuming that $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is the sequence of speech features extracted from the given audio stream and \mathbf{x} is drawn from an independent multivariate Gaussian process:

$$\mathbf{x}_i \sim N(\mu_i, \Sigma_i), \quad (7)$$

where μ_i, Σ_i is mean and covariance matrix of the i -th feature window, two hypothesis H_0 and H_1 can be denoted as follows:

$$H_0 : \mathbf{x}_1 \cdots \mathbf{x}_N \sim N(\mu, \Sigma) \quad (8)$$

$$H_1 : \mathbf{x}_1 \cdots \mathbf{x}_i \sim N(\mu_1, \Sigma_1) \quad (9)$$

$$\mathbf{x}_{i+1} \cdots \mathbf{x}_N \sim N(\mu_2, \Sigma_2) \quad (10)$$

Thus, hypothesis H_0 models two sample windows with one Gaussian while hypothesis H_1 models two sample windows with two Gaussians. Using the Eq. (8), the maximum likelihood ratio statistics can be expressed as

$$R(i) = N \log |\Sigma| - N_1 \log |\Sigma_1| - N_2 \log |\Sigma_2|, \quad (11)$$

where the sample covariance Σ is from $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, Σ_1 is from $\{\mathbf{x}_1, \dots, \mathbf{x}_i\}$ and Σ_2 is from $\{\mathbf{x}_{i+1}, \dots, \mathbf{x}_N\}$. Finally, a BIC value between two models is expressed:

$$BIC(i) = R(i) - \lambda P, \quad (12)$$

where P is the penalty term [110] defined as

$$P = \frac{1}{2} \left(d + \frac{1}{2} d(d+1) \right) \log N, \quad (13)$$

and d is dimension of the feature. The penalty weight λ is generally set to $\lambda = 1$. The change point is set when the following equation becomes true,

$$\left\{ \max_i BIC(i) \right\} > 0. \quad (14)$$

As described above, the speaker change points can be detected by using hypothesis testing based on BIC values or other methods such as KL2 [10], GLR [109]. However, if speech segmentation is done by speaker change point detection method, the length of each segment is not consistent. Therefore, after the advent of i-vector [51] and DNN-based embeddings [61] the segmentation based on speaker change point detection was mostly replaced by uniform segmentation [112, 113, 49], since varying length of the segment created an additional variability into the speaker representation and deteriorated the fidelity of the speaker representations. In uniform segmentation schemes, the given audio stream input is segmented with a fixed window length and overlap length. Thus, the length of the unit of speaker diarization result is remains fixed.

However, the process of uniformly segmenting the input signals for diarization poses some potential problems. First, uniform segmentation introduces a trade-off error related to the segment length: segments need to be sufficiently short to safely assume that they do not contain multiple speakers but at the same time it is necessary to capture enough acoustic information to extract a meaningful speaker representation x_j .

2.4. Speaker Representations and Speaker Embeddings

In this section, we explain a few popular methods for measuring the similarity of speech segments. These methods are paired with clustering algorithms, which will be explained in the next section. We first introduce GMM based hypothesis testing approaches which are usually employed with segmentation approaches based on a speaker change point detection. We then introduce well-known speaker representations for speaker diarization systems that are usually employed with the uniform segmentation method in Section 2.4.2 and Section 2.4.3.

2.4.1. GMM speaker model for similarity measure

The early days of speaker diarization systems were based on a GMM built on acoustic features such as the MFCCs. Along with GMM based method, AHC was also employed for clustering, resulting in the speaker homogeneous clusters. While there are many hypothesis testing methods for speech segment clustering process such as greedy BIC [110], GLR [114] and KL [115] methods, greedy BIC method was the most popular approach. While greedy BIC method also employs BIC value as in speaker change point detection, in greedy BIC method, BIC value is used for measuring the similarity between two nodes during the AHC process. For the given nodes to be clustered, $\mathcal{S} = \{s_1, \dots, s_k\}$, greedy BIC method model each node s_i as a multivariate Gaussian distribution $N(\mu_i, \Sigma_i)$ where μ_i and Σ_i are mean and covariance matrix of the merged samples in the node s_i . BIC value for merging the node s_1 and s_2 is calculated as

$$BIC = n \log |\Sigma| - n_1 \log |\Sigma_1| - n_2 \log |\Sigma_2| - \lambda P, \quad (15)$$

where λ and P value are identical to Eq. (12) and n is sample size of the merged node ($n = n_1 + n_2$). During the clustering process, we merge the nodes if Eq. (15) is negative. GMM based hypothesis testing method with bottom-up hierarchical clustering method was popularly used until i-vector and DNN-based speaker representations dominate the speaker diarization research scene.

2.4.2. Joint Factor Analysis and i-vector

Before the advent of speaker representations such as i-vector [51] or x-vector [50], “Universal Background Model” (UBM) [116] framework showed success for speaker recognition tasks by employing a large mixture of Gaussians, while covering a fairly large amount of speech data. The idea of modeling and testing the similarity of voice characteristics with GMM-UBM [116] is largely improved by JFA [117, 118]. GMM-UBM based hypothesis testing had a problem of Maximum a Posterior (MAP) adaptation that is not only affected

by speaker-specific characteristics but also other nuisance factors such as channel and background noise. Therefore, the concept of supervector generated by GMM-UBM method was not ideal. JFA tackles this problem and decompose a supervector into speaker independent, speaker dependent, channel dependent and residual components. Thus, the ideal speaker supervector s can be decomposed as in the Eq. (16). A term \mathbf{m} denotes speaker independent component, \mathbf{U} denotes channel dependent component matrix, and \mathbf{D} denotes speaker-dependent residual component matrix. Along with these component matrices, vector \mathbf{y} is for the speaker factors, vector \mathbf{x} is for the channel factors and vector \mathbf{z} is for the speaker-specific residual factors. All of these vectors have a prior distribution of $N(0, 1)$.

$$\mathbf{M}(s) = \mathbf{m} + \mathbf{V}\mathbf{y} + \mathbf{U}\mathbf{x} + \mathbf{D}\mathbf{z}. \quad (16)$$

The idea of JFA approach is further simplified by employing the so called ‘‘Total Variability’’ matrix \mathbf{T} modeling both the channel and the speaker variability, and the vector \mathbf{w} which is referred to as the ‘‘i-vector’’ [51]. The supervector \mathbf{M} is modeled as:

$$\mathbf{M} = \mathbf{m} + \mathbf{T}\mathbf{w}, \quad (17)$$

In Eq. (17), \mathbf{m} is the session and channel-independent component of the mean supervector. Similarly to JFA, \mathbf{w} is assumed to follow standard normal distribution and calculated by MAP estimation, in [119]. The notion of speaker representation is popularized by i-vectors, where the speaker representation vector can contain a numerical feature that characterize the vocal tract of each speaker. The i-vector speaker representations have employed in not only speaker recognition studies but also in numerous speaker diarization studies [112, 120, 121] and showed superior performance over GMM-based hypothesis testing methods.

2.4.3. Neural Network Based Speaker Representations

Speaker representations for speaker diarization has also been heavily affected by the rise of neural networks and deep learning approaches. The idea of representation learning was first introduced for face recognition tasks [122, 123]. The fundamental idea of neural network-based representations is that we can use deep neural network architecture to map the input signal source (an image or an audio clip) to a dense vector by sampling the activations of a layer in the neural network model. The neural network based representation does not require eigenvalue decomposition or factor analysis model that involves hand-crafted design of the intrinsic factor. Also, there is no assumption or requirement of Gaussianity for the input data. Thus, the representation learning process has become more straight-forward and the inference speed has been also improved compared to the traditional factor analysis based methods.

Among many of the neural network based speaker representations, d-vector [61] remains one of the most prominent speaker representation extraction frameworks. The d-vector employs stacked filterbank features that include context frames as an input feature and trains a multiple fully connected layers

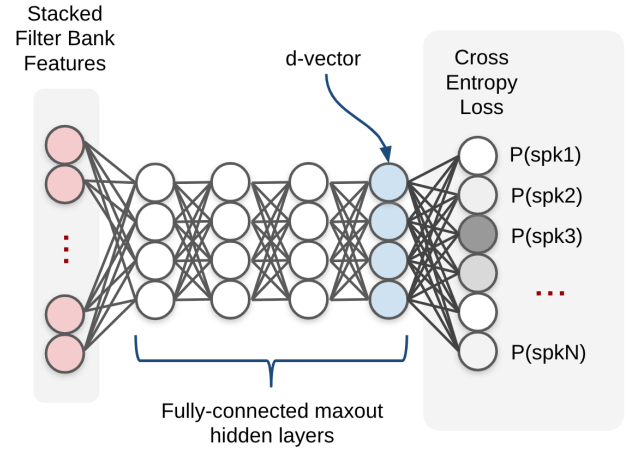


Fig. 2: Diagram of d-vector model.

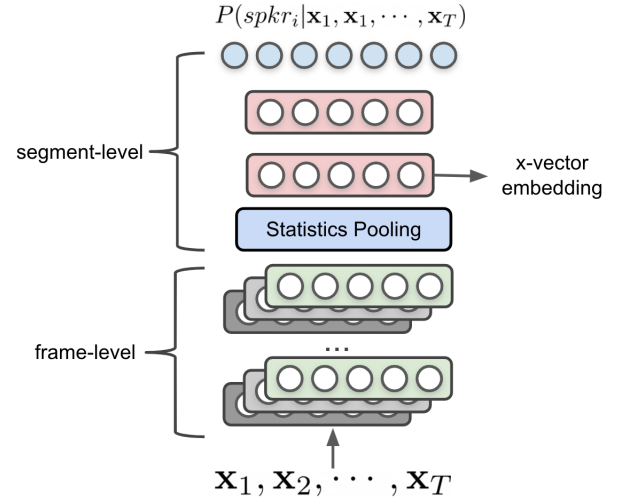


Fig. 3: Diagram of x-vector embedding extractor.

with the cross entropy loss. The d-vector embeddings are obtained in the last fully connected layer as in Fig. 2. The d-vector scheme appears in numerous speaker diarization papers, e.g., in [49, 55].

DNN-based speaker representations are even more improved by x-vector [62, 50]. The x-vector showed a superior performance by winning the NIST speaker recognition challenge [124] and the first DIHARD challenge [84]. Fig. 3 shows the structure of x-vector framework. The time-delay architecture and statistical pooling layer differentiate x-vector architecture from d-vector while statistical pooling layer mitigates the effect of the input length. This is especially advantageous when it comes to speaker diarization since the speaker diarization systems are bound to process segments that are shorter than the regular window length.

For speaker diarization tasks, ‘‘Probabilistic Linear Discriminant Analysis’’ (PLDA) has been frequently used along with x-vector or i-vector to measure the affinity between two speech segments. PLDA employs the following modeling for the given

speaker representation ϕ_{ij} of the i -th speaker and j -th session as below:

$$\phi_{ij} = \mu + \mathbf{F}\mathbf{h}_i + \mathbf{G}\mathbf{w}_{ij} + \epsilon_{ij}. \quad (18)$$

Here, \mathbf{m} is mean vector, \mathbf{F} is speaker variability matrix, \mathbf{G} is channel variability matrix and ϵ is residual component. The terms \mathbf{h}_i and \mathbf{w}_{ij} are latent variable for \mathbf{F} and \mathbf{G} respectively. During the training process of PLDA, \mathbf{m} , $\mathbf{\Sigma}$, \mathbf{F} and \mathbf{G} are estimated using expectation maximization (EM) algorithm where $\mathbf{\Sigma}$ is a covariance matrix. Based on the estimated variability matrices and the latent variables \mathbf{h}_i and \mathbf{w}_{ij} , two hypotheses are tested: hypothesis H_0 for the case that two samples are from the same speaker and hypothesis H_1 for the case that two samples are from different speakers. The hypothesis H_0 can be written as follows:

$$\begin{bmatrix} \phi_1 \\ \phi_2 \end{bmatrix} = \begin{bmatrix} \mu \\ \mu \end{bmatrix} + \begin{bmatrix} \mathbf{F} & \mathbf{G} & 0 \\ \mathbf{F} & 0 & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{h}_{12} \\ \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}. \quad (19)$$

On the other hand, The hypothesis H_1 can be modeled as the following equation.

$$\begin{bmatrix} \phi_1 \\ \phi_2 \end{bmatrix} = \begin{bmatrix} \mu \\ \mu \end{bmatrix} + \begin{bmatrix} \mathbf{F} & \mathbf{G} & 0 & 0 \\ 0 & 0 & \mathbf{F} & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{w}_1 \\ \mathbf{h}_2 \\ \mathbf{w}_2 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \end{bmatrix}. \quad (20)$$

The PLDA model projects the given speaker representation onto the subspace \mathbf{F} to co-vary the most while de-emphasizing the subspace \mathbf{G} pertaining to channel variability. Using the above hypotheses, we can calculate a log likelihood ratio.

$$s(\phi_1, \phi_2) = \log p(\phi_1, \phi_2 | H_0) - \log p(\phi_1, \phi_2 | H_1). \quad (21)$$

Ideally, stopping criterion should be 0, but in practice it varies from around zero values and the stopping criterion needs to be tuned on development set. The stopping criterion largely affects the estimated number of speakers because the clustering process stops when the distance between closest samples reaches threshold and the number of clusters is determined by the number of remaining clusters at the step where clustering is stopped.

2.5. Clustering

After generating the speaker representations for each segment, clustering algorithm is applied to make clusters of segments. We introduce the most commonly used clustering methods for speaker diarization task.

2.5.1. Mean-shift

Mean-shift [125] is a clustering algorithm that assigns the given data points to the clusters iteratively by finding the modes in a non-parametric distribution. Mean-shift algorithm follows the following steps:

1. Start with the data points assigned to a cluster of their own.
2. Compute a mean for the each group.

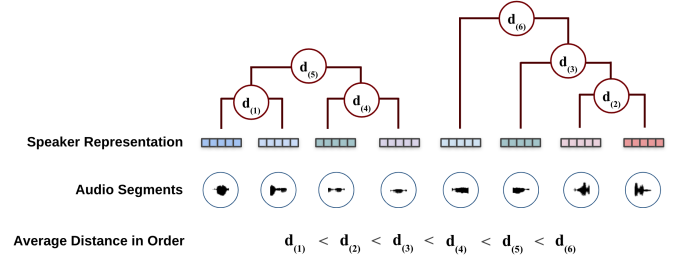


Fig. 4: Agglomerative Hierarchical Clustering.

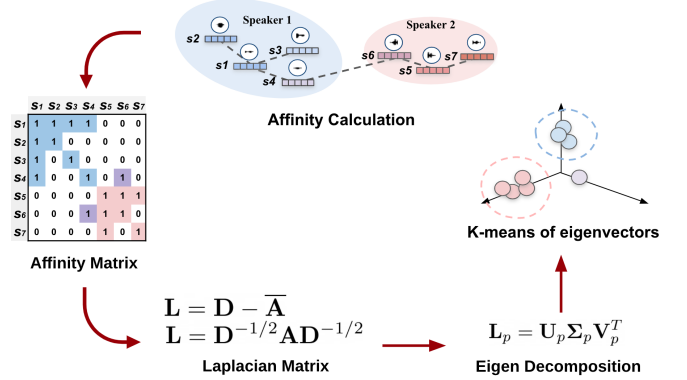


Fig. 5: General steps of spectral clustering.

3. Shift the search window to the new mean.
4. Repeat the process until convergence.

Mean-shift clustering algorithm was applied to speaker diarization task with KL distance [126], i-vector and cosine distance in [112, 127] and i-vector and PLDA [128]. The advantage of mean-shift clustering algorithm is that the clustering algorithm does not require the number of clusters in advance unlike k-means clustering methods. This becomes a significant advantage in speaker diarization tasks where the number of speakers is unknown as in most of the applications.

2.5.2. Agglomerative Hierarchical Clustering (AHC)

AHC is a clustering method that has been constantly employed in many speaker diarization systems with a number of different distance metric such as BIC [110, 129], KL [115] and PLDA [84, 90, 130]. AHC is an iterative process of merging the existing clusters until the clustering process meets a criterion. AHC process starts by calculating the similarity between N singleton clusters. At each step, a pair of clusters that has the highest similarity is merged. The iterative merging process of AHC produces a dendrogram which is depicted in Fig. 4.

One of the most important aspect of AHC is the stopping criterion. For speaker diarization task, AHC process can be stopped using either a similarity threshold or a target number of clusters. Ideally, if PLDA is employed as distance metric, the AHC process should be stopped at $s(\phi_1, \phi_2) = 0$ in Eq. (18). However, it is widely employed that the stopping metric is adjusted to get an accurate number of clusters based on a

development set. On the other hand, if the number of speakers is known or estimated by other methods, AHC process can be stopped when the clusters created by AHC process reaches the pre-determined number of speaker K .

2.5.3. Spectral Clustering

Spectral Clustering is another popular clustering approach for speaker diarization. While there are many variations, spectral clustering involves the following steps.

- i. Affinity Matrix Calculation: There are many ways to generate an affinity matrix \mathbf{A} depending on the way the affinity value is processed. The raw affinity value d is processed by kernel such as $\exp(-d^2/\sigma^2)$ where σ is a scaling parameter. On the other hand, the raw affinity value d could also be masked by zeroing the values below a threshold to only keep the prominent values.
- ii. Laplacian Matrix Calculation [131]: The graph Laplacian can be calculated in the following two types; normalized and unnormalized. The degree matrix \mathbf{D} contains diagonal elements $d_i = \sum_{j=1}^n a_{ij}$ where a_{ij} is the element of the i -th row and j -th column in an affinity matrix \mathbf{A} .
 - (a) Normalized Graph Laplacian:

$$\mathbf{L} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}. \quad (22)$$
 - (b) Unnormalized Graph Laplacian:

$$\mathbf{L} = \mathbf{D} - \mathbf{A}. \quad (23)$$
- iii. Eigen Decomposition: The graph Laplacian matrix \mathbf{L} is decomposed into the eigenvector matrix \mathbf{X} and the diagonal matrix that contains eigenvalues. Thus, $\mathbf{L} = \mathbf{X} \mathbf{\Lambda} \mathbf{X}^T$.
- iv. Re-normalization (optional) : the rows of \mathbf{X} is normalized so that $y_{ij} = x_{ij} / (\sum_j x_{ij}^2)^{1/2}$ where x_{ij} and y_{ij} are the elements of the i -th row and j -th column in matrix \mathbf{X} and \mathbf{Y} , respectively.
- v. Speaker Counting: Speaker number is estimated by finding the maximum eigengap.
- vi. k-means Clustering: The k -smallest eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ and the corresponding eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ are used to make $\mathbf{U} \in \mathbb{R}^{m \times n}$ where m is dimension of the row vectors in \mathbf{U} . Finally, the row vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$ are clustered by k-means algorithm.

Among many variations of spectral clustering algorithm, Ng-Jordan-Weiss (NJW) algorithm [132] is often employed for speaker diarization task. NJW algorithm employs a kernel $\exp(-d^2/\sigma^2)$ where d is a raw distance for calculating an affinity matrix. The affinity matrix is used for calculating a normalized graph Laplacian. In addition, NJW algorithm involves renormalization before the k-means clustering process. The speaker diarization system in [133] employed NJW algorithm

while choosing σ by using predefined scalar value β and variance values from the data points while the speaker diarization system in [134] did not use β value for NJW algorithm. On the other hand, in the speaker diarization system in [52], $\sigma^2 = 0.5$ for NJW algorithm.

Aside from NJW spectral clustering algorithm, many other types of spectral clustering were successfully applied to speaker diarization task. The speaker diarization system in [49] employed Gaussian blur for affinity values, diffusion process $\mathbf{Y} = \mathbf{X} \mathbf{X}^T$ and row-wise max normalization ($Y_{ij} = X_{ij} / \max_k X_{ik}$). In the spectral clustering approach appeared in [135], similarity values that are calculated from a neural network model were used without any kernel, and the unnormalized graph Laplacian is employed to perform spectral clustering. More recently, auto-tuning spectral clustering method was proposed for speaker diarization task [136] where the proposed clustering method does not require parameter tuning on a separate development set. The work in [136] employs binarized affinity matrix with the row-wise count p and the binarization parameter p is selected by choosing the minimum value of $r(p) = p/g_p$ where g_p represents the maximum eigengap from the unnormalized graph Laplacian matrix. Thus, $r(p)$ represents how clear the clusters are for the given value p and p could be automatically selected to perform spectral clustering without tuning the p -value.

2.6. Post-processing

2.6.1. Resegmentation

Resegmentation is a process to refine the speaker boundary that is roughly estimated by the clustering procedure. In [137], Viterbi resegmentation method based on the Baum-Welch algorithm was introduced. In this method, estimation of Gaussian mixture model corresponding to each speaker and Viterbi-algorithm-based resgmentation by using the estimated speaker GMM are alternately applied.

Later, a method to represent the diarization process based on Variational Bayesian Hidden Markov Model (VB-HMM) was proposed, and was shown to be superior as a resegmentation method compared to Viterbi resegmentation [138, 139, 140]. In the VB-HMM framework, the speech feature $\mathbf{X} = (\mathbf{x}_t | t = 1, \dots, T)$ is assumed to be generated from HMM where each HMM state corresponds to one of K possible speakers. Given we have M HMM states, M -dimensional variable $\mathbf{Z} = (\mathbf{z}_t | t = 1, \dots, T)$ is introduced where k -th element of \mathbf{z}_t is 1 if k -th speaker is speaking at the time index t , and 0 otherwise. At the same time, the distribution of \mathbf{x}_t is modeled based on a hidden variable $\mathbf{Y} = \{\mathbf{y}_k | k = 1, \dots, K\}$, where \mathbf{y}_k is a low dimensional vector for k -th speaker. Given these notation, the joint probability of \mathbf{X} , \mathbf{Y} , and \mathbf{Z} is decomposed as

$$P(\mathbf{X}, \mathbf{Z}, \mathbf{Y}) = P(\mathbf{X} | \mathbf{Z}, \mathbf{Y}) P(\mathbf{Z}) P(\mathbf{Y}), \quad (24)$$

where $P(\mathbf{X} | \mathbf{Z}, \mathbf{Y})$ is the emission probability modeled by GMM whose mean vector is represented by \mathbf{Y} , $P(\mathbf{Z})$ is the transition probability of the HMM, and $P(\mathbf{Y})$ is the prior distribution of \mathbf{Y} . Because \mathbf{Z} represents the trajectory of speakers, the diarization problem can be expressed as the inference problem of \mathbf{Z} that maximize the posterior distribution $P(\mathbf{Z} | \mathbf{X}) = \int P(\mathbf{Z}, \mathbf{Y} | \mathbf{X}) d\mathbf{Y}$.

Since it is intractable to directly solve this problem, Variational Bayes method is used to estimate the model parameters that approximate $P(\mathbf{Z}, \mathbf{Y}|\mathbf{X})$ [139, 141]. The VB-HMM framework was originally designed as a standalone diarization framework. However, it requires the parameter initialization to start VB estimation, and the parameters are usually initialized based on the result of speaker clustering. In that context, VB-HMM can be seen as a resegmentation method, and widely used as the final step of speaker diarization (e.g., [142, 113]).

2.6.2. System Fusion

As another direction of post processing, there have been a series of studies on the fusion method of multiple diarization results to improve the diarization accuracy. While it is widely known that the system combination generally yields better result for various systems (e.g., speech recognition [143] or speaker recognition [144]), combining multiple diarization hypotheses has several unique problems. Firstly, the speaker labeling is not standardized among different diarization systems. Secondly, the estimated number of speakers may differ among different diarization systems. Finally, the estimated time boundaries may be also different among multiple diarization systems. System combination methods for speaker diarization systems need to handle these problems during the fusion process of multiple hypotheses.

In [145], a method to select the best diarization result among multiple diarization systems were proposed. In this method, AHC is applied on the set of diarization results where the distance of two diarization results are measured by symmetric DER. AHC is executed until the number of groups becomes two, and the diarization result that has the smallest distance to all other results in the biggest group is selected as the final diarization result. In [146], two diarization systems are combined by finding the matching between two speaker clusters, and then performing the resegmentation based on the matching result.

More recently, DOVER (diarization output voting error reduction) method [147] was proposed to combine multiple diarization results based on the voting scheme. In the DOVER method, speaker labels among different diarization systems are aligned one by one to minimize DER between the hypotheses (the processes 2 and 3 of Fig. 6). After every hypotheses are aligned, each system votes its speaker label to each segmented region (each system may have different weight for voting), and the speaker label that gains the highest voting weight is selected for each segmented region (the process 4 of Fig. 6). In case of multiple speaker labels get the same voting weight, a heuristic to break the ties (such as selecting the result from the first system) is used.

The DOVER method has an implicit assumption that there is no overlapping speech, i.e., at most only 1 speaker is assigned for each time index. To combine the diarization hypotheses with overlapping speakers, two methods were recently proposed. In [104], the authors proposed the modified DOVER method, where the speaker labels in different diarization results are first aligned with a root hypothesis, and the speech activity of each speaker is estimated based on the weighted voting score for each speaker for each small segment. Raj et al. [148]

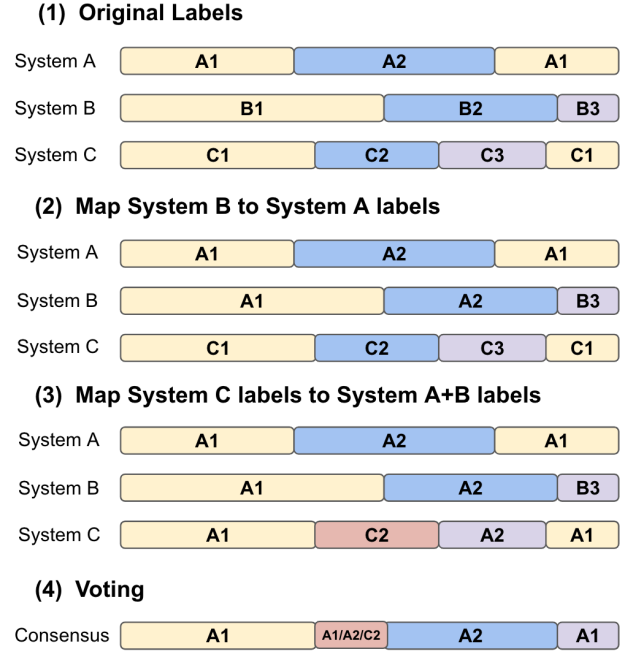


Fig. 6: Example of DOVER system.

proposed a method called DOVER-Lap, in which the speakers of multiple hypothesis are aligned by a weighted k-partite graph matching, and the number of speakers K for each small segment is estimated based on the weighted average of multiple systems to select the top- K voted speaker labels. Both the modified DOVER and DOVER-Lap showed the improvement of DER for the speaker diarization result with speaker overlaps.

3. Recent Advances in Speaker Diarization using Deep Learning

This section introduces various recent efforts toward deep learning-based speaker diarization techniques. Firstly, methods that incorporate deep learning into a single component of speaker diarization, such as clustering or post-processing, are introduced in Section 3.1. Then, methods that unify several components of speaker diarization into a single neural network are introduced in Section 3.2.

3.1. Single-module optimization

3.1.1. Speaker clustering enhanced by deep learning

Several methods that enhance the speaker clustering based on deep learning were proposed. A deep-learning based clustering algorithm, called Improved Deep Embedded Clustering (IDEC) is proposed in [149]. The goal is to transform the input features, herein speaker embeddings, to become more separable, given the number of clusters/speakers. The key idea is that each embedding has a probability of “belonging” to each of the available speaker cluster [150, 64],

$$q_{ij} = \frac{(1 + \|z_i - \mu_j\|^2/a)^{-\frac{a+1}{a}}}{\sum_l (1 + \|z_i - \mu_l\|^2/a)^{-\frac{a+1}{a}}}, \quad p_{ij} = \frac{q_{ij}^2/f_i}{\sum_l q_{il}^2/f_l} \quad (25)$$

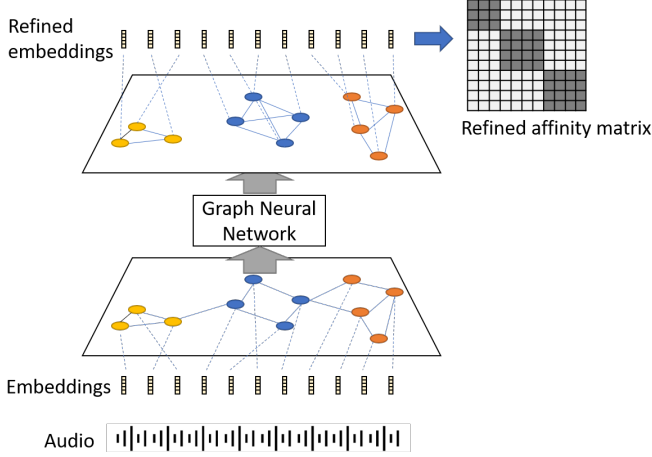


Fig. 7: Speaker diarization with graph neural network

where z_i are the bottleneck features, μ_i is the centroid of i -th cluster and f_i is the soft cluster frequency with $f_i = \sum q_{ij}$. The clusters are iteratively refined based on a target distribution [150] based on bottleneck features estimated using an autoencoder.

The initial DEC approach presented some problems. As such, improved versions of the algorithm have been proposed, where the possibility of trivial (empty) clusters is addressed (under the assumption that the distribution of speaker turns is uniform across all speakers, i.e. all speakers contribute equally to the session). This assumption is not realistic in real meeting environments but it constrains the solution space enough to avoid the empty clusters without affecting overall performance. An additional loss term penalizes the distance from the centroids μ_i , bringing the behavior of the algorithm closer to k-means [149].

Based on these improvements, the loss function of the revisited DEC algorithm consists of three different loss components, i.e. L_c the clustering error, L_u the uniform “speaker air-time” distribution constraint and L_{MSE} the distance of the bottleneck features from the centroids [149],

$$L = \alpha L_c + \beta L_r + \gamma L_u + \delta L_{MSE} \quad (26)$$

allowing for the different loss functions to be weighted differently and the weights α, β, γ and δ can be fine-tuned on some held-out data.

In [65], a different approach that purify the similarity matrix for the spectral clustering based on the graph neural network (GNN) was proposed (Fig. 7). Given a sequence of speaker embeddings $\{\mathbf{e}_1, \dots, \mathbf{e}_N\}$ where N is the length of sequence. The first layer of the GNN takes the input $\{\mathbf{x}_i^0 = \mathbf{e}_i | i = 1, \dots, N\}$. The GNN then computes the output of the p -th layer $\{\mathbf{x}_i^{(p)} | i = 1, \dots, N\}$ as followings.

$$\mathbf{x}_i^{(p)} = \sigma(\mathbf{W} \sum_j \mathbf{L}_{i,j} \mathbf{x}_j^{(p-1)}), \quad (27)$$

where \mathbf{L} represents a normalized affinity matrix added by self-connection, \mathbf{W} is a trainable weight matrix for the p -th layer,

and $\sigma(\cdot)$ is a nonlinear function. GNN was optimized by minimizing the distance between the reference affinity matrix and estimated affinity matrix, where the distance was calculated by a combination of histogram loss [151] and nuclear norm.

There are also several different approaches to generate the affinity matrix. In [152], self-attention-based network was introduced to directly generate a similarity matrix from a sequence of speaker embeddings. In [153], several affinity matrices with different temporal resolutions were fused into single affinity matrix based on a neural network.

3.1.2. Learning the Distance Estimator

Data-driven techniques perform remarkably well on a wide variety of tasks [154]. However, traditional DL architectures may fail when the problem involves relational information between observations [155]. Recently, Relational Recurrent Neural Networks (RRNN) were introduced by [155, 156, 157] to solve this “relational information learning” task. Speaker diarization can be seen as a member of this class of tasks, since the final decision depends on the distance relations between speech segments and speaker profiles or centroids.

The challenges of audio segmentation are detailed in Section 2.3. Further, speaker embeddings are usually extracted from a network trained to distinguish speakers among thousands of candidates [50]. However, a different level of granularity in the speaker space is required, since only a small number of participants is typically involved in an interactive meeting scenario. In addition to that, the distance metric used is often heuristic and/or dependent on certain assumptions which do not necessarily hold, e.g., assuming Gaussianity in the case of PLDA [158], etc. Finally, the audio chunks are treated independently and any temporal information about the past and future is simply ignored. Most of these issues can be addressed with the RRNNs in [159], where a data-driven, memory-based approach is bridging the performance gap between the heuristic and the trainable distance estimating approaches. The RRNNs have shown great success on several problems requiring relational reasoning [156, 155, 159], and specifically using the Relational Memory Core (RMC) [155].

In this context, a novel approach of learning the distance between such centroids (or speaker profiles) and the embeddings was proposed in [159] (Fig. 8). The diarization process can be seen as a classification task on already segmented audio, Section 2.3, where the audio signal is first segmented either uniformly [160] or based on estimated speaker change points [161]. As these segments are assumed to be speaker-homogeneous, speaker embeddings x_j for each segment are extracted and then compared against all the available speaker profiles or speaker centroids. By minimizing a particular distance metric, the most suitable speaker label is assigned to the segment. The final decision relies on a distance estimation, either the cosine [51] or the PLDA [158] distance, or the distance based on RRNNs as proposed in [159]. The later method based on memory networks has shown consistent improvements in performance.

3.1.3. Deep learning-based post processing

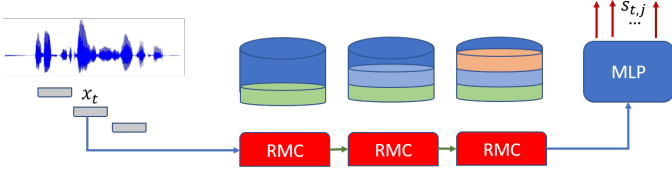


Fig. 8: Continuous speaker identification system based on RMC. The speech signal is segmented uniformly and each segment x_t is compared against all the available speaker profiles according to a distance metric $d(\cdot, \cdot)$. A speaker label $s_{t,j}$ is assigned to each x_t minimizing this metric.

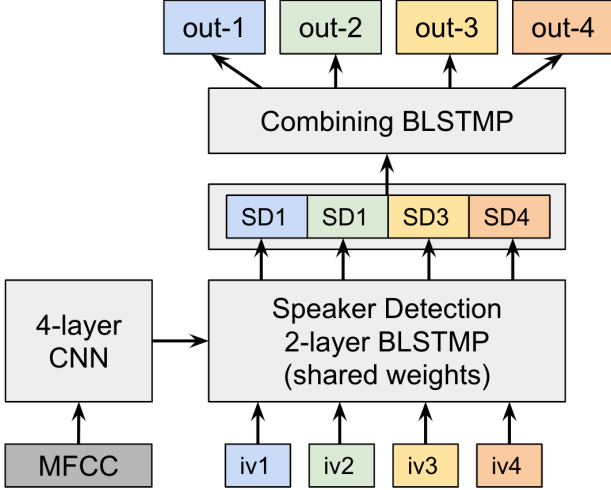


Fig. 9: Target Speaker Voice Activity Detector

There are a few recent studies to train a neural network that is applied on top of the result of a clustering-based speaker diarization. These method can be categorized as an extension of the post processing.

Medennikov et al. proposed the Target-Speaker Voice Activity Detection (TS-VAD) to achieve accurate speaker diarization even with many speaker overlaps noisy conditions [91, 66]. As shown in Fig. 9, TS-VAD takes the input of acoustic feature (MFCC) as well as the i-vector of all target speakers. The model has an output layer where i -th element becomes 1 at time frame t if i -th speaker is speaking at the time frame, and 0 otherwise. To convert the raw output into a sequence of segment, a further post-processing based on heuristics (median filtering, binarization with the threshold, etc.) or HMM-based decoding with states representing silence, non-overlapping speech of each speaker, and overlapping speech from all possible pairs of speakers is used. Prior to inference, TS-VAD requires the i-vector of all target speakers. The i-vectors are initialized based on the conventional clustering-based speaker diarization result. After initializing the i-vector, the inference by TS-VAD and refinement of i-vector based on the TS-VAD result can be repeated until it converges. TS-VAD showed a significantly better DER compared with the conventional clustering based approach [91, 88]. On the other hand, it has a constraint that the maximum number of speakers that the model can handle is limited by the number of element of the output layer.

As a different approach, Horiguchi et al. proposed to apply the EEND model (detailed in Section 3.2.4) to refine the result of a clustering-based speaker diarization [162]. A clustering-based speaker diarization method can handle a large number of speakers while it is not good at handling the overlapped speech. On the other hand, EEND has the opposite characteristics. To complementary use two methods, they first apply a conventional clustering method. Then, the two-speaker EEND model is iteratively applied for each pair of detected speakers to refine the time boundary of overlapped regions.

3.2. Joint optimization for speaker diarization

3.2.1. Joint segmentation and clustering

A model called Unbounded Interleaved-State Recurrent Neural Networks (UIS-RNN) was proposed that replaces the segmentation and clustering procedure into a trainable model [55]. Given the input sequence of embeddings $\mathbf{X} = (\mathbf{x}_t \in \mathbb{R}^d | t = 1, \dots, T)$, UIS-RNN generates the diarization result $\mathbf{Y} = (y_t \in \mathbb{N} | t = 1, \dots, T)$ as a sequence of speaker index for each time frame. The joint probability of \mathbf{X} and \mathbf{Y} can be decomposed by the chain rule as follows.

$$P(\mathbf{X}, \mathbf{Y}) = P(\mathbf{x}_1, y_1) \prod_{t=2}^T P(\mathbf{x}_t, y_t | \mathbf{x}_{1:t-1}, y_{1:t-1}). \quad (28)$$

To model the distribution of speaker change, UIS-RNN then introduce a latent variable $\mathbf{Z} = (z_t \in \{0, 1\} | t = 2, \dots, T)$, where z_t becomes 1 if the speaker indices at time $t-1$ and t are different, and 0 otherwise. The joint probability including \mathbf{Z} is then decomposed as follows.

$$P(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = P(\mathbf{x}_1, y_1) \prod_{t=2}^T P(\mathbf{x}_t, y_t, z_t | \mathbf{x}_{1:t-1}, y_{1:t-1}, z_{1:t-1}) \quad (29)$$

Finally, the term $P(\mathbf{x}_t, y_t, z_t | \mathbf{x}_{1:t-1}, y_{1:t-1}, z_{1:t-1})$ is further decomposed into three components.

$$P(\mathbf{x}_t, y_t, z_t | \mathbf{x}_{1:t-1}, y_{1:t-1}, z_{1:t-1}) = P(\mathbf{x}_t | \mathbf{x}_{1:t-1}, y_{1:t}) P(y_t | z_t, y_{1:t-1}) P(z_t | z_{1:t-1}) \quad (30)$$

Here, $P(\mathbf{x}_t | \mathbf{x}_{1:t-1}, y_{1:t})$ represents the sequence generation probability, and modeled by gated recurrent unit (GRU)-based recurrent neural network. $P(y_t | z_t, y_{1:t-1})$ represents the speaker assignment probability, and modeled by a distant dependent Chinese restaurant process [163], which can model the distribution of unbounded number of speakers. Finally, $P(z_t | z_{1:t-1})$ represents the speaker change probability, and modeled by Bernoulli distribution. Since all models are represented by trainable models, the UIS-RNN can be trained in a supervised way by finding parameters that maximizes $\log P(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$ over training data. The inference can be conducted by finding \mathbf{Y} that maximizes $\log P(\mathbf{X}, \mathbf{Y})$ given \mathbf{X} based on the beam search in an online fashion. While UIS-RNN works in an online fashion, UIS-RNN showed better DER than that of the offline system based on the spectral clustering.

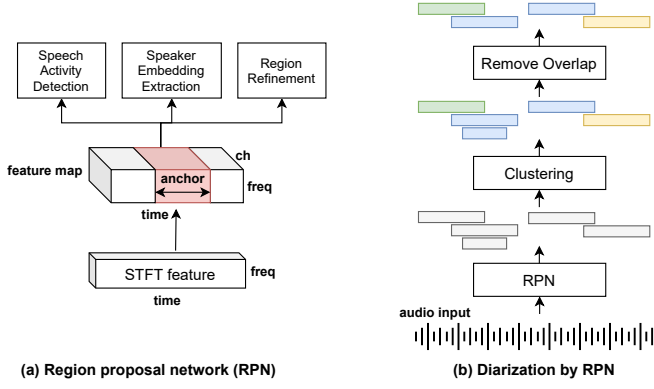


Fig. 10: (a) RPN for speaker diarization, (b) diarization procedure based on RPN.

3.2.2. Joint segmentation, embedding extraction, and re-segmentation

A speaker diarization method based on the Region Proposal Networks (RPN) was proposed to jointly perform segmentation, speaker embedding extraction, and re-segmentation procedures by a single neural network [75]. The RPN was originally proposed to detect multiple objects from a 2-d image [164], and 1-d variant of the RPN is used for speaker diarization along with the time-axis. RPN works on the Short-Term Fourier Transform (STFT) features. A neural network converts the STFT feature into the feature map (Fig. 10 (a)). Then, for each candidates of time region of speech activity, called an anchor, the neural network jointly perform three tasks to (i) estimate whether the anchor includes speech activity or not, (ii) extract a speaker embedding corresponding to the anchor, and (iii) estimate the difference of the duration and center position of the anchor and the reference speech activity. The first, second, and third tasks corresponds to the segmentation, speaker embedding extraction, and re-segmentation, respectively.

The inference procedure by RPN is depicted in Fig. 10 (b). The RPN is firstly applied to every anchors on the test audio, and the regions with speech activity probability higher than a pre-determined threshold are listed as a candidate time regions. Estimated regions are then clustered by using a conventional clustering method (e.g., k-means) based on the speaker embeddings corresponding to each region. Finally, a procedure called non-maximum suppression is applied to remove highly overlapped segments.

The RPN-based speaker diarization has the advantage that it can handle overlapped speech with possibly any number of speakers. Also, it is much simpler than the conventional speaker diarization system. It was shown in multiple dataset that the RPN-based speaker diarization system achieved significantly better DER than the conventional clustering-based speaker diarization system [75, 88].

3.2.3. Joint speech separation and diarization

There are also recent researches to jointly perform speech separation and speaker diarization. Kounades-Bastian et al. [165, 166] proposed to incorporate a speech activity model into

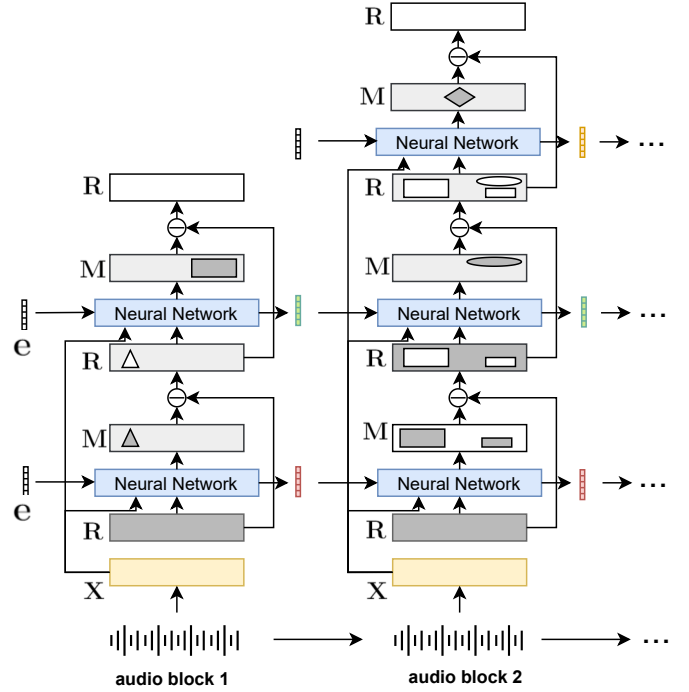


Fig. 11: Joint speech separation, speaker counting, and speaker diarization model.

speech separation based on the spatial covariance model with non-negative matrix factorization. They derived the EM algorithm to estimate separated speech and speech activity of each speaker from the multi-channel overlapped speech. While their method jointly perform speaker diarization and speech separation, their method is based on a statistical modeling, and estimation was conducted solely based on the observation, i.e. without any model training.

Neumann et al. [76, 167] later proposed a trainable model, called online Recurrent Selective Attention Network (online RSAN), for joint speech separation, speaker counting, and speaker diarization based on a single neural network (Fig. 11). Their neural network takes the input of spectrogram $X \in \mathbb{R}^{T \times F}$, a speaker embedding $e \in \mathbb{R}^d$, and a residual mask $R \in \mathbb{R}^{T \times F}$, where T and F is the maximum time index and the maximum frequency bin of the spectrogram, respectively. It output the speech mask $M \in \mathbb{R}^{T \times F}$ and an updated speaker embedding for the speaker corresponding to e . The neural network is firstly applied with R whose element is all 1, and e whose element is all 0. After the first inference of M , R is updated as $R \leftarrow \max(R - M, 0)$, and the neural network is again applied with the updated R . This procedure is repeated until sum of R becomes less than a threshold. A separated speech can be obtained by $M \odot X$ where \odot is the element-wise multiplication. The speaker embedding is used to keep track the speaker of adjacent blocks. Thanks to the iterative approach, this neural network can cope with variable number of speakers while jointly performing speech separation and speaker diarization.

3.2.4. Fully end-to-end neural diarization

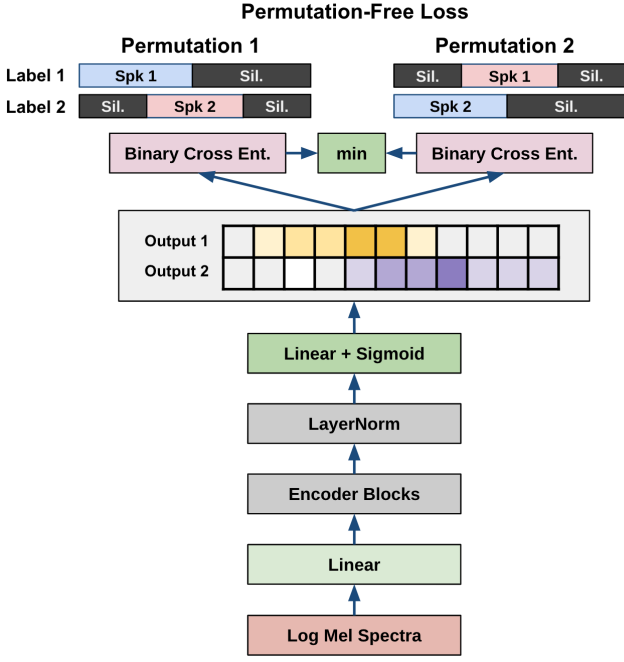


Fig. 12: Two-speaker end-to-end neural diarization model

Recently, the framework called End-to-End Neural Diarization (EEND) was proposed [56, 57], which performs all speaker diarization procedure based on a single neural network. The architecture of EEND is shown in Fig. 12. An input to the EEND model is a T -length sequence of acoustic features (e.g., log mel filterbank), $\mathbf{X} = (\mathbf{x}_t \in \mathbb{R}^F | t = 1, \dots, T)$. A neural network then outputs the corresponding speaker label sequence $\mathbf{Y} = (\mathbf{y}_t | t = 1, \dots, T)$ where $\mathbf{y}_t = [y_{t,k} \in \{0, 1\} | k = 1, \dots, K]$. Here, $y_{t,k} = 1$ represents the speech activity of the speaker k at the time frame t , and K is the maximum number of speakers that the neural network can output. Importantly, $y_{t,k}$ and $y_{t,k'}$ can be both 1 for different speakers k and k' , which represents that two speakers k and k' is speaking simultaneously (i.e. overlapping speech). The neural network is trained to maximize $\log P(\mathbf{Y}|\mathbf{X}) \sim \sum_t \sum_k \log P(y_{t,k}|\mathbf{X})$ over the training data by assuming the conditional independence of the output $y_{t,k}$. Because there can be multiple candidates of the reference label \mathbf{Y} by swapping the speaker index k , the loss function is calculated for all possible reference labels and the reference label that has the minimum loss is used for the error back-propagation, which is inspired by the permutation free objective used in speech separation [59]. EEND was initially proposed with a bidirectional long short-term memory (BLSTM) network [56], and was soon extended to the self-attention-based network [57] by showing the state-of-the-art DER for CALLHOME dataset (LDC2001S97) and Corpus of Spontaneous Japanese [168].

There are multiple advantages of EEND. Firstly, it can handle overlapping speech in a sound way. Secondly, the network is directly optimized towards maximizing diarization accuracy, by which we can expect a high accuracy. Thirdly, it can be retrained by a real data (i.e. not synthetic data) just by feeding a reference diarization label while it is often not strait-

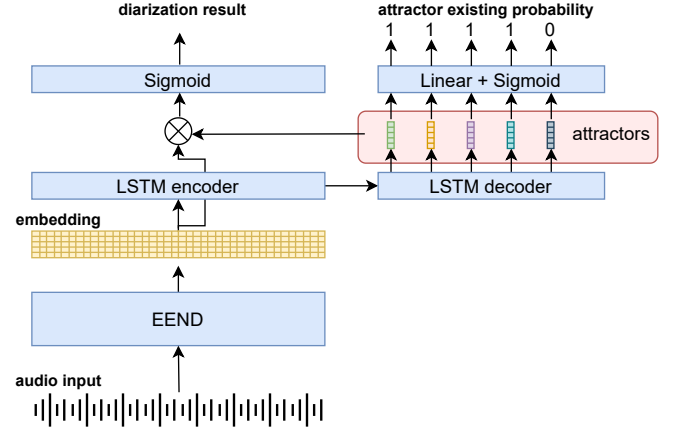


Fig. 13: EEND with encoder-decoder-based attractor (EDA).

forward for the prior works. On the other hand, several limitations are also known for EEND. Firstly, the model architecture constrains the maximum number of speakers that the model can cope with. Secondly, EEND consists of BLSTM or self-attention neural networks, which makes it difficult to do online processing. Thirdly, it was empirically suggested that EEND tends to overfit to the distribution of the training data [56].

To cope with an unbounded number of speakers, several extensions of EEND have been investigated. Horiguchi et al. [169] proposed an extension of EEND with the encoder-decoder-based attractor (EDA) (Fig. 13). This method applies an LSTM-based encoder-decoder on the output of EEND to generate multiple attractors. Attractors are generated until the attractor existing probability becomes less than a threshold. Then, each attractor is multiplied with the embeddings generated from EEND to calculate the speech activity for each speaker. On the other hand, Fujita et al. [170] proposed another approach to output the speech activity one after another by using a conditional speaker chain rule. In this method, a neural network is trained to produce a posterior probability $P(\mathbf{y}_k | \mathbf{y}_1, \dots, \mathbf{y}_{k-1}, \mathbf{X})$, where $\mathbf{y}_k = (y_{t,k} \in \{0, 1\} | t = 1, \dots, T)$ is the speech activity for k -th speaker. Then, the joint speech activity probability of all speakers can be estimated from the following speaker-wise conditional chain rule as:

$$P(\mathbf{y}_1, \dots, \mathbf{y}_K | \mathbf{X}) = \prod_{k=1}^K P(\mathbf{y}_k | \mathbf{y}_1, \dots, \mathbf{y}_{k-1}, \mathbf{X}). \quad (31)$$

During inference, the neural network is repeatedly applied until the speech activity y_k for the last estimated speaker approaches zero. Kinoshita et al. [171] proposed a different approach that combines EEND and speaker clustering. In their method, a neural network is trained to generate speaker embeddings as well as the speech activity probability. Speaker clustering constrained by the estimated speech activity by EEND is applied to align the estimated speakers among different processing blocks.

There are also a few recent trials to extend the EEND for online processing. Xue et al. [172] proposed a method with a speaker tracing buffer to better align the speaker labels of adjacent processing blocks. Han et al. [173] proposed a block on-

line version of EDA-EEND [169] by carrying the hidden state of the LSTM-encoder to generate attractors block by block.

4. Speaker Diarization in the context of ASR

From a conventional perspective, speaker diarization is considered a pre-processing step for ASR. In the traditional system structure for speaker diarization as depicted in Fig. 1, speech inputs are processed sequentially across the diarization components without considering the ASR objective, which corresponds to minimize word error rate (WER). One issue is that the tight boundaries of speech segments as the outcomes of speaker diarization have a high chance of causing unexpected word truncation or deletion errors in ASR decoding. In this section we discuss how speaker diarization systems have been developed in the context of ASR, not only resulting in better WER by preventing speaker diarization from hurting ASR performance, but also benefiting from ASR artifacts to enhance diarization performance. More recently, there have been a few pioneering proposals made for joint modeling of speaker diarization and ASR, which we will introduce in the section as well.

4.1. Early Works

The lexical information from ASR output has been employed for speaker diarization system in a few different ways. First, the earliest approach was RT03 evaluation [1] which used word boundary information for segmentation purpose. In [1], a general ASR system for broadcast news data was built where the basic components are segmentation, speaker clustering, speaker adaptation and system combination after ASR decoding from the two sub-systems with the different adaptation methods. To understand the impact of the word boundary information, they used ASR outputs to replace the segmentation part and compared the diarization performance of the each system. In addition, ASR result was also used for refining SAD in IBM’s submission [174] for RT07 evaluation. The system appeared in [174] incorporates word alignments from speaker independent ASR module and refines SAD result to reduce false alarms so that the speaker diarization system can have better clustering quality. The segmentation system in [175] also takes advantage of word alignments from ASR. The authors in [175] focused on the word-breakage problem where the words from ASR output are truncated by segmentation results since segmentation result and decoded word sequence are not aligned. Therefore, word-breakage (WB) ratio was proposed to measure the rate of change-points that are detected inside intervals corresponding to words. The DER and WB were reported together to measure the influence of word truncation problem. While the fore-mentioned early works of speaker diarization systems that are leveraging ASR output are focusing on the word alignment information to refine the SAD or segmentation result, the speaker diarization system in [176] created a dictionary for the phrases that commonly appear in broadcast news. The phrases in this dictionary provide identity of who is speaking, who will speak and who spoke in the broadcast news scenario. For example,

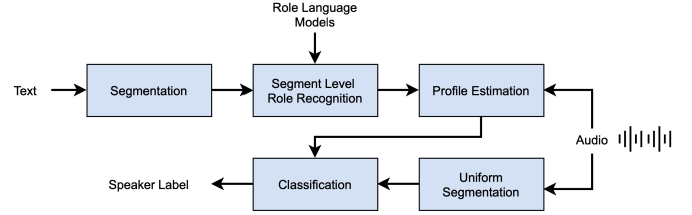


Fig. 14: Integration of lexical information and acoustic information.

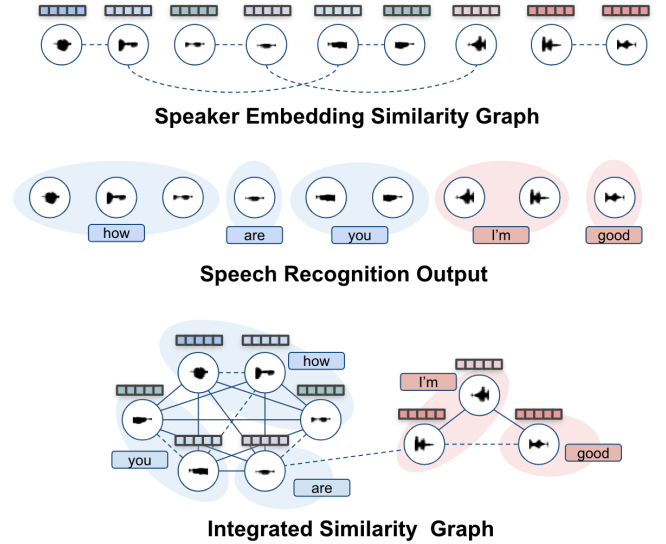


Fig. 15: Integration of lexical information and acoustic information.

“This is [name]” indicates who was the speaker of the broadcast news section. Although the early speaker diarization studies did not fully leverage the lexical information to drastically improve DER, the idea of integrating the information from ASR output has been employed by many studies to refine or improve the speaker diarization output.

4.2. Using lexical information from ASR

The more recent speaker diarization systems that take advantage of the ASR transcript have employed a DNN model to capture the linguistic pattern in the given ASR output to enhance the speaker diarization result. The authors in [177] proposed a way of using the linguistic information for the speaker diarization task where participants have distinct roles that are known to the speaker diarization system. Fig. 14 shows the diagram of speaker diarization system appeared in [177]. In this system, a neural text-based speaker change detector and a text-based role recognizer are employed. By employing both linguistic and acoustic information, DER was significantly improved compared to the acoustic only system.

Lexical information from ASR output was also utilized for speaker segmentation [178] by employing a sequence to sequence model that outputs speaker turn tokens. Based on the estimated speaker turn, the input utterance is segmented accordingly. The experimental results in [178] show that using both

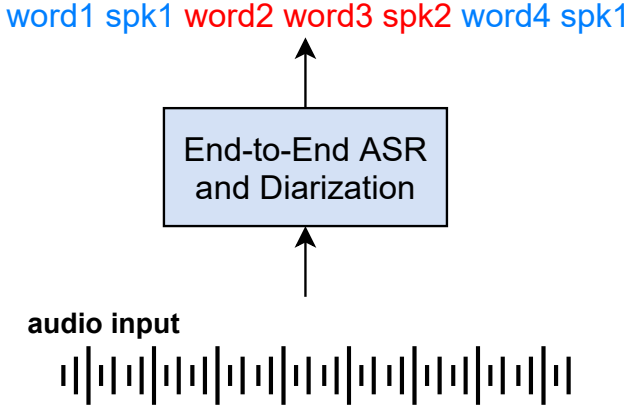


Fig. 16: Joint ASR and diarization by inserting a speaker tag in the transcription.

acoustic and lexical information can get an extra advantage owing to the word boundaries we get from the ASR output.

[179] presented follow-up research within the above thread. Unlike the system in [178], lexical information from the ASR module was integrated with the speech segment clustering process by employing an integrated adjacency matrix. The adjacency matrix is obtained from max operation between acoustic information created from affinities among audio segments and lexical information matrix created by segmenting the word sequence into word chunks that are likely to be spoken by the same speaker. Fig. 15 shows a diagram that explains how lexical information is integrated in an affinity matrix with acoustic information. The integrated adjacency matrix leads to an improved speaker diarization performance for CALLHOME American English dataset.

4.3. Joint ASR and speaker diarization with deep learning

Motivated by the recent success of deep learning and end-to-end modeling, several models have been proposed to jointly perform ASR and speaker diarization. As with the previous section, ASR results contain a strong cue to improve speaker diarization. On the other hand, speaker diarization results can be used to improve the ASR accuracy, for example, by adapting the ASR model towards each estimated speaker. Joint modeling can leverage such inter-dependency to improve both ASR and speaker diarization. In the evaluation, a word error rate (WER) metric that is affected by both ASR errors and speaker attribution errors, such as speaker-attributed WER [180] or cpWER [89], is often used. ASR-specific metrics (e.g., speaker-agnostic WER) or diarization-specific metrics (e.g., DER) is also used complementary.

A first line of approaches is introducing a speaker tag in the transcription of end-to-end ASR models (Fig. 16). Shafey et al. [77] proposed to insert a speaker role tag (e.g., <doctor> and <patient>) in the output of a recurrent neural network-transducer (RNN-T)-based ASR system. Similarly, Mao et al. [78] proposed to insert a speaker identity tag in the output of an attention-based encoder-decoder ASR system. These method have been shown to be able to perform both ASR and speaker

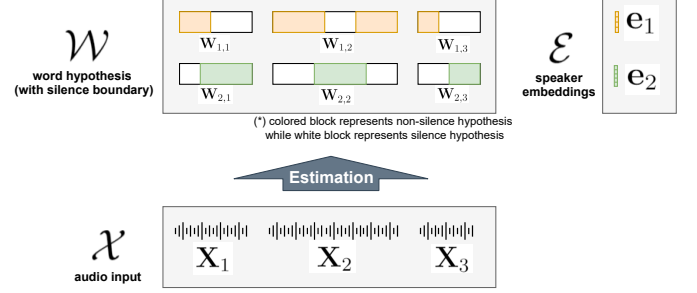


Fig. 17: Joint decoding framework for ASR and speaker diarization.

diarization in their experiments. On the other hand, the speaker roles or speaker identity tags needs to be determined and fixed during training, so it is difficult to cope with an arbitrary number of speakers with this approach.

A second approach is a MAP-based joint decoding framework. Kanda et al. [79] formulated the joint decoding of ASR and speaker diarization as followings (see also Fig. 17). Assume that a sequence of observations is represented by $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_U\}$, where U is the number of segments (e.g., generated by applying VAD on a long audio) and \mathbf{X}_u is the acoustic feature sequence of the u -th segment. Further assume that word hypotheses with time boundary information is represented by $\mathcal{W} = \{\mathbf{W}_1, \dots, \mathbf{W}_U\}$ where \mathbf{W}_u is the speech recognition hypotheses corresponding to the segment u . Here, $\mathbf{W}_u = (\mathbf{W}_{1,u}, \dots, \mathbf{W}_{K,u})$ contains all speakers' hypotheses in the segment u where K is the number of speakers, and $\mathbf{W}_{k,u}$ represents the speech recognition hypothesis of the k -th speaker of the segment u . Finally, a tuple of speaker embeddings $\mathcal{E} = (\mathbf{e}_1, \dots, \mathbf{e}_K)$, where $\mathbf{e}_j \in \mathbb{R}^d$ is d -dimensional speaker embeddings of k -th speaker, is also assumed. With all these notations, the joint decoding framework of multi-speaker ASR and diarization can be formulated as a problem to find most likely $\hat{\mathcal{W}}$ as,

$$\hat{\mathcal{W}} = \underset{\mathcal{W}}{\operatorname{argmax}} P(\mathcal{W}|\mathcal{X}) \quad (32)$$

$$= \underset{\mathcal{W}}{\operatorname{argmax}} \left\{ \sum_{\mathcal{E}} P(\mathcal{W}, \mathcal{E}|\mathcal{X}) \right\} \quad (33)$$

$$\approx \underset{\mathcal{W}}{\operatorname{argmax}} \left\{ \max_{\mathcal{E}} P(\mathcal{W}, \mathcal{E}|\mathcal{X}) \right\}, \quad (34)$$

where we use the Viterbi approximation to obtain the final equation. This maximization problem is further decomposed into two iterative problems as,

$$\hat{\mathcal{W}}^{(i)} = \underset{\mathcal{W}}{\operatorname{argmax}} P(\mathcal{W}|\hat{\mathcal{E}}^{(i-1)}, \mathcal{X}), \quad (35)$$

$$\hat{\mathcal{E}}^{(i)} = \underset{\mathcal{E}}{\operatorname{argmax}} P(\mathcal{E}|\hat{\mathcal{W}}^{(i)}, \mathcal{X}), \quad (36)$$

where i is the iteration index of the procedure. In [79], Eq. (35) is modeled by the target speaker ASR [181, 182, 183, 71] and Eq. (36) is modeled by the overlap-aware speaker embedding estimation. This method shows a similar speaker-attributed WER compared to that of the target speaker ASR with oracle speaker embeddings. On the other hand, it requires an iterative

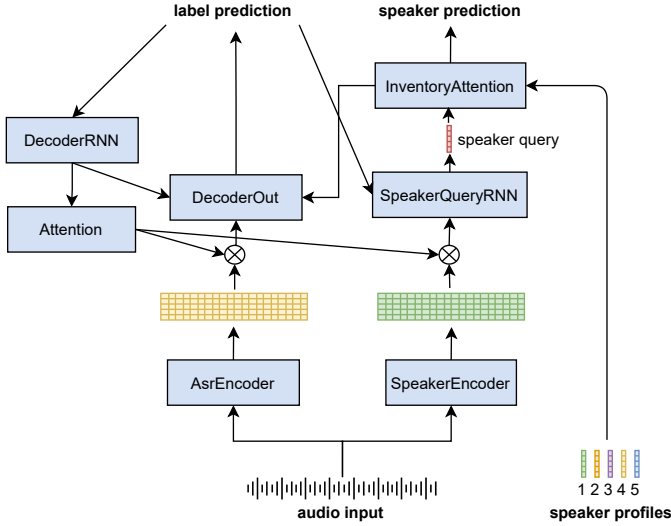


Fig. 18: End-to-end speaker-attributed ASR

application of the target-speaker ASR and speaker embedding extraction, which makes it challenging to apply the method in online mode.

As a third line of approaches, End-to-End (E2E) Speaker-Attributed ASR (SA-ASR) model was recently proposed to jointly perform speaker counting, multi-speaker ASR, and speaker identification [184, 185]. Different from the first two approaches, the E2E SA-ASR model takes the additional input of speaker profiles and identifies the index of speaker profiles based on the attention mechanism (Fig. 18). Thanks to the attention mechanism for speaker identification and multi-talker ASR capability based on serialized output training [186], there is no limitation of a maximum number of speakers that the model can cope with. In case relevant speaker profiles are supplied in the inference, the E2E SA-ASR model can automatically transcribe the utterance while identifying the speaker of each utterance based on the supplied speaker profiles. On the other hand, in case the relevant speaker profiles cannot be used prior to the inference, the E2E SA-ASR model can still be applied with example profiles, and speaker clustering on the internal speaker embeddings of the E2E SA-ASR model (“speaker query” in Fig. 18) is used to diarize the speaker [80].

5. Evaluation of Speaker Diarization

This section describes the evaluation scheme for speaker diarization. The dataset that is widely used for the evaluation of speaker diarization is first introduced in Section 5.1. Then, the evaluation metric for speaker diarization is introduced in Section 5.2. Finally, international efforts to evaluate diarization systems are introduced in Section 5.3. The summary of the dataset is shown in Table 2.

5.1. Diarization Evaluation Datasets

5.1.1. CALLHOME: NIST SRE 2000 (LDC2001S97)

NIST SRE 2000 (Disk-8), often referred to as CALLHOME dataset, has been the most widely used dataset for speaker di-

arization in the recent papers. CALLHOME dataset contains 500 sessions of multilingual telephonic speech. Each session has 2 to 7 speakers while there are two dominant speakers in each conversation.

5.1.2. AMI

The AMI database [187] includes 100 hours of meeting recordings from multiple sites in 171 meeting sessions. AMI database provides audio source recorded with lapel microphones which are separately recorded and amplified for each speaker. Another audio source is recorded with omnidirectional microphone arrays that are mounted on the table while meeting. AMI database is a suitable dataset for evaluating speaker diarization system integrated with ASR module since AMI provides forced alignment data which contains word and phoneme level timings along with the transcript and speaker label. Each meeting session contains 3 to 5 speakers.

5.1.3. ICSI meeting Corpus

The ICSI meeting corpus [188] contains 75 meeting corpus with 4 meeting types. ICSI meeting corpus provides word level timing along with the transcript and speaker label. The audio source is recorded with close-talking individual microphone and six tabletop microphones to provide speaker-specific channel and multi-channel recording. Each meeting has 3 to 10 participants.

5.1.4. DIHARD Challenge dataset

DIHARD challenge dataset is created for DIHARD challenge 1, 2 and 3 [189, 85, 190] while focusing on very challenging domains. DIHARD challenge development set and evaluation set include clinical interviews, web videos, speech in the wild (e.g., recordings in restaurants). DIHARD challenge dataset also includes relatively less challenging datasets such as conversational telephonic speech (CTS) and audio books to diversify the domains in development set and evaluation set. Contrary to other speaker diarization datasets, domains such as restaurant conversation and web videos contain significantly lower signal to noise ratio (SNR) that makes DER way higher. The first DIHARD challenge, DIHARD 1, started with track 1 for diarization beginning from oracle SAD and track 2 diarization from scratch using system SAD. Unlike DIHARD 1, DIHARD 2 included multichannel speaker diarization task in track 3 (oracle SAD) and track 4 (system SAD) adding the recordings drawn from CHiME-5 corpus [191]. In the latest DIHARD challenge, DIHARD 3, CTS dataset was added to DIHARD 3 dev set and eval set and DIHARD 3 removed track 3 and track 4 while keeping only track 1 (oracle SAD) and track 2 (system SAD).

5.1.5. CHiME-5/6 challenge corpus

The CHiME-5 corpus [191] includes 50 hours of multi-party real conversations in the every-day home environment. It contains speaker labels, segmentation, and corresponding transcriptions. All of them are manually annotated. The audio source is recorded by multiple 4-channel microphone arrays located in the kitchen and dining/living rooms in a house, and also

Table 2: Diarization Evaluation Datasets

	Size (hr)	Style	# speakers
CALLHOME	20	Conversation	2–7
AMI	100	Meeting	3–5
ICSI meeting	72	Meeting	3–10
DIHARD I Track 1,2	19(dev), 21(eval)	Miscellaneous	1–7
DIHARD II Track 1,2	24(dev), 22(eval)	Miscellaneous	1–8
DIHARD II Track 3,4	262(dev), 31(eval)	Miscellaneous	4
DIHARD III Track 1,2	34(dev), 33(eval)	Miscellaneous	1–7
CHiME-5/6	50	Conversation	4
VoxConverse	74	YouTube video	1–21
LibriCSS	10	Read speech	8

recorded by binaural microphones worn by participants. The number of participants is fixed as four. The CHiME-6 challenge uses the same CHiME-5 corpus, but track 2 includes the speaker diarization problem in the challenge (i.e., no speaker labels and segmentation are given). The CHiME-5 corpus was also used as one track in the DIHARD 2 challenge.

5.1.6. VoxConverse

The VoxConverse dataset [192] contains 74 hours of human conversation extracted from YouTube video. The dataset is divided into development set (20.3 hours, 216 recordings), and test set (53.5 hours, 310 recordings). The number of speakers in each recording has a wide range of variety from 1 speaker to 21 speakers. The audio includes various types of noises such as background music, laughter etc. It also contains noticeable portion of overlapping speech from 0% to 30.1% dependent on the recording. While the dataset contains the visual information as well as audio, as of January 2021, only the audio of the development set is released under a Creative Commons Attribution 4.0 International License for research purpose. The audio of the evaluation set was used at the track 4 of the VoxCeleb Speaker Recognition Challenge 2020 (Section 5.3) as a blind test set.

5.1.7. LibriCSS

The LibriCSS corpus [87] is 10 hours of multi-channel recordings designed for the research of speech separation, speech recognition, and speaker diarization. It was made by playing back the audio in the LibriSpeech corpus [193] in a real meeting room, and recorded by a 7-ch microphone array. It consists of 10 sessions, each of which is further decomposed to six 10-min mini-sessions. Each mini-session was made by audio of 8 speakers and designed to have different overlap ratio from 0% to 40%. To facilitate the research, the baseline system for speech separation and ASR [87] and the baseline system that integrates speech separation, speaker diarization and ASR [88] has been developed and released.

5.2. Diarization Evaluation Metrics

5.2.1. DER

The accuracy of speaker diarization system is measured by Diarization Error Rate (DER) [194] where DER is sum of three

different error types: False alarm (FA) of speech, missed detection of speech and confusion between speaker labels.

$$\text{DER} = \frac{\text{FA} + \text{Missed} + \text{Speaker-Confusion}}{\text{Total Duration of Time}} \quad (37)$$

To establish a one-to-one mapping between the hypothesis outputs and the reference transcript, Hungarian algorithm [195] is employed. In Rich Transcription 2006 evaluation [194], 0.25 second of “no score” collar is set around every boundary of reference segment to mitigate the effect of inconsistent annotation and human errors in reference transcript and this evaluation scheme has been most widely used in speaker diarization studies.

5.2.2. JER

Jaccard Error Rate (JER) was first introduced in DIHARD II evaluation. The goal of JER is to evaluate each speaker with equal weight. Unlike DER, JER does not use speaker error to obtain the error value.

$$\text{JER} = \frac{1}{N} \sum_i \frac{\text{FA}_i + \text{MISS}_i}{\text{TOTAL}_i} \quad (38)$$

In Eq. (38), TOTAL is union of i -th speaker’s speaking time in reference transcript and i -th speaker’s speaking time in the hypotheses. The sum of FA and MISS divided by TOTAL value is then averaged over N_{ref} -speakers in the reference script. Since JER is using union operation between reference and the hypotheses, JER never exceeds 100% while DER can sometimes reach way over 100%. DER and JER are highly correlated but if a subset of speakers are dominant in the given audio recording, JER tends to get higher than ordinary case.

5.2.3. WDER

While DER is based on the duration of speaking time of each speaker, Word-level DER (WDER) is designed to measure the error that is caused in the lexical(output transcription) side. The motivation of WDER is the discrepancy between DER and the accuracy of final transcript output since DER relies on the duration of speaking time that is not always aligned with the word boundaries. The concept of word-breakage was proposed in Silovsky et al. [175] where WB shares the similar idea with

WDER. Unlike WDER, WB measures the number of speaker change point occurs inside a word boundary. The work in Park and Georgiou [196] suggested the term WDER, evaluating the diarization output with ground-truth transcription. More recently, the joint ASR and speaker diarization system was evaluated in WDER format in Shafey et al. [77]. Although the way of calculating WDER would differ over the studies but the underlying idea is that the diarization error is calculated by counting the correctly or incorrectly labeled words.

5.3. Diarization Evaluation Series

The Rich Transcription (RT) Evaluation [20] is the pioneering evaluation series of initiating deeper investigation on speaker diarization in relation with ASR. The main purpose of this effort was to create ASR technologies that would produce transcriptions with descriptive metadata, like who said when, where speaker diarization plays in. Thus the main tasks in the evaluation were naturally ASR and speaker diarization. The domains of the data of interest were broadcast news, CTS and meeting recordings with multiple participants. Throughout the period 2002 to 2009, the RT evaluation series promoted and gauged advances in speaker diarization as well as ASR technology.

DIHARD challenge [189, 85] is the most recent evaluation that focuses on challenging diarization tasks. DIHARD challenge data contains many different challenging and diverse domains including the recordings from restaurants, meetings, interview videos and court room. DIHARD evaluation focuses on the performance gap of state-of-the-art diarization systems on challenging domains (e.g. recordings from outdoors) and relatively clean speech (e.g. telephonic speech). DIHARD challenge employs a stricter evaluation scheme where the scoring rule does not have “no score” collar and also evaluates overlapped regions. In addition, DIHARD challenge also employed JER.

The CHiME-6 challenge [89] track 2 revisits the previous CHiME-5 challenge [191] and further considers the problem of distant multi-microphone conversational speech diarization and recognition in everyday home environments. Although the final evaluation criterion is ranked with the WER, the challenge participants in this track also need to submit the diarization result. The evaluation metrics of the diarization follow the DIHARD challenge, i.e., “no score” collar and it also evaluates overlapped regions when computing the DER and JER.

The VoxCeleb Speaker Recognition Challenge (VoxSRC) is the recent evaluation series for speaker recognition systems [197, 105]. The goal of VoxSRC is to probe how well the current technology can cope with the speech “in the wild”. The evaluation data is obtained from YouTube videos of various domains, such as celebrity interviews, news shows, talk shows, and debates. The audio includes various types of background noises, laughter as well as noticeable portion of overlapping speech, all of which make the task very challenging. This evaluation series initially started with a pure speaker verification task [197], and the diarization task was added as the track 4 at the latest evaluation at the VoxCeleb Speaker Recognition Challenge 2020 (VoxSRC-20) [105]. The VoxConverse dataset

[192] was used for evaluation with DER as a primary metric to determine the ranking of submitted systems. JER was also measured as a secondary metric.

6. Applications

6.1. Meeting Transcription

The goal of meeting transcription is to automatically generate speaker-attributed transcripts during real-life meetings based on their audio and optionally video recordings. Accurate meeting transcriptions are the one of the processing steps in a pipeline for several tasks like summarization, topic extraction, etc. Similarly, the same transcription system can be used in other domains such as healthcare [198]. Although this task was introduced by NIST in the Rich Transcription Evaluation series back in 2003 [180, 188, 199], the initial systems had very poor performance, and consequently commercialization of the technology was not possible. However, recent advances in the areas of Speech Recognition [200, 201], far-field speech processing [202, 203, 204], Speaker ID and diarization [205, 206, 113], have greatly improved the speaker-attributed transcription accuracy, enabling such commercialization. Bi-modal processing combining cameras with microphone arrays has further improved the overall performance [207, 208]. As such, these latest trends motivated us to include an end-to-end audio-visual meeting transcription system overview in this paper.

Reflecting the variety of application scenarios, customer needs, and business scope, different constraints may be imposed on meeting transcription systems. For example, it is most often required to provide the resulting transcriptions in low latency, making the diarization and recognition even more challenging. On the other hand, the architecture of the transcription system can substantially improve the overall performance, e.g., employing microphone arrays of known geometry as the input device. Also, in the case where the expected meeting attendees are known beforehand, the transcription system can further improve speaker attribution, all while providing the exact name of the speaker, instead of a randomly generated discrete speaker labels.

Two different scenarios in this space are presented: first, a fix-geometry microphone array combined with a fish-eye camera system, and second, an ad-hoc geometry microphone array system without a camera. In both scenarios, a “non-binding” list of participants and their corresponding speaker profiles are considered known. In more detail, the transcription system has access to the invitees’ names and profiles, however the actual attendees may not accurately match those invited. As such, there is an option to either include “unannounced” participants. Also, some of the invitees may not have profiles. In both scenarios, there is a constraint of low-latency transcriptions, where initial results need to be shown with low latency. The finalized results can be updated later in an offline fashion.

Some of the technical challenges to overcome are [209]:

1. Although ASR on overlapping speech is one of the main challenges in meeting transcription, limited progress has been made over the years. Numerous multi-channel

speech separation methods have been proposed based on Independent Component Analysis(ICA) or Spatial Clustering [210, 211, 212, 213, 214, 215], but applying them to a meeting setup had limited success. In addition, neural network-based separation methods like Permutation Invariant Training (PIT) [59] or deep clustering (DC) [58] cannot adequately address reverberation and background noise [216].

2. Flexible framework: It is desirable that the transcription system can process all the available information, such as the multi-channel audio and the visual cues. The system needs to process a dynamically changing number of audio channels without loss of performance. As such, the architecture needs to be modular enough to encompass the different settings.
3. Speaker-Attributed ASR of natural meetings requires on-line/streaming ASR, audio pre-processing such as dereverberation, and accurate diarization and speaker identification. These multiple processing steps are usually optimized separately and thus, the overall pipeline is most frequently inefficient.
4. Using multiple, not-synchronized audio streams, e.g., audio capturing with mobile devices, adds complexity to the meeting setup and processing. In return, we gain potentially better spatial coverage since the devices are usually distributed around the room and near the speakers. As part of the application scenario, the meeting participants bring their personal devices, which can be re-purposed to improve the overall meeting transcription quality. On the other hand, while there are several pioneering studies [217], it is unclear what the best strategies are for consolidating multiple asynchronous audio streams and to what extent they work for natural meetings in online and offline setups.

Based on these considerations, an architecture of meeting transcription system with asynchronous distant microphones have been proposed in [161]. In this work, various fusion strategies have been investigating: from early fusion beamforming the audio signals, to mid-fusion combining senones per channel, to late fusion combining the diarization and ASR results [147]. The resulting system performance was benchmarked on real-world meeting recordings against fix-geometry systems. As mentioned above, the requirement of speaker-attributed transcriptions with low latency was adhered, as well. In addition to the end-to-end system analysis, the paper [161] proposed the idea of “leave-one-out beamforming” in the asynchronous multi-microphone setup, enriching the “diversity” of the resulting signals, as proposed in [218]. Finally, it is described how an online, incremental version of ROVER can process both the ASR and diarization outputs, enhancing the overall speaker-attributed ASR performance.

6.2. Conversational Interaction Analysis and Behavioral Modeling

Speech and spoken language are central to conversational interactions and carry crucial information about a speaker’s intent, emotions, identity, age and other individual and interpersonal trait and state variables including health state, and computational advances are increasingly allowing for accessing such rich information [219, 220]. For example, knowing how much, and how, a child speaks in an interaction reveals critical information about the developmental state, and offers clues to clinicians in diagnosing disorders such as Autism [221]. Such analyses are made possible by capturing and processing the audio recordings of the interactions, often involving two or more people. An important foundational step is identifying and associating the speech portions belonging to specific individuals involved in the conversation. The technologies that provide this capability are speech activity detection (SAD) and speaker diarization. Speech portions segmented with speaker-specific information provided by speaker diarization, by itself without any explicit lexical transcription, can offer important information to domain experts who can take advantage of speaker diarization results for quantitative turn-taking analysis.

A domain that is the most relevant such analyses of spoken conversational interactions relates to behavioral signal processing (BSP) [222, 219] which refers to the technology and algorithms for modeling and understanding human communicative, affective and social behavior. For example, these may include analyzing how positive or negative a person is, how empathic an individual toward another, what does the behavior patterns reveal about the relationship status, and health condition of an individual [220]. BSP involves addressing all the complexities of spontaneous interactions in conversations with additional challenges involved in handling and understanding emotional, social and interpersonal behavioral dynamics revealed through vocal verbal and nonverbal cues of the interaction participants. Therefore, the knowledge of speaker specific vocal information plays a significant role in BSP, requiring highly accurate speaker diarization performance. For example, speaker diarization module is employed as a pre-processing module for analyzing psychotherapy mechanisms and quality [223], and suicide risk assessment [224].

Another popular application of speaker diarization for conversation interaction analysis is the medical doctor-patient interactions. In the system described in [225], the nature of memory problem of a patient is detected from the conversations between neurologists and patients. Speech and language features extracted from ASR transcripts combined with speaker diarization results are used to predict the type of disorder. An automated assistant system for medical domain transcription is proposed in [226] which includes speaker diarization module, ASR module and natural language generation (NLG) module. The automated assistant module accepts the audio clip and outputs grammatically correct sentences that describe the topic of the conversation, subject and subject’s symptom.

6.3. Audio indexing

Content-based audio indexing is a well known application domain for speaker diarization. It can provide meta information such as the content or data type of a given audio data to make

information retrieval efficient since search query by machines would be limited by such metadata. The more diverse information were available, the better efficiency we could achieve in retrieving audio contents from a database.

One useful piece of information for the audio indexing would be ASR transcripts to understand the content of speech portions in the audio data. Speaker diarization can augment those transcripts in terms of “who spoke when”, which was the main purpose of the Rich Transcription evaluation series [20] as we discussed in Sections 4.1 and 5.3. The aggregated spoken utterances from speakers by a speaker diarization system also enable per-speaker summary or keyword list-up, which can be used for another query values to retrieve relevant contents from the database. In [227], we can peek a view of how speaker diarization outputs can be linked for information searching in consumer facing applications.

6.4. Conversational AI

Thanks to the advance of ASR technology, the applications of ASR are evolved from simple voice command recognition systems to conversational AI systems. Conversational AI systems, as opposed to voice command recognition systems, have features that voice command recognition systems are lack of. The fundamental idea of conversational AI is making a machine that humans can talk to and interact with the system. In this sense, focusing on an interested speaker in multi-party setting is one of the most important feature of conversational AI and speaker diarization becomes essential feature for conversational AI. For example, conversational AI equipped in a car can pay attention to a specific speaker that is demanding a piece of information from the navigation system by applying speaker diarization along with ASR.

Smart speakers and voice assistants are the most popular products where speaker diarization plays a significant role for conversational AI. Since response time and online processing are the crucial factors in real-life settings, the demand for end-to-end speaker diarization system integrated into ASR pipeline is growing. The performance of incremental (online) ASR and speaker diarization of the commercial ASR services are evaluated and compared in [228]. It is expected that the real-time and low latency aspect of speaker diarization will be more emphasized in the speaker diarization systems in the future since the performance of online diarization and online ASR still have much room for improvement.

7. Challenges and the Future of Speaker Diarization

This paper has provided a comprehensive overview of speaker diarization techniques, highlighting the recent development of deep learning-based diarization approaches. In the early days, a speaker diarization system was developed as a pipeline of sub-modules including front-end processing, speech activity detection, segmentation, speaker embedding extraction, clustering, and post-processing, leading to a standalone system without much connection to other components in a given speech application. As the rise of the deep learning technology,

more and more advancements have been made for speaker diarization, from a method that replaces a single module into a deep-learning-based one, to a fully end-to-end neural diarization. Furthermore, as the speech recognition technology becomes more accessible, a trend to tightly integrate speaker diarization and ASR systems has emerged, such as benefiting from the ASR output to improve speaker diarization accuracy. As of late, joint modeling for speaker diarization and speech recognition is investigated in an attempt to enhance the overall performance. Thanks to these great achievement, speaker diarization systems have already been deployed in many applications, including meeting transcription, conversational interaction analysis, audio indexing, and conversational AI systems.

As we have seen, tremendous progress has been made for speaker diarization systems. Nevertheless, there are still much room for improvement. As the final remark, we conclude this paper by listing up the remaining challenges for speaker diarization towards future research and development.

Online processing of speaker diarization. Most speaker diarization methods assume that an entire recording can be observed to execute speaker diarization. However, many applications such as meeting transcription systems or smart agents require only short latency for assigning the speaker. While there have been several attempts to make online speaker diarization system both for clustering-based systems (e.g., [205]) and neural network-based diarization systems (e.g., [55, 172, 173]), it’s still remaining as a challenging problem.

Domain mismatch. A model that is trained on a data in a specific domain often works poorly on a data in another domain. For example, it is experimentally known that the EEND model tends to overfit to the distribution of the speaker overlaps of the training data [56]. Such domain mismatch issue is universal for any training-based method. Given the growing interest for trainable speaker diarization systems, it will become more important to assess the ability for handling the variety of inputs. The international evaluation efforts for speaker diarization such as the DIHARD challenge [189, 85, 190] or VoxSRC [197, 105] will also have great importance for that direction.

Speaker overlap. Overlap of multi-talker speech is inevitable nature of conversation. For example, average 12% to 15% of speaker overlap was observed for meeting recordings [229, 102], and it can become higher for daily conversations [230, 191, 89]. Nevertheless, many conventional speaker diarization systems, especially clustering-based systems, treated only non-overlapped region of recordings sometimes even for the evaluation metric. While the topic has been studied for long years (e.g. early works [231, 232]), there is a growing interest for handling the speaker overlaps towards better speaker diarization, including the application of speech separation [104], post-processing [233, 162], and joint modeling of speech separation and speaker diarization [76, 184].

Integration with ASR. Not all but many applications require ASR results along with speaker diarization results. In the line of

the modular combination of speaker diarization and ASR, some systems put a speaker diarization system before ASR [91] while some systems put a diarization system after ASR [209]. Both types of systems showed a strong performance for a specific task, and it is still an open problem that what kind of system architecture is the best for the speaker diarization and ASR tasks [88]. Furthermore, there is another line of research to jointly perform speaker diarization and ASR [77, 78, 79, 184] as introduced in Section 4. The joint modeling approach could leverage the inter-dependency between speaker diarization and ASR to better perform both tasks. However, it has not yet fully investigated whether such joint frameworks perform better than the well-tuned modular systems. Overall, the integration of speaker diarization and ASR is one of the hottest topics that has still been pursued.

Audio visual modeling. Visual information contains a strong clue to identify speakers. For example, the video captured by a fisheye camera was used to improve the speaker diarization accuracy in a meeting transcription task [209]. The visual information was also used to significantly improve the speaker diarization accuracy for speaker diarization on YouTube video [192]. While these studies showed the effectiveness of visual information, the audio-visual speaker diarization has yet been rarely investigated compared with audio-only speaker diarization, and there will be many rooms for the improvement.

References

- [1] S. E. Tranter, K. Yu, D. A. Reynolds, G. Evermann, D. Y. Kim, P. C. Woodland, An investigation into the the interactions between speaker diarisation systems and automatic speech transcription, CUED/F-INFENG/TR-464 (2003).
- [2] S. E. Tranter, D. A. Reynolds, An overview of automatic speaker diarization systems, IEEE Transactions on Audio, Speech, and Language Processing 14 (2006) 1557–1565.
- [3] X. Anguera, S. Bozonnet, N. Evans, C. Fredouille, G. Friedland, O. Vinyals, Speaker diarization: A review of recent research, IEEE Transactions on Audio, Speech, and Language Processing 20 (2012) 356–370.
- [4] H. Gish, M. . Siu, R. Rohlicek, Segregation of speakers for speech recognition and speaker identification, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 1991, pp. 873–876.
- [5] M.-H. Siu, Y. George, H. Gish, An unsupervised, sequential learning algorithm for segmentation for speech waveforms with multiple speakers, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 1992, pp. 189–192.
- [6] J. R. Rohlicek, D. Ayuso, M. Bates, R. Bobrow, A. Boulanger, H. Gish, P. Jeanrenaud, M. Meteer, M. Siu, Gisting conversational speech, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 1992, pp. 113–116.
- [7] M. Sugiyama, J. Murakami, H. Watanabe, Speech segmentation and clustering based on speaker features, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 1993, pp. 395–398.
- [8] U. Jain, M. A. Siegler, S.-J. Doh, E. Gouvea, J. Huerta, P. J. Moreno, B. Raj, R. M. Stern, Recognition of continuous broadcast news with multiple unknown speakers and environments, in: Proceedings of ARPA Spoken Language Technology Workshop, 1996, pp. 61–66.
- [9] M. Padmanabhan, L. R. Bahl, D. Nahamoo, M. A. Picheny, Speaker clustering and transformation for speaker adaptation in large-vocabulary speech recognition systems, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 1996, pp. 701–704.
- [10] M. A. Siegler, U. Jain, B. Raj, R. M. Stern, Automatic segmentation, classification and clustering of broadcast news audio, in: Proceedings of DARPA Speech Recognition Workshop, 1997, pp. 97–99.
- [11] H. Jin, F. Kubala, R. Schwartz, Automatic speaker clustering, in: Proceedings of Speech Recognition Workshop, 1997.
- [12] H. S. Beigi, S. H. Maes, Speaker, channel and environment change detection, in: Proceedings of World Congress of Automation, 1998.
- [13] S. S. Chen, P. S. Gopalakrishnan, Speaker, environment and channel change detection and clustering via the Bayesian Information Criterion, in: Tech. Rep., IBM T. J. Watson Research Center, 1998, pp. 127–132.
- [14] A. Solomonoff, A. Mielke, M. Schmidt, H. Gish, Clustering speakers by their voices, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 1998, pp. 757–760.
- [15] J.-L. Gauvain, G. Adda, L. Lamel, M. Adda-Decker, Transcription of broadcast news: The LIMSI Nov 96 Hub4 system, in: Proceedings of ARPA Speech Recognition Workshop, 1997, pp. 56–63.
- [16] J.-L. Gauvain, L. Lamel, G. Adda, The LIMSI 1997 Hub-4E transcription system, in: Proceedings of DARPA News Transcription and Understanding Workshop, 1998, pp. 75–79.
- [17] J.-L. Gauvain, L. Lamel, G. Adda, Partitioning and transcription of broadcast news data, in: Proceedings of the International Conference on Spoken Language Processing, 1998, pp. 1335–1338.
- [18] D. Liu, F. Kubala, Fast speaker change detection for broadcast news transcription and indexing, in: Proceedings of the International Conference on Spoken Language Processing, 1999, pp. 1031–1034.
- [19] AMI Consortium. <http://www.amiproject.org/index.html>.
- [20] NIST, Rich Transcription Evaluation. <https://www.nist.gov/itl/iad/mig/rich-transcription-evaluation>.
- [21] J. Ajmera, C. Wooters, A robust speaker clustering algorithm, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2003, pp. 411–416.
- [22] S. E. Tranter, D. A. Reynolds, Speaker diarisation for broadcast news, in: Proceedings of Odyssey Speaker and Language Recognition Workshop, 2004, pp. 337–344.
- [23] C. Wooters, J. Fung, B. Peskin, X. Anguera, Toward robust speaker segmentation: The ICSI-SRI Fall 2004 diarization system, in: Proceedings of Fall 2004 Rich Transcription Workshop, 2004, pp. 402–414.
- [24] D. A. Reynolds, P. Torres-Carrasquillo, The MIT Lincoln Laboratory RT-04F diarization systems: Applications to broadcast audio and telephone conversations, in: Proceedings of Fall 2004 Rich Transcription Workshop, 2004.
- [25] D. A. Reynolds, P. Torres-Carrasquillo, Approaches and applications of audio diarization, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2005, pp. 953–956.
- [26] X. Zhu, C. Barras, S. Meignier, J.-L. Gauvain, Combining speaker identification and BIC for speaker diarization, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2005, pp. 2441–2444.
- [27] C. Barras, Xuan Zhu, S. Meignier, J.-L. Gauvain, Multistage speaker diarization of broadcast news, IEEE Transactions on Audio, Speech, and Language Processing 14 (2006) 1505–1512.
- [28] N. Mirghafori, C. Wooters, Nuts and flakes: A study of data characteristics in speaker diarization, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2006, pp. 1017–1020.
- [29] S. Meignier, D. Moraru, C. Fredouille, J.-F. Bonastre, L. Besacier, Step-by-step and integrated approaches in broadcast news speaker diarization, Computer, Speech & Language 20 (2006) 303–330.
- [30] A. E. Rosenberg, A. Gorin, Z. Liu, P. Parthasarathy, Unsupervised speaker segmentation of telephone conversations, in: Proceedings of the International Conference on Spoken Language Processing, 2002, pp. 565–568.
- [31] D. Liu, F. Kubala, A cross-channel modeling approach for automatic segmentation of conversational telephone speech, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2003, pp. 333–338.
- [32] S. E. Tranter, K. Yu, G. Evermann, P. C. Woodland, Generating and evaluating for automatic speech recognition of conversational telephone speech, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2004, pp. 753–756.
- [33] D. A. Reynolds, P. Kenny, F. Castaldo, A study of new approaches

- to speaker diarization, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2009, pp. 1047–1050.
- [34] P. Kenny, D. Reynolds, F. Castaldo, Diarization of telephone conversations using factor analysis, *IEEE Journal of Selected Topics in Signal Processing* 4 (2010) 1059–1070.
- [35] T. Pfau, D. Ellis, A. Stolcke, Multispeaker speech activity detection for the ICSI meeting recorder, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2001, pp. 107–110.
- [36] J. Ajmera, G. Lathoud, L. McCowan, Clustering and segmenting speakers and their locations in meetings, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2004, pp. 605–608.
- [37] Q. Jin, K. Laskowski, T. Schultz, A. Waibel, Speaker segmentation and clustering in meetings, in: Proceedings of the International Conference on Spoken Language Processing, 2004, pp. 597–600.
- [38] X. Anguera, C. Wooters, B. Peskin, M. Aguilo, Robust speaker segmentation for meetings: The ICSI-SRI Spring 2005 diarization system, in: Proceedings of Machine Learning for Multimodal Interaction Workshop, 2005, pp. 402–414.
- [39] X. Anguera, C. Wooters, J. Hernando, Purity algorithms for speaker diarization of meetings data, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, volume 1, 2006, pp. 1025–1028.
- [40] D. Istrate, C. Fredouille, S. Meignier, L. Besacier, J.-F. Bonastre, NIST RT05S evaluation: Pre-processing techniques and speaker diarization on multiple microphone meetings, in: Proceedings of Machine Learning for Multimodal Interaction Workshop, 2006.
- [41] D. A. V. Leeuwen, M. Konecny, Progress in the AMIDA speaker diarization system for meeting data, in: Proceedings of International Evaluation Workshops CLEAR 2007 and RT 2007, 2007, pp. 475–483.
- [42] X. Anguera, C. Wooters, J. Hernando, Acoustic beamforming for speaker diarization of meetings, *IEEE Transactions on Audio, Speech, and Language Processing* 15 (2007) 2011–2023.
- [43] X. Zhu, C. Barras, L. Lamel, J.-L. Gauvain, Multi-stage speaker diarization for conference and lecture meetings, in: Proceedings of International Evaluation Workshops CLEAR 2007 and RT 2007, 2007, pp. 533–542.
- [44] D. Vijayasenan, F. Valente, H. Bourlard, An information theoretic approach to speaker diarization of meeting data, *IEEE Transactions on Audio, Speech, and Language Processing* 17 (2009) 1382–1393.
- [45] F. Valente, P. Motlicek, D. Vijayasenan, Variational Bayesian speaker diarization of meeting recordings, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2010, pp. 4954–4957.
- [46] P. Kenny, G. Boulianne, P. Ouellet, P. Dumouchel, Joint factor analysis versus eigenchannels in speaker recognition, *IEEE Transactions on Audio, Speech, and Language Processing* 15 (2007) 1435–1447.
- [47] E. Varni, X. Lei, E. McDermott, I. L. Moreno, J. G-Dominguez, Deep neural networks for small footprint text-dependent speaker verification, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2014, pp. 4052–4056.
- [48] G. Heigold, I. Moreno, S. Bengio, N. Shazeer, End-to-end text-dependent speaker verification, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2016, pp. 5115–5119.
- [49] Q. Wang, C. Downey, L. Wan, P. A. Mansfield, I. L. Moreno, Speaker diarization with LSTM, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2018, pp. 5239–5243.
- [50] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, S. Khudanpur, X-vectors: Robust DNN embeddings for speaker recognition, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2018, pp. 5329–5333.
- [51] N. Dehak, P. Kenny, R. Dehak, P. Dumouchel, P. Ouellet, Front-end factor analysis for speaker verification, *IEEE Transactions on Audio, Speech, and Language Processing* 19 (2011).
- [52] S. Shum, N. Dehak, J. Glass, On the use of spectral and iterative methods for speaker diarization, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2012, pp. 482–485.
- [53] G. Dupuy, M. Rouvier, S. Meignier, Y. Esteve, i-Vectors and ILP clustering adapted to cross-show speaker diarization, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2012, pp. 2174–2177.
- [54] S. H. Shum, N. Dehak, R. Dehak, J. R. Glass, Unsupervised methods for speaker diarization: An integrated and iterative approach, *IEEE Transactions on Audio, Speech, and Language Processing* 21 (2013).
- [55] A. Zhang, Q. Wang, Z. Zhu, J. Paisley, C. Wang, Fully supervised speaker diarization, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2019, pp. 6301–6305.
- [56] Y. Fujita, N. Kanda, S. Horiguchi, K. Nagamatsu, S. Watanabe, End-to-end neural speaker diarization with permutation-free objectives, Proceedings of the Annual Conference of the International Speech Communication Association (2019) 4300–4304.
- [57] Y. Fujita, N. Kanda, S. Horiguchi, Y. Xue, K. Nagamatsu, S. Watanabe, End-to-end neural speaker diarization with self-attention, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, IEEE, 2019, pp. 296–303.
- [58] J. R. Hershey, Z. Chen, J. Le Roux, S. Watanabe, Deep clustering: Discriminative embeddings for segmentation and separation, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2016, pp. 31–35.
- [59] M. Kolbæk, D. Yu, Z.-H. Tan, J. Jensen, Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks, *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 25 (2017) 1901–1913.
- [60] Y. Luo, N. Mesgarani, Conv-tasnet: Surpassing ideal time-frequency magnitude masking for speech separation, *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 27 (2019) 1256–1266.
- [61] E. Varni, X. Lei, E. McDermott, I. L. Moreno, J. Gonzalez-Dominguez, Deep neural networks for small footprint text-dependent speaker verification, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2014, pp. 4052–4056.
- [62] D. Snyder, D. Garcia-Romero, D. Povey, S. Khudanpur, Deep neural network embeddings for text-independent speaker verification, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2017, pp. 999–1003.
- [63] T. Drugman, Y. Stylianou, Y. Kida, M. Akamine, Voice activity detection: Merging source and filter-based information, *IEEE Signal Processing Letters* 23 (2015) 252–256.
- [64] X. Guo, L. Gao, X. Liu, J. Yin, Improved deep embedded clustering with local structure preservation, in: Proceedings of International Joint Conference on Artificial Intelligence, 2017, pp. 1753–1759.
- [65] J. Wang, X. Xiao, J. Wu, R. Ramamurthy, F. Rudzicz, M. Brudno, Speaker diarization with session-level speaker embedding refinement using graph neural networks, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2020, pp. 7109–7113.
- [66] I. Medennikov, M. Korenevsky, T. Prisyach, Y. Khokhlov, M. Korenevskaya, I. Sorokin, T. Timofeeva, A. Mitrofanov, A. Andrusenko, I. Podluzhny, A. Laptev, A. Romanenko, Target-speaker voice activity detection: a novel approach for multi-speaker diarization in a dinner party scenario, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2020, pp. 274–278.
- [67] D. Yu, X. Chang, Y. Qian, Recognizing multi-talker speech with permutation invariant training, Proceedings of the Annual Conference of the International Speech Communication Association (2017) 2456–2460.
- [68] H. Seki, T. Hori, S. Watanabe, J. Le Roux, J. R. Hershey, A purely end-to-end system for multi-speaker speech recognition, 2018, pp. 2620–2630.
- [69] X. Chang, Y. Qian, K. Yu, S. Watanabe, End-to-end monaural multi-speaker ASR system without pretraining, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2019, pp. 6256–6260.
- [70] N. Kanda, Y. Fujita, S. Horiguchi, R. Ikeshita, K. Nagamatsu, S. Watanabe, Acoustic modeling for distant multi-talker speech recognition with single- and multi-channel branches, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2019, pp. 6630–6634.
- [71] N. Kanda, S. Horiguchi, R. Takashima, Y. Fujita, K. Nagamatsu, S. Watanabe, Auxiliary interference speaker loss for target-speaker speech recognition, in: Proceedings of the Annual Conference of the

- International Speech Communication Association, 2019, pp. 236–240.
- [72] X. Wang, N. Kanda, Y. Gaur, Z. Chen, Z. Meng, T. Yoshioka, Exploring end-to-end multi-channel asr with bias information for meeting transcription, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2021.
- [73] P. Wang, Z. Chen, X. Xiao, Z. Meng, T. Yoshioka, T. Zhou, L. Lu, J. Li, Speech separation using speaker inventory, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2019, pp. 230–236.
- [74] C. Han, Y. Luo, C. Li, T. Zhou, K. Kinoshita, S. Watanabe, M. Delcroix, H. Erdogan, J. R. Hershey, N. Mesgarani, et al., Continuous speech separation using speaker inventory for long multi-talker recording, arXiv preprint arXiv:2012.09727 (2020).
- [75] Z. Huang, S. Watanabe, Y. Fujita, P. García, Y. Shao, D. Povey, S. Khudanpur, Speaker diarization with region proposal network, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2020, pp. 6514–6518.
- [76] T. von Neumann, K. Kinoshita, M. Delcroix, S. Araki, T. Nakatani, R. Haeb-Umbach, All-neural online source separation, counting, and diarization for meeting analysis, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2019, pp. 91–95.
- [77] L. E. Shafey, H. Soltau, I. Shafran, Joint Speech Recognition and Speaker Diarization via Sequence Transduction, in: Proceedings of the Annual Conference of the International Speech Communication Association, ISCA, 2019, pp. 396–400.
- [78] H. H. Mao, S. Li, J. McAuley, G. Cottrell, Speech recognition and multi-speaker diarization of long conversations, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2020, pp. 691–695.
- [79] N. Kanda, S. Horiguchi, Y. Fujita, Y. Xue, K. Nagamatsu, S. Watanabe, Simultaneous speech recognition and speaker diarization for monaural dialogue recordings with target-speaker acoustic models, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2019, pp. 31–38.
- [80] N. Kanda, X. Chang, Y. Gaur, X. Wang, Z. Meng, Z. Chen, T. Yoshioka, Investigation of end-to-end speaker-attributed ASR for continuous multi-talker recordings, in: Proceedings of IEEE Spoken Language Technology Workshop, 2021.
- [81] R. Haeb-Umbach, S. Watanabe, T. Nakatani, M. Bacchiani, B. Hoffmeister, M. L. Seltzer, H. Zen, M. Souden, Speech processing for digital home assistants: Combining signal processing with deep-learning techniques, IEEE Signal Processing Magazine 36 (2019) 111–124.
- [82] E. Vincent, T. Virtanen, S. Gannot, Audio source separation and speech enhancement, John Wiley & Sons, 2018.
- [83] D. Wang, J. Chen, Supervised speech separation based on deep learning: An overview, IEEE/ACM Transactions on Audio, Speech, and Language Processing 26 (2018) 1702–1726.
- [84] G. Sell, D. Snyder, A. McCree, D. Garcia-Romero, J. Villalba, M. Maciejewski, V. Manohar, N. Dehak, D. Povey, S. Watanabe, et al., Diarization is hard: Some experiences and lessons learned for the JHU team in the inaugural DIHARD challenge., in: Proceedings of the Annual Conference of the International Speech Communication Association, 2018, pp. 2808–2812.
- [85] N. Ryant, K. Church, C. Cieri, A. Cristia, J. Du, S. Ganapathy, M. Liberman, The second DIHARD diarization challenge: Dataset, task, and baselines, Proceedings of the Annual Conference of the International Speech Communication Association (2019) 978–982.
- [86] M. Diez, F. Landini, L. Burget, J. Rohdin, A. Silnova, K. Zmolíková, O. Novotný, K. Veselý, O. Glembek, O. Plchot, et al., BUT system for DIHARD speech diarization challenge 2018., in: Proceedings of the Annual Conference of the International Speech Communication Association, 2018, pp. 2798–2802.
- [87] Z. Chen, T. Yoshioka, L. Lu, T. Zhou, Z. Meng, Y. Luo, J. Wu, X. Xiao, J. Li, Continuous speech separation: Dataset and analysis, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2020, pp. 7284–7288.
- [88] D. Raj, P. Denisov, Z. Chen, H. Erdogan, Z. Huang, M. He, S. Watanabe, J. Du, T. Yoshioka, Y. Luo, N. Kanda, J. Li, S. Wisdom, J. R. Hershey, Integration of speech separation, diarization, and recognition for multi-speaker meetings: System description, comparison, and analysis, in: Proceedings of IEEE Spoken Language Technology Workshop, 2021.
- [89] S. Watanabe, M. Mandel, J. Barker, E. Vincent, A. Arora, X. Chang, S. Khudanpur, V. Manohar, D. Povey, D. Raj, et al., CHiME-6 challenge: Tackling multispeaker speech recognition for unsegmented recordings, in: 6th International Workshop on Speech Processing in Everyday Environments (CHiME 2020), 2020.
- [90] A. Arora, D. Raj, A. S. Subramanian, K. Li, B. Ben-Yair, M. Maciejewski, P. Zelasko, P. Garcia, S. Watanabe, S. Khudanpur, The JHU multi-microphone multi-speaker asr system for the CHiME-6 challenge, arXiv preprint arXiv:2006.07898 (2020).
- [91] I. Medennikov, M. Korenevsky, T. Prisyach, Y. Khokhlov, M. Korenevskaya, I. Sorokin, T. Timofeeva, A. Mitrofanov, A. Andrusenko, I. Podluzhny, et al., The STC system for the CHiME-6 challenge, in: CHiME 2020 Workshop on Speech Processing in Everyday Environments, 2020.
- [92] X. Lu, Y. Tsao, S. Matsuda, C. Hori, Speech enhancement based on deep denoising autoencoder., in: Proceedings of the Annual Conference of the International Speech Communication Association, 2013, pp. 436–440.
- [93] Y. Xu, J. Du, L.-R. Dai, C.-H. Lee, A regression approach to speech enhancement based on deep neural networks, IEEE/ACM Transactions on Audio, Speech, and Language Processing 23 (2014) 7–19.
- [94] H. Erdogan, J. R. Hershey, S. Watanabe, J. Le Roux, Phase-sensitive and recognition-boosted speech separation using deep recurrent neural networks, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2015, pp. 708–712.
- [95] P. C. Loizou, Speech enhancement: theory and practice, CRC press, 2013.
- [96] T. Gao, J. Du, L.-R. Dai, C.-H. Lee, Densely connected progressive learning for lstm-based speech enhancement, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2018, pp. 5054–5058.
- [97] J. Heymann, L. Drude, R. Haeb-Umbach, Neural network based spectral mask estimation for acoustic beamforming, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2016, pp. 196–200.
- [98] H. Erdogan, J. R. Hershey, S. Watanabe, M. I. Mandel, J. Le Roux, Improved MVDR beamforming using single-channel mask prediction networks, Proceedings of the Annual Conference of the International Speech Communication Association (2016) 1981–1985.
- [99] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, B.-H. Juang, Speech dereverberation based on variance-normalized delayed linear prediction, IEEE Transactions on Audio, Speech, and Language Processing 18 (2010) 1717–1731.
- [100] T. Yoshioka, T. Nakatani, Generalization of multi-channel linear prediction methods for blind mimo impulse response shortening, IEEE Transactions on Audio, Speech, and Language Processing 20 (2012) 2707–2720.
- [101] L. Drude, J. Heymann, C. Boeddeker, R. Haeb-Umbach, NARA-WPE: A python package for weighted prediction error dereverberation in numpy and tensorflow for online and offline processing, in: Speech Communication; 13th ITG-Symposium, VDE, 2018, pp. 1–5.
- [102] T. Yoshioka, H. Erdogan, Z. Chen, X. Xiao, F. Alleva, Recognizing overlapped speech in meetings: A multichannel separation approach using neural networks, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2018, pp. 3038–3042.
- [103] C. Boeddecker, J. Heitkaemper, J. Schmalenstroeer, L. Drude, J. Heymann, R. Haeb-Umbach, Front-end processing for the CHiME-5 dinner party scenario, in: Proceedings of CHiME 2018 Workshop on Speech Processing in Everyday Environments, 2018, pp. 35–40.
- [104] X. Xiao, N. Kanda, Z. Chen, T. Zhou, T. Yoshioka, Y. Zhao, G. Liu, J. Wu, J. Li, Y. Gong, Microsoft speaker diarization system for the voxceleb speaker recognition challenge 2020, arXiv preprint arXiv:2010.11458 (2020).
- [105] A. Nagrani, J. S. Chung, J. Huh, A. Brown, E. Coto, W. Xie, M. McLaren, D. A. Reynolds, A. Zisserman, VoxSRC 2020: The second VoxCeleb speaker recognition challenge, arXiv preprint arXiv:2012.06867 (2020).
- [106] T. Ng, B. Zhang, L. Nguyen, S. Matsoukas, X. Zhou, N. Mesgarani, K. Veselý, P. Matějka, Developing a speech activity detection system for the darpa rats program, in: Proceedings of the Annual Conference of

- the International Speech Communication Association, 2012, pp. 1969–1972.
- [107] R. Sarikaya, J. H. Hansen, Robust detection of speech activity in the presence of noise, in: *Proceedings of the International Conference on Spoken Language Processing*, volume 4, Citeseer, 1998, pp. 1455–8.
 - [108] D. Haws, D. Dimitriadis, G. Saon, S. Thomas, M. Picheny, On the importance of event detection for asr, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, 2016.
 - [109] S. Meignier, D. Moraru, C. Fredouille, J.-F. Bonastre, L. Besacier, Step-by-step and integrated approaches in broadcast news speaker diarization, *Computer Speech and Language* 20 (2006) 303–330.
 - [110] S. Chen, P. Gopalakrishnan, et al., Speaker, environment and channel change detection and clustering via the bayesian information criterion, in: *Proceedings DARPA broadcast news transcription and understanding workshop*, volume 8, Virginia, USA, 1998, pp. 127–132.
 - [111] P. Delacourt, C. J. Wellekens, Distbic: A speaker-based segmentation for audio data indexing, *Speech Communication* 32 (2000) 111–126.
 - [112] M. Senoussaoui, P. Kenny, T. Stafylakis, P. Dumouchel, A study of the cosine distance-based mean shift for telephone speech diarization, *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 22 (2013) 217–227.
 - [113] G. Sell, D. Snyder, A. McCree, D. Garcia-Romero, J. Villalba, M. Maciejewski, V. Manohar, N. Dehak, D. Povey, S. Watanabe, S. Khudanpur, Diarization is hard: some experiences and lessons learned for the JHU team in the inaugural DIHARD challenge, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2018, pp. 2808–2812.
 - [114] W.-H. Tsai, S.-S. Cheng, H.-M. Wang, Speaker clustering of speech utterances using a voice characteristic reference space, in: *Proceedings of the International Conference on Spoken Language Processing*, 2004.
 - [115] J. E. Rougui, M. Rziza, D. Aboutajdine, M. Gelgon, J. Martinez, Fast incremental clustering of gaussian mixture speaker models for scaling up retrieval in on-line broadcast, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, volume 5, IEEE, 2006, pp. V–V.
 - [116] D. A. Reynolds, T. F. Quatieri, R. B. Dunn, Speaker verification using adapted gaussian mixture models, *Digital signal processing* 10 (2000) 19–41.
 - [117] P. Kenny, G. Boulianne, P. Ouellet, P. Dumouchel, Speaker and session variability in gmm-based speaker verification, *IEEE Transactions on Audio, Speech, and Language Processing* 15 (2007) 1448–1460.
 - [118] P. Kenny, P. Ouellet, N. Dehak, V. Gupta, P. Dumouchel, A study of interspeaker variability in speaker verification, *IEEE Transactions on Audio, Speech, and Language Processing* 16 (2008) 980–988.
 - [119] P. Kenny, G. Boulianne, P. Dumouchel, Eigenvoice modeling with sparse training data, *IEEE Transactions on Speech and Audio Processing* 13 (2005) 345–354.
 - [120] G. Sell, D. Garcia-Romero, Speaker diarization with plda i-vector scoring and unsupervised calibration, in: *Proceedings of IEEE Spoken Language Technology Workshop*, IEEE, 2014, pp. 413–417.
 - [121] W. Zhu, J. Pelecanos, Online speaker diarization using adapted i-vector transforms, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2016, pp. 5045–5049.
 - [122] Y. Sun, X. Wang, X. Tang, Deep learning face representation from predicting 10,000 classes, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1891–1898.
 - [123] Y. Taigman, M. Yang, M. Ranzato, L. Wolf, Deepface: Closing the gap to human-level performance in face verification, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1701–1708.
 - [124] J. Villalba, N. Chen, D. Snyder, D. Garcia-Romero, A. McCree, G. Sell, J. Borgstrom, F. Richardson, S. Shon, F. Grondin, et al., State-of-the-art speaker recognition for telephone and video speech: The JHU-MIT submission for NIST SRE18., in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 1488–1492.
 - [125] D. Comaniciu, P. Meer, Mean shift: A robust approach toward feature space analysis, *IEEE Transactions on pattern analysis and machine intelligence* 24 (2002) 603–619.
 - [126] T. Stafylakis, V. Katsouros, G. Carayannis, Speaker clustering via the mean shift algorithm, *Recall* 2 (2010) 7.
 - [127] M. Senoussaoui, P. Kenny, P. Dumouchel, T. Stafylakis, Efficient iterative mean shift based cosine dissimilarity for multi-recording speaker clustering, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2013, pp. 7712–7715.
 - [128] I. Salmun, I. Shapiro, I. Opher, I. Lapidot, Plda-based mean shift speakers’ short segments clustering, *Computer Speech and Language* 45 (2017) 411–436.
 - [129] K. J. Han, S. S. Narayanan, A robust stopping criterion for agglomerative hierarchical clustering in a speaker diarization system, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2007.
 - [130] S. Novoselov, A. Gusev, A. Ivanov, T. Pekhovsky, A. Shulipa, A. Avdeeva, A. Gorlanov, A. Kozlov, Speaker diarization with deep speaker embeddings for dihard challenge ii., in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 1003–1007.
 - [131] U. Von Luxburg, A tutorial on spectral clustering, *Statist. and Comput.* 17 (2007) 395–416.
 - [132] A. Ng, M. Jordan, Y. Weiss, On spectral clustering: Analysis and an algorithm, *Advances in neural information processing systems* 14 (2001) 849–856.
 - [133] H. Ning, M. Liu, H. Tang, T. S. Huang, A spectral clustering approach to speaker diarization, in: *Proceedings of the International Conference on Spoken Language Processing*, 2006, pp. 2178–2181.
 - [134] J. Luque, J. Hernando, On the use of agglomerative and spectral clustering in speaker diarization of meetings, in: *Proceedings of Odyssey: The Speaker and Language Recognition Workshop*, 2012, pp. 130–137.
 - [135] Q. Lin, R. Yin, M. Li, H. Bredin, C. Barras, LSTM based similarity measurement with spectral clustering for speaker diarization, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 366–370.
 - [136] T. J. Park, K. J. Han, M. Kumar, S. Narayanan, Auto-tuning spectral clustering for speaker diarization using normalized maximum eigengap, *IEEE Signal Processing Letters* 27 (2019) 381–385.
 - [137] P. Kenny, D. Reynolds, F. Castaldo, Diarization of telephone conversations using factor analysis, *IEEE Journal of Selected Topics in Signal Processing* 4 (2010) 1059–1070.
 - [138] M. Diez, L. Burget, P. Matejka, Speaker diarization based on bayesian hmm with eigenvoice priors., in: *Proceedings of Odyssey: The Speaker and Language Recognition Workshop*, 2018, pp. 147–154.
 - [139] M. Diez, L. Burget, F. Landini, J. Černocký, Analysis of speaker diarization based on bayesian hmm with eigenvoice priors, *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 28 (2019) 355–368.
 - [140] M. Diez, L. Burget, S. Wang, J. Rohdin, J. Černocký, Bayesian hmm based x-vector clustering for speaker diarization., in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 346–350.
 - [141] F. Landini, J. Profant, M. Diez, L. Burget, Bayesian hmm clustering of x-vector sequences (vbx) in speaker diarization: theory, implementation and analysis on standard tasks, *arXiv preprint arXiv:2006.07898* (2020).
 - [142] G. Sell, D. Garcia-Romero, Diarization resegmentation in the factor analysis subspace, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2015, pp. 4794–4798.
 - [143] J. G. Fiscus, A post-processing system to yield reduced word error rates: Recognizer output voting error reduction (ROVER), in: *Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding*, IEEE, 1997, pp. 347–354.
 - [144] N. Brummer, L. Burget, J. Černocký, O. Glembek, F. Grezl, M. Karafiat, D. A. van Leeuwen, P. Matejka, P. Schwarz, A. Strasheim, Fusion of heterogeneous speaker recognition systems in the STBU submission for the NIST speaker recognition evaluation 2006, *IEEE Transactions on Audio, Speech, and Language Processing* 15 (2007) 2072–2084.
 - [145] M. Huijbregts, D. van Leeuwen, F. Jong, The majority wins: a method for combining speaker diarization systems, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, ISCA, 2009, pp. 924–927.
 - [146] S. Bozonnet, N. Evans, X. Anguera, O. Vinyals, G. Friedland, C. Fredouille, System output combination for improved speaker diarization, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, ISCA, 2010, pp. 2642–2645.
 - [147] A. Stolcke, T. Yoshioka, DOVER: A method for combining diariza-

- tion outputs, in: *Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding*, IEEE, 2019, pp. 757–763.
- [148] D. Raj, L. P. Garcia-Perera, Z. Huang, S. Watanabe, D. Povey, A. Stolcke, S. Khudanpur, DOVER-Lap: A method for combining overlap-aware diarization outputs, in: *Proceedings of IEEE Spoken Language Technology Workshop*, 2021.
- [149] D. Dimitriadis, Enhancements for Audio-only Diarization Systems, arXiv preprint arXiv:1909.00082 (2019).
- [150] J. Xie, R. Girshick, A. Farhadi, Unsupervised deep embedding for clustering analysis, in: *Proceedings of International Conference on Machine Learning*, 2016, pp. 478–487.
- [151] E. Ustinova, V. Lempitsky, Learning deep embeddings with histogram loss, *Proceedings of Advances in Neural Information Processing Systems* 29 (2016) 4170–4178.
- [152] Q. Lin, Y. Hou, M. Li, Self-attentive similarity measurement strategies in speaker diarization, *Proceedings of the Annual Conference of the International Speech Communication Association* (2020) 284–288.
- [153] T. J. Park, M. Kumar, S. Narayanan, Multi-scale speaker diarization with neural affinity score fusion, arXiv preprint arXiv:2011.10527 (2020).
- [154] Y. LeCun, Y. Bengio, G. Hinton, Deep Learning, *Nature* 521 (2015) 436.
- [155] A. Santoro, R. Faulkner, D. Raposo, J. Rae, M. Chrzanowski, T. Weber, D. Wierstra, O. Vinyals, R. Pascanu, T. Lillicrap, Relational Recurrent Neural Networks, in: *Proceedings of Advances in Neural Information Processing Systems*, 2018, pp. 7299–7310.
- [156] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, T. Lillicrap, Meta-learning with Memory-Augmented Neural Networks, in: *Proceedings of International Conference on Machine Learning*, 2016, pp. 1842–1850.
- [157] S. Sukhbaatar, J. Weston, R. Fergus, et al., End-to-End Memory Networks, in: *Proceedings of Advances in Neural Information Processing Systems*, 2015, pp. 2440–2448.
- [158] D. Garcia-Romero, C. Y. Espy-Wilson, Analysis of i-vector Length Normalization in Speaker Recognition Systems, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2011, pp. 249–252.
- [159] N. Flemotomos, D. Dimitriadis, A Memory Augmented Architecture for Continuous Speaker Identification in Meetings, arXiv preprint arXiv:2001.05118 (2020).
- [160] Z. Zajíc, M. Kunešová, V. Radová, Investigation of Segmentation in i-vector Based Speaker Diarization of Telephone Speech, in: *International Conference on Speech and Computer*, 2016, pp. 411–418.
- [161] T. Yoshioka, D. Dimitriadis, A. Stolcke, W. Hinthorn, Z. Chen, M. Zeng, H. Xuedong, Meeting Transcription Using Asynchronous Distant Microphones, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 2968–2972.
- [162] S. Horiguchi, P. Garcia, Y. Fujita, S. Watanabe, K. Nagamatsu, End-to-end speaker diarization as post-processing, arXiv preprint arXiv:2012.10055 (2020).
- [163] D. M. Blei, P. I. Frazier, Distance dependent chinese restaurant processes., *Journal of Machine Learning Research* 12 (2011).
- [164] S. Ren, K. He, R. Girshick, J. Sun, Faster R-CNN: Towards real-time object detection with region proposal networks, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39 (2016) 1137–1149.
- [165] D. Kounades-Bastian, L. Girin, X. Alameda-Pineda, S. Gannot, R. Horaud, An EM algorithm for joint source separation and diarisation of multichannel convolutive speech mixtures, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2017, pp. 16–20.
- [166] D. Kounades-Bastian, L. Girin, X. Alameda-Pineda, R. Horaud, S. Gannot, Exploiting the intermittency of speech for joint separation and diarization, in: *Proceedings of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, IEEE, 2017, pp. 41–45.
- [167] K. Kinoshita, M. Delcroix, S. Araki, T. Nakatani, Tackling real noisy reverberant meetings with all-neural source separation, counting, and diarization system, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2020, pp. 381–385.
- [168] K. Maekawa, Corpus of spontaneous japanese: Its design and evaluation, in: *ISCA & IEEE Workshop on Spontaneous Speech Processing and Recognition*, 2003, pp. 7–12.
- [169] S. Horiguchi, Y. Fujita, S. Watanabe, Y. Xue, K. Nagamatsu, End-to-end speaker diarization for an unknown number of speakers with encoder-decoder based attractors, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2020, pp. 269–273.
- [170] Y. Fujita, S. Watanabe, S. Horiguchi, Y. Xue, J. Shi, K. Nagamatsu, Neural speaker diarization with speaker-wise chain rule, arXiv preprint arXiv:2006.01796 (2020).
- [171] K. Kinoshita, M. Delcroix, N. Tawara, Integrating end-to-end neural and clustering-based diarization: Getting the best of both worlds, arXiv preprint arXiv:2010.13366 (2020).
- [172] Y. Xue, S. Horiguchi, Y. Fujita, S. Watanabe, K. Nagamatsu, Online end-to-end neural diarization with speaker-tracing buffer, arXiv preprint arXiv:2006.02616 (2020).
- [173] E. Han, C. Lee, A. Stolcke, BW-EDA-EEND: Streaming end-to-end neural speaker diarization for a variable number of speakers, arXiv preprint arXiv:2011.02678 (2020).
- [174] J. Huang, E. Marcheret, K. Visweswariah, G. Potamianos, The ibm rt07 evaluation systems for speaker diarization on lecture meetings, in: *Multimodal Technologies for Perception of Humans*, Springer, 2007, pp. 497–508.
- [175] J. Silovsky, J. Zdansky, J. Nouza, P. Cerva, J. Prazak, Incorporation of the asr output in speaker segmentation and clustering within the task of speaker diarization of broadcast streams, in: *International Workshop on Multimedia Signal Processing*, IEEE, 2012, pp. 118–123.
- [176] L. Canseco-Rodriguez, L. Lamel, J.-L. Gauvain, Speaker diarization from speech transcripts, in: *Proceedings of the International Conference on Spoken Language Processing*, volume 4, 2004, pp. 3–7.
- [177] N. Flemotomos, P. Georgiou, S. Narayanan, Linguistically aided speaker diarization using speaker role information, arXiv (2019) arXiv:1911.
- [178] T. J. Park, P. Georgiou, Multimodal speaker segmentation and diarization using lexical and acoustic cues via sequence to sequence neural networks, *Proceedings of the Annual Conference of the International Speech Communication Association* (2018) 1373–1377.
- [179] T. J. Park, K. J. Han, J. Huang, X. He, B. Zhou, P. Georgiou, S. Narayanan, Speaker diarization with lexical information, *Proceedings of the Annual Conference of the International Speech Communication Association* (2019) 391–395.
- [180] J. Fiscus, J. Ajot, J. Garofolo, The Rich Transcription 2007 meeting recognition evaluation, 2007, pp. 373–389.
- [181] K. Zmolikova, M. Delcroix, K. Kinoshita, T. Higuchi, A. Ogawa, T. Nakatani, Speaker-aware neural network based beamformer for speaker extraction in speech mixtures., in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2017, pp. 2655–2659.
- [182] M. Delcroix, K. Zmolikova, K. Kinoshita, A. Ogawa, T. Nakatani, Single channel target speaker extraction and recognition with speaker beam, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2018, pp. 5554–5558.
- [183] M. Delcroix, S. Watanabe, T. Ochiai, K. Kinoshita, S. Karita, A. Ogawa, T. Nakatani, End-to-end SpeakerBeam for single channel target speech recognition., in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2019, pp. 451–455.
- [184] N. Kanda, Y. Gaur, X. Wang, Z. Meng, Z. Chen, T. Zhou, T. Yoshioka, Joint speaker counting, speech recognition, and speaker identification for overlapped speech of any number of speakers, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2020, pp. 36–40.
- [185] N. Kanda, Z. Meng, L. Lu, Y. Gaur, X. Wang, Z. Chen, T. Yoshioka, Minimum bayes risk training for end-to-end speaker-attributed asr, arXiv preprint arXiv:2011.02921 (2020).
- [186] N. Kanda, Y. Gaur, X. Wang, Z. Meng, T. Yoshioka, Serialized output training for end-to-end overlapped speech recognition, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2020, pp. 2797–2801.
- [187] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillelot, T. Hain, J. Kadlec, V. Karaikos, W. Kraaij, M. Kronenthal, et al., The ami meeting corpus: A pre-announcement, in: *International workshop on machine learning for multimodal interaction*, Springer, 2005, pp. 28–39.
- [188] A. Janin, D. Baron, J. Edwards, D. Ellis, D. Gelbart, N. Morgan, B. Piskin, T. Pfau, E. Shriberg, A. Stolcke, C. Wooters, The ICSI meeting corpus, in: *Proceedings of IEEE International Conference on Acoustics,*

- Speech and Signal Processing, 2003, pp. I–364–I–367.
- [189] N. Ryant, K. Church, C. Cieri, A. Cristia, J. Du, S. Ganapathy, M. Liberman, The first dihard speech diarization challenge, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2018.
 - [190] N. Ryant, K. Church, C. Cieri, J. Du, S. Ganapathy, M. Liberman, Third dihard challenge evaluation plan, arXiv preprint arXiv:2006.05815 (2020).
 - [191] J. Barker, S. Watanabe, E. Vincent, J. Trmal, The fifth ‘chime’ speech separation and recognition challenge: Dataset, task and baselines, Proceedings of the Annual Conference of the International Speech Communication Association (2018) 1561–1565.
 - [192] J. S. Chung, J. Huh, A. Nagrani, T. Afouras, A. Zisserman, Spot the conversation: Speaker diarisation in the wild, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2020, pp. 299–303.
 - [193] V. Panayotov, G. Chen, D. Povey, S. Khudanpur, LibriSpeech: an ASR corpus based on public domain audio books, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2015, pp. 5206–5210.
 - [194] J. G. Fiscus, J. Ajot, M. Michel, J. S. Garofolo, The rich transcription 2006 spring meeting recognition evaluation, in: Proceedings of International Workshop on Machine Learning and Multimodal Interaction, May 2006, pp. 309–322.
 - [195] P. E. Black, Hungarian algorithm, 2019. <https://xlinux.nist.gov/dads/HTML/HungarianAlgorithm.html>.
 - [196] T. J. Park, P. Georgiou, Multimodal speaker segmentation and diarization using lexical and acoustic cues via sequence to sequence neural networks, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2018, pp. 1373–1377. URL: <http://dx.doi.org/10.21437/Interspeech.2018-1364>. doi:10.21437/Interspeech.2018-1364.
 - [197] J. S. Chung, A. Nagrani, E. Coto, W. Xie, M. McLaren, D. A. Reynolds, A. Zisserman, VoxSRC 2019: The first VoxCeleb speaker recognition challenge, arXiv preprint arXiv:1912.02522 (2019).
 - [198] C. Chiu, A. Tripathi, K. Chou, C. Co, N. Jaitly, D. Jaunzeikare, A. Kannan, P. Nguyen, H. Sak, A. Sankar, J. Tansuwan, N. Wan, Y. Wu, X. Zhang, Speech recognition for medical conversations, CoRR abs/1711.07274 (2017). URL: <http://arxiv.org/abs/1711.07274>. arXiv:1711.07274.
 - [199] J. Carletta, S. Ashby, S. Bourban, M. Flynn, M. Guillemot, T. Hain, J. Kadlec, V. Karaiskos, W. Kraaij, M. Kronenthal, G. Lathoud, M. Lincoln, A. Lisowska, I. McCowan, W. P. and D. Reidsma, P. Wellner, The AMI meeting corpus: a pre-announcement, in: Proceedings of Int. Worksh. Machine Learning for Multimodal Interaction, 2006, pp. 28–39.
 - [200] W. Xiong, J. Droppo, X. Huang, F. Seide, M. Seltzer, A. Stolcke, D. Yu, G. Zweig, Achieving human parity in conversational speech recognition, CoRR abs/1610.05256 (2016). URL: <http://arxiv.org/abs/1610.05256>. arXiv:1610.05256.
 - [201] G. Saon, G. Kurata, T. Sercu, K. Audhkhasi, S. Thomas, D. Dimitriadis, X. Cui, B. Ramabhadran, M. Picheny, L. Lim, B. Roomi, P. Hall, English conversational telephone speech recognition by humans and machines, CoRR abs/1703.02136 (2017). URL: <http://arxiv.org/abs/1703.02136>. arXiv:1703.02136.
 - [202] T. Yoshioka, N. Ito, M. Delcroix, A. Ogawa, K. Kinoshita, M. Fujimoto, C. Yu, W. Fabian, M. Espi, T. Higuchi, S. Araki, T. Nakatani, The NTT CHiME-3 system: advances in speech enhancement and recognition for mobile multi-microphone devices, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2015, pp. 436–443.
 - [203] J. Du, Y. Tu, L. Sun, F. Ma, H. Wang, J. Pan, C. Liu, J. Chen, C. Lee, The USTC-iFlytek system for CHiME-4 challenge, in: Proceedings of CHiME-4 Workshop, 2016, pp. 36–38.
 - [204] B. Li, T. N. Sainath, A. Narayanan, J. Caroselli, M. Bacchiani, A. Misra, I. Shafran, H. Sak, G. Punduk, K. Chin, K. C. Sim, R. J. Weiss, K. W. Wilson, E. Variiani, C. Kim, O. Siohan, M. Weintrauba, E. McDermott, R. Rose, M. Shannon, Acoustic modeling for Google Home, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2017, pp. 399–403.
 - [205] D. Dimitriadis, P. Fousek, Developing on-line speaker diarization system, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2017, pp. 2739–2743.
 - [206] A. Zhang, Q. Wan, Z. Zhu, J. Paisley, C. Wang, Fully supervised speaker diarization, arXiv preprint arXiv:1810.04719 (2018).
 - [207] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: IEEE Conf. Computer Vision, Pattern Recognition, 2016, pp. 770–778. doi:10.1109/CVPR.2016.90.
 - [208] K. He, G. Gkioxari, P. Dollár, R. Girshick, Mask R-CNN, CoRR abs/1703.06870 (2017). URL: <http://arxiv.org/abs/1703.06870>. arXiv:1703.06870.
 - [209] T. Yoshioka, I. Abramovski, C. Aksoylar, Z. Chen, M. David, D. Dimitriadis, Y. Gong, I. Gurvich, X. Huang, Y. Huang, A. Hurvitz, L. Jiang, S. Koubi, E. Krupka, I. Leichter, C. Liu, P. Parthasarathy, A. Vinnikov, L. Wu, X. Xiao, W. Xiong, H. Wang, Z. Wang, J. Zhang, Y. Zhao, T. Zhou, Advances in Online Audio-Visual Meeting Transcription, in: Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding, 2019, pp. 276–283.
 - [210] H. Buchner, R. Aichner, W. Kellermann, A generalization of blind source separation algorithms for convolutive mixtures based on second-order statistics, IEEE Transactions on Speech and Audio Processing 13 (2005) 120–134.
 - [211] H. Sawada, S. Araki, S. Makino, Measuring dependence of bin-wise separated signals for permutation alignment in frequency-domain BSS, in: Int. Symp. Circ., Syst., 2007, pp. 3247–3250.
 - [212] F. Nesta, P. Svaizer, M. Omologo, Convolutional bss of short mixtures by ica recursively regularized across frequencies, IEEE Transactions on Audio, Speech, and Language Processing 19 (2011) 624–639.
 - [213] H. Sawada, S. Araki, S. Makino, Underdetermined convolutive blind source separation via frequency bin-wise clustering and permutation alignment, IEEE Transactions on Audio, Speech, and Language Processing 19 (2011) 516–527.
 - [214] N. Ito, S. Araki, T. Yoshioka, T. Nakatani, Relaxed disjointness based clustering for joint blind source separation and dereverberation, in: Proceedings of International Workshop on Acoustic Echo and Noise Control, 2014, pp. 268–272.
 - [215] L. Drude, R. Haeb-Umbach, Tight integration of spatial and spectral features for BSS with deep clustering embeddings, in: Proceedings of the Annual Conference of the International Speech Communication Association, 2017, pp. 2650–2654.
 - [216] M. Maciejewski, G. Sell, L. P. Garcia-Perera, S. Watanabe, S. Khudanpur, Building corpora for single-channel speech separation across multiple domains, CoRR abs/1811.02641 (2018). URL: <http://arxiv.org/abs/1811.02641>. arXiv:1811.02641.
 - [217] S. Araki, N. Ono, K. Kinoshita, M. Delcroix, Meeting recognition with asynchronous distributed microphone array using block-wise refinement of mask-based MVDR beamformer, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2018, pp. 5694–5698.
 - [218] A. Stolcke, Making the most from multiple microphones in meeting recordings, in: Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing, 2011, pp. 4992–4995.
 - [219] S. Narayanan, P. G. Georgiou, Behavioral signal processing: Deriving human behavioral informatics from speech and language, Proceedings of the IEEE 101 (2013) 1203–1233.
 - [220] D. Bone, C.-C. Lee, T. Chaspari, J. Gibson, S. Narayanan, Signal processing and machine learning for mental health research and clinical applications, IEEE Signal Processing Magazine 34 (2017) 189–196.
 - [221] M. Kumar, S. H. Kim, C. Lord, S. Narayanan, Speaker diarization for naturalistic child-adult conversational interactions using contextual information, Journal of the Acoustical Society of America 147 (2020) EL196–EL200. doi:10.1121/10.0000736.
 - [222] P. G. Georgiou, M. P. Black, S. S. Narayanan, Behavioral signal processing for understanding (distressed) dyadic interactions: some recent developments, in: Proceedings of the joint ACM workshop on Human gesture and behavior understanding, 2011, pp. 7–12.
 - [223] B. Xiao, C. Huang, Z. E. Imel, D. C. Atkins, P. Georgiou, S. S. Narayanan, A technology prototype system for rating therapist empathy from audio recordings in addiction counseling, PeerJ Computer Science 2 (2016) e59.
 - [224] S. N. Chakravarthula, M. Nasir, S.-Y. Tseng, H. Li, T. J. Park, B. Baucom, C. J. Bryan, S. Narayanan, P. Georgiou, Automatic prediction

- of suicidal risk in military couples using multimodal interaction cues from couples conversations, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2020, pp. 6539–6543.
- [225] B. Mirheidari, D. Blackburn, K. Harkness, T. Walker, A. Venneri, M. Reuber, H. Christensen, Toward the automation of diagnostic conversation analysis in patients with memory complaints, *Journal of Alzheimer's Disease* 58 (2017) 373–387.
- [226] G. P. Finley, E. Edwards, A. Robinson, N. Sadoughi, J. Fone, M. Miller, D. Suendermann-Oeft, M. Brenndorfer, N. Axtmann, An automated assistant for medical scribes., in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2018, pp. 3212–3213.
- [227] A. Guo, A. Faria, J. Riedhammer, Remeeting – Deep insights to conversations, in: *Proceedings of the Annual Conference of the International Speech Communication Association*, 2016, pp. 1964–1965.
- [228] A. Addelee, Y. Yu, A. Eshghi, A comprehensive evaluation of incremental speech recognition and diarization for conversational ai, in: *Proceedings of the International Conference on Computational Linguistics*, 2020, pp. 3492–3503.
- [229] O. Cetin, E. Shriberg, Speaker overlaps and ASR errors in meetings: Effects before, during, and after the overlap, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, volume 1, IEEE, 2006, pp. 357–360.
- [230] N. Kanda, C. Boeddeker, J. Heitkaemper, Y. Fujita, S. Horiguchi, K. Nagamatsu, R. Haeb-Umbach, Guided source separation meets a strong ASR backend: Hitachi/Paderborn University joint investigation for dinner party ASR, *Proceedings of the Annual Conference of the International Speech Communication Association* (2019) 1248–1252.
- [231] S. Otterson, M. Ostendorf, Efficient use of overlap information in speaker diarization, in: *Proceedings of IEEE Workshop on Automatic Speech Recognition and Understanding*, IEEE, 2007, pp. 683–686.
- [232] K. Boakye, B. Trueba-Hornero, O. Vinyals, G. Friedland, Overlapped speech detection for improved speaker diarization in multiparty meetings, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2008, pp. 4353–4356.
- [233] L. Bullock, H. Bredin, L. P. Garcia-Perera, Overlap-aware diarization: Resegmentation using neural end-to-end overlapped speech detection, in: *Proceedings of IEEE International Conference on Acoustics, Speech and Signal Processing*, IEEE, 2020, pp. 7114–7118.