

Neural Target Speech Extraction: An Overview

Katerina Zmolikova, Marc Delcroix, Tsubasa Ochiai, Keisuke Kinoshita, Jan Černocký, Dong Yu

Abstract—Humans can listen to a target speaker even in challenging acoustic conditions that have noise, reverberation, and interfering speakers. This phenomenon is known as the cocktail-party effect. For decades, researchers have focused on approaching the listening ability of humans. One critical issue is handling interfering speakers because the target and non-target speech signals share similar characteristics, complicating their discrimination. Target speech/speaker extraction (TSE) isolates the speech signal of a target speaker from a mixture of several speakers with or without noises and reverberations using clues that identify the speaker in the mixture. Such clues might be a spatial clue indicating the direction of the target speaker, a video of the speaker’s lips, or a pre-recorded enrollment utterance from which their voice characteristics can be derived. TSE is an emerging field of research that has received increased attention in recent years because it offers a practical approach to the cocktail-party problem and involves such aspects of signal processing as audio, visual, array processing, and deep learning. This paper focuses on recent neural-based approaches and presents an in-depth overview of TSE. We guide readers through the different major approaches, emphasizing the similarities among frameworks and discussing potential future directions.

Index Terms—Speech processing, target speech extraction, speech enhancement, multi-modal, deep learning

I. INTRODUCTION

In everyday life, we are constantly immersed in complex acoustic scenes consisting of multiple sounds, such as a mixture of speech signals from multiple speakers and background noise from air-conditioners or music. Humans naturally extract relevant information from such noisy signals as they enter our ears. The cocktail-party problem is a typical example [1], where we can follow the conversation of a speaker of interest (target speaker) in a noisy room with multiple interfering speakers. Humans can manage this complex task due to selective attention or a selective hearing mechanism that allows us to focus our attention on a target speaker’s voice and ignore others. Although the mechanisms of human selective hearing are not fully understood yet, many studies have identified essential cues exploited by humans to attend to a target speaker in a speech mixture: spatial, spectral (audio), visual, or semantic cues [1]. One long-lasting goal of speech processing research is designing machines that can achieve similar listening abilities as humans, i.e., selectively extracting the speech of a desired speaker based on auxiliary cues.

In this paper, we present an overview of recent developments in target speech/speaker extraction (TSE), which estimates the speech signal of a target speaker in a mixture

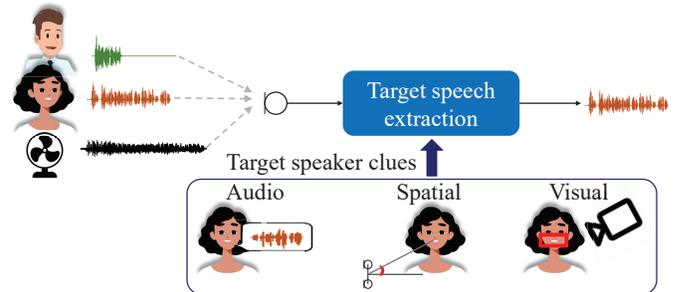


Fig. 1. TSE problem and examples of clues

of several speakers, given auxiliary cues to identify the target¹. In the following, we call auxiliary cues, clues, since they represent hints for identifying the target speaker in the mixture. Fig. 1 illustrates the TSE problem and shows that by exploiting the clues, TSE can focus on the voice of the target speaker while ignoring other speakers or noise. Inspired by psychoacoustic studies [1], several clues have been explored to tackle the TSE problem, such as spatial clues that provide the direction of the target speaker [2], [3], visual clues from video of their face [4]–[9], or audio clues extracted from pre-recorded enrollment recording of their voice [10]–[12].

The TSE problem is directly related to human selective hearing, although we approach it from an engineering point of view and do not try to precisely mimic human mechanisms. TSE is related to other speech and audio-processing tasks such as noise reduction and blind source separation (BSS) that do not use clues about the target speaker. Although noise reduction does suppress the background noise, it cannot handle well interfering speakers. BSS estimates each speech source signal in a mixture, which usually requires estimating the number of sources, a step that is often challenging. Moreover, it estimates the source signals without identifying them, which leads to global permutation ambiguity at its output; it remains ambiguous which of the estimated source signals corresponds to the target speaker. In contrast, TSE focuses on the target speaker’s speech by exploiting clues without assuming knowledge of the number of speakers in the mixture and avoids global permutation ambiguity. It thus offers a practical alternative to noise reduction or BSS when the use case requires focusing on a desired speaker’s voice.

Solving the TSE problem promises real implications for the development of many applications: (1) robust voice user interfaces or voice-controlled smart devices that only respond to a specific user; (2) teleconferencing systems that can remove

Katerina Zmolikova and Jan Černocký are with Brno University of Technology, Speech@FIT. Marc Delcroix, Tsubasa Ochiai and Keisuke Kinoshita are with NTT Corporation. Dong Yu is with Tencent, AI Lab.

¹Alternative terms in the literature for TSE include informed source separation, personalized speech enhancement, or audio-visual speech separation, depending on the context and the modalities involved.

interfering speakers close by; (3) hearing aids/hearables that can emphasize the voice of a desired interlocutor.

TSE ideas can be traced back to early works on beamformers [2]. Several works also extended BSS approaches to exploit clues about the target speaker [4], [5], [12]. Most of these approaches required a microphone array [5] or models trained on a relatively large amount of speech data from the target speaker [4]. The introduction of neural networks (NNs) enabled the building of powerful models that learn to perform complex conditioning on various clues by leveraging large amounts of speech data of various speakers. This evolution resulted in impressive extraction performance. Moreover, neural TSE systems can operate with a single microphone and with speakers unseen during the training of the models, allowing more flexibility.

This overview paper covers recent TSE development and focuses on neural approaches. Its remaining sections are organized as follows. In Section II, we formalize the TSE problem and its relation to noise reduction and BSS and introduce its historical context. We then present a taxonomy of TSE approaches and motivate the focus of this overview paper in Section III. We describe a general neural TSE framework in Section IV. The later sections (V, VI, and VII) introduce implementations of TSE with different clues, such as audio, visual, and spatial clues. We discuss extensions to other tasks in Section VIII. Finally, we conclude by describing the outlook on remaining issues in Section IX and provide pointers to available resources for experimenting with TSE in Section X.

II. PROBLEM DEFINITION

A. Speech recorded with a distant microphone

Imagine recording a target speaker's voice in a living room using a microphone placed on a table. This scenario represents a typical use case of a voice-controlled smart device or a video-conferencing device in a remote-work situation. Many sounds may co-occur while the speaker is speaking, e.g., a vacuum cleaner, music, children screaming, voices from another conversation, or from a TV. The speech signal captured at a microphone thus consists of a mixture of the target speaker's speech and interference from the speech of other speakers and background noise². We can express the mixture signal recorded at a microphone as

$$\mathbf{y}^m = \mathbf{x}_s^m + \underbrace{\sum_{k \neq s} \mathbf{x}_k^m}_{\triangleq \mathbf{i}^m} + \mathbf{v}^m, \quad (1)$$

where $\mathbf{y}^m = [y^m[0], \dots, y^m[T]] \in \mathbb{R}^T$, $\mathbf{x}_s^m \in \mathbb{R}^T$, $\mathbf{x}_k^m \in \mathbb{R}^T$, and $\mathbf{v}^m \in \mathbb{R}^T$ are the time-domain signal of the mixture, the target speech, the interference speech, and noise signals, respectively. Variable T represents the duration (number of samples) of the signals, m is the index of the microphone in an array of microphones, s represents the index of the target speaker and k is the index for the other speech sources.

²In this paper, we do not explicitly consider the effect of reverberation caused by the reflection of sounds on the walls and surfaces in a room, which also corrupt the recorded signal. Some of the approaches we discussed implicitly handle reverberation.

We drop microphone index m whenever we deal with single microphone approaches. In the TSE problem, we are interested in only recovering the target speech of speaker s , \mathbf{x}_s^m , and view all the other sources as undesired signals to be suppressed. We can thus define the interference signal as $\mathbf{i}^m \in \mathbb{R}^T$. Note that we make no explicit hypotheses about the number of interfering speakers.

B. TSE problem and its relation to BSS and noise reduction

The TSE problem is to estimate the target speech, given a clue, \mathbf{C}_s , as

$$\hat{\mathbf{x}}_s = \text{TSE}(\mathbf{y}, \mathbf{C}_s; \theta^{\text{TSE}}), \quad (2)$$

where $\hat{\mathbf{x}}_s$ is the estimate of the target speech, $\text{TSE}(\cdot; \theta^{\text{TSE}})$ represents a TSE system with parameters θ^{TSE} . The clue, \mathbf{C}_s , allows identifying the target speaker in the mixture. It can be of various types, such as a pre-recorded enrollment utterance, $\mathbf{C}_s^{(a)}$, a video signal capturing the face or lips movements of the target speaker, $\mathbf{C}_s^{(v)}$, or such spatial information as the direction of arrival (DOA) of the speech of the target speaker, $\mathbf{C}_s^{(d)}$.

In the later sections, we expand on how to design TSE systems. Here, we first emphasize the key difference between TSE and BSS and noise reduction. Fig. 2 compares these three problems.

BSS [13], [14] estimates all the source signals in a mixture without requiring clues:

$$\{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_K\} = \text{BSS}(\mathbf{y}; \theta^{\text{BSS}}), \quad (3)$$

where $\text{BSS}(\cdot; \theta^{\text{BSS}})$ represents a separation system with parameters θ^{BSS} , $\hat{\mathbf{x}}_k$ are the estimates of the speech sources, and K is the number of sources in the mixture. As seen in Eq. (3), BSS does not and cannot differentiate the target speech from other speech sources. Therefore, we cannot know in advance which output corresponds to the target speech, i.e., there is a global permutation ambiguity problem between the outputs and the speakers. Besides, since the number of outputs is given by the number of sources, the number of sources K must be known or estimated. Comparing Eqs. (2) and (3) emphasizes the fundamental difference between TSE and BSS: (1) TSE estimates only the target speech signal, while BSS estimates all the signals, and (2) TSE is conditioned on speaker clue \mathbf{C}_s , while BSS only relies on the observed mixture³. Typical use cases for BSS include applications that require estimating speech signals of every speaker, such as automatic meeting transcription systems.

Noise reduction is another related problem. It assumes that the interference only consists of background noise, i.e., $\mathbf{i} = \mathbf{v}$, and can thus enhance the target speech without requiring clues:

$$\hat{\mathbf{x}}_s = \text{Denoise}(\mathbf{y}; \theta^{\text{Denoise}}), \quad (4)$$

where $\text{Denoise}(\cdot; \theta^{\text{Denoise}})$ represents a noise reduction system with parameters θ^{Denoise} . Unlike BSS, a noise reduction system's output only consists of target speech $\hat{\mathbf{x}}_s$, and there

³Another setup sitting between TSE and BSS is a task that extracts multiple target speakers, e.g., extracting the speech of all the meeting attendees given such information about them as enrollment or videos of all the speakers.

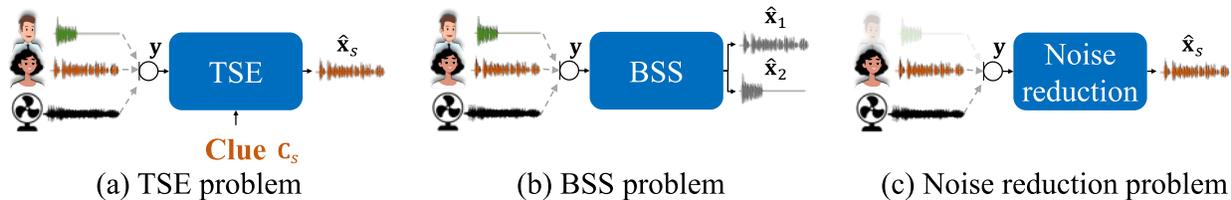


Fig. 2. Comparison of TSE with BSS and noise reduction

is thus no global permutation ambiguity. This is possible if the background noise and speech have distinct characteristics. For example, we can assume that ambient noise and speech signals exhibit different spectro-temporal characteristics that enable their discrimination. However, noise reduction cannot suppress interfering speakers because it cannot discriminate among different speakers in a mixture without clues⁴. Noise reduction is often used, e.g., in video-conferencing systems or hearing aids.

TSE is an alternative to BSS and noise reduction, which uses a clue to simplify the problem. Like BSS, it can handle speech mixtures. Like noise reduction, it only estimates the target speaker, thus avoiding global permutation ambiguity and the need to estimate the number of sources. However, TSE requires access to clues, unlike BSS and noise reduction. Moreover, it must internally perform two sub-tasks: (1) identifying the target speaker and (2) estimating the speech of that speaker in the mixture. TSE is thus a challenging problem that introduces specific issues and requires dedicated solutions.

A straightforward way to achieve TSE using BSS methods is to first apply BSS and next select the target speaker among the estimated sources. Such a cascade system allows the separate development of BSS and speaker identification modules. However, this scheme is usually computationally more expensive and imports some disadvantages of BSS, such as the need to estimate the number of speakers in the mixture. Therefore, we focus on approaches that directly exploit the clues in the extraction process. Nevertheless, most TSE research is rooted in BSS, as argued in the following discussion on the historical context.

C. Historical context

The first studies related to TSE were performed in the 1980s. Flanagan et al. [2] explored enhancing a target speaker’s voice in a speech mixture, assuming that the target speech originated from a fixed and known direction. They employed a microphone array to record speech and designed a fixed beamformer that enhanced the signals from the target direction [2], [16]. We consider that this work represents an early TSE system that relies on spatial clues.

In the mid-1990s, the BSS problem gained attention with pioneering works on independent component analysis (ICA).

⁴Some works propose to exploit clues for noise reduction and apply similar ideas of TSE to reduce background noise (and sometimes interfering speakers). In the literature, this is called personalized speech enhancement, which in this paper, we view as a special case of the TSE problem, where only the target speaker is actively speaking [15].

ICA estimates spatial filters that separate the sources by relying on the assumption of the independence of the sources in the mixture and the fact that speech signals are non-Gaussian [13]. A frequency-domain ICA suffers from a frequency permutation problem because it treats each frequency independently. In the mid-2000s, independent vector analysis (IVA) addressed the frequency-permutation problem by working on vectors spanning all frequency bins, which allowed modeling dependency among frequencies [13]. Several works have extended ICA and IVA to perform TSE, which simplifies inference by focusing on a single target source. For example, in the late 2000s, TSE systems were designed by incorporating the voice activity information of the target speaker derived from video signals to the ICA criterion, allowing identification and extraction of only the target source [5]. In the late 2010s, independent vector extraction (IVE) extended IVA to extract a single source out of the mixture. In particular, IVE exploits clues to guide the extraction process, such as the enrollment of the target speaker to achieve TSE [12]. All these approaches require a microphone array to capture speech.

In the first decade of the 2000s, single-channel approaches for BSS emerged, such as factorial hidden Markov model (F-HMM) [17] and non-negative matrix factorization (NMF) [18]. These approaches relied on pre-trained spectral models of speech signals learned on clean speech data. An F-HMM is a model of speech mixtures, where the speech of each speaker in the mixture is explicitly modeled using a separate hidden Markov model (HMM). The parameters of each speaker-HMM are learned on the clean speech data of that speaker. The separation process involves inferring the most likely HMM state sequence associated with each speaker-HMM, which requires approximations to make inference tractable. This approach was the first to achieve super-human performance using only single-channel speech [17]. In the early 2000s, F-HMM was also among the first approaches to exploit visual clues [4]⁵. In NMF, the spectrogram of each source is modeled as a multiplication of pre-learned bases, representing the basic spectral patterns and their time-varying activations. NMF methods have also been extended to multi-channel signals [13] and used to extract a target speaker [19] with a flexible multi-source model of the background. The main shortcoming of the F-HMM and NMF methods is that they require pre-trained source models and thus struggle with unseen speakers. Furthermore,

⁵This framework needs having clues for all of the speakers, a requirement that negates some of the advantages of TSE, e.g., the number of speakers must be known beforehand. Despite that, the method does not suffer from global permutation ambiguity, since visual clues identify the target speaker, and we thus include this work in the broader view of TSE methods.

the inference employs a computationally expensive iterative optimization.

In the mid-2010s, deep NNs (DNNs) were first introduced to address the BSS problem. These approaches rapidly gained attention with the success of deep-clustering and permutation invariant training (PIT) [20], [21], which showed that single-channel speaker-open⁶ BSS was possible. In particular, the introduction of DNNs enabled more accurate and flexible spectrum modeling and computationally efficient inference. These advances were facilitated by supervised training methods that can exploit a large amount of data.

Neural BSS rapidly influenced TSE research. For example, Du et al. [22] trained a speaker-close NN to extract the speech of a target speaker using training data with mixed various interfering speakers. This work is an initial neural TSE system using audio clues. However, using speaker-close models requires a significant amount of data from the target speaker and cannot be extended to speakers unseen during training. Subsequently, the introduction of TSE systems conditioned on speaker characteristics derived from an enrollment utterance significantly mitigated this requirement [10], [11], [23]. Enrollment consists of a recording of a target speaker’s voice, which amounts to a few seconds of speech. With these approaches, audio clue-based TSE became possible for speakers unseen during training as long as an enrollment utterance was available. Furthermore, the flexibility of NNs to integrate different modalities combined with the high modeling capability of face recognition or lip-reading systems offered new possibilities for speaker-open visual clue-based TSE [7], [8]. More recently, neural approaches have also been introduced for spatial-clue-based TSE [3], [24].

TSE has gained increased attention. For example, dedicated tasks were part of such recent evaluation campaigns as the deep noise suppression (DNS)⁷ and Clarity⁸ challenges. Many works have extended TSE to other tasks, such as a direct automatic speech recognition (ASR) of a target speaker from a mixture, which is called target speaker ASR (TS-ASR) [25], [26], or personalized voice activity detection (VAD)/diarization [27], [28]. Notably, target speaker VAD (TS-VAD)-based diarization [28] has been very successful in such evaluation campaigns as CHiME-6⁹ or DIHARD-3¹⁰, outperforming state-of-the-art diarization approaches in challenging conditions.

III. TSE TAXONOMY

TSE is a vast research area spanning a multitude of approaches. This section organizes them to emphasize their relations and differences. We categorized the techniques using four criteria: 1) type of clues, 2) number of channels, 3) speaker-close vs. open, and 4) generative vs. discriminative. Table I summarizes the taxonomy; the works in the scope of this overview paper are emphasized in red.

⁶BSS is possible for speakers unseen during training, i.e., not present in the training data.

⁷<https://www.microsoft.com/en-us/research/academic-program/deep-noise-suppression-challenge-icassp-2022/>

⁸https://claritychallenge.github.io/clarity_CC_doc

⁹<https://chimechallenge.github.io/chime6/results.html>

¹⁰<https://dihardchallenge.github.io/dihard3/results>

A. Type of clue

The type of clue used to determine the target speaker is an important factor in distinguishing among TSE approaches. The most prominent types are audio, visual, and spatial clues. This classification also defines the main organization of this article, which covers such approaches in Sections V, VI, and VII. Other types have and could be proposed, as we briefly discuss in Section IX.

An *audio clue* consists of a recording of a speech signal of the target speaker. Such a clue can be helpful, e.g., in the use case of personal devices, where the user can pre-record an example of their voice. Alternatively, for long recordings, such as meetings, clues can be obtained directly from part of the recording. The interest in audio clues sharply increased recently with the usage of neural models for TSE [10]–[12]. Audio clues are perhaps the most universal, because they do not require using any additional devices, such as multiple microphones or a camera. However, the performance may be limited compared to other clues, since discriminating speakers based only on their voice characteristics is prone to errors due to inter- and intra-speaker variability. For example, the voice characteristics of different speakers, such as family members, often closely resemble each other. On the other hand, the voice characteristics of one speaker may change depending on such factors as emotions, health, or age.

A *visual clue* consists of a video of the target speaker talking. This type is often constrained to the speaker’s face, sometimes just to the lip area. Unlike audio clues, visual clues are typically synchronized with audio signals that are processed, i.e., not pre-recorded. A few works also explored just using a photo of the speaker [37]. Visual clues have been employed to infer the activity pattern and location of the target speaker [5] or to jointly model audio and visual signals [4], [5]. Recent works usually use visual clues to guide discriminative models toward extracting the target speaker [7]–[9]. Visual clues are especially useful when speakers in the recording have similar voices [8]. However, they might be sensitive to physical obstructions of the speaker in the video.

A *spatial clue* refers to the target speaker’s location, e.g., the angle from the recording devices. The location can be inferred in practice from a video of the room or a recording of a speaker in the same position. Extracting the speaker based on their location has been researched from mid 1980’s, with beamforming techniques that pioneered this topic [2], [16]. More recent IVE models use location for initialization [12]. Finally, several works have shown that NNs informed by location can also achieve promising performance [3], [24]. Spatial clues are inherently applicable only when a recording from multiple microphones is available. However, they can identify the target speaker in the mixture rather reliably, especially when the speakers are stationary.

Different clues may work better in different situations. For example, the performance with audio clues might depend on the similarity of voices of the present speakers, and obstructions in the video may influence visual clues. As such, it is advantageous to use multiple clues simultaneously to combine their strengths. Many works have combined audio and visual

TABLE I
TAXONOMY OF TSE WORKS: APPROACHES WITHIN SCOPE OF THIS OVERVIEW PAPER ARE EMPHASIZED IN RED.

	Representative approaches	References	Year	Type of clues			Number of mic.		Speaker-close/open	
				Audio	Visual	Spatial	Single	Multi	Close	Open
	Fixed beamforming	[2], [16] ¹¹	1985	-	-	✓	-	✓	-	✓
Generative	Audio-visual F-HMM	[4]	2001	✓ ¹²	✓	-	✓	-	✓	-
	ICA with visual voice activity	[5]	2007	-	✓	-	-	✓	-	✓
	Multi-channel NMF	[19]	2011	✓ ¹²	-	-	-	✓	✓	-
	IVE with x-vectors	[12]	2020	✓	-	-	-	✓	-	✓
	Audio-visual VAE	[29]	2020	-	✓	-	✓	-	-	✓
Discriminative	Speaker-specific network	[22]	2014	✓ ¹²	-	-	✓	-	✓	-
	Multi-channel SpeakerBeam	[10], [30]	2017	✓	-	-	-	✓	-	✓
	SpeakerBeam	[10]	2019	✓	-	-	✓	-	-	✓
	VoiceFilter	[11]	2019	✓	-	-	✓	-	-	✓
	SpEx	[31]	2020	✓	-	-	✓	-	-	✓
	The conversation	[7]	2018	-	✓	-	✓	-	-	✓
	Looking-to-listen	[8]	2018	-	✓	-	✓	-	-	✓
	On/off-screen audio-visual separation	[9]	2018	-	✓	-	✓	-	-	✓
	Landmark-based AV speech enh.	[32]	2019	-	✓	-	✓	-	-	✓
	Multi-modal SpeakerBeam	[33], [34]	2019	✓	✓	-	✓	-	-	✓
	AV speech enh. through obstructions	[35]	2019	✓	✓	-	✓	-	-	✓
	Neural spatial filter	[3]	2019	✓	-	✓	-	✓	-	✓
	Spatial speaker extractor	[24]	2019	✓	-	✓	-	✓	-	✓
	Multi-channel multi-modal TSE	[36]	2020	✓	✓	✓	-	✓	-	✓

clues [4], [33], and some have even added spatial clues [36].

B. Number of microphones

Another way to categorize the TSE approaches is based on the number of microphones (channels) they use. Multiple channels allow the spatial diversity of the sources to be exploited to help discriminate the target speaker from interference. Such an approach also closely follows human audition, where binaural signals are crucial for solving the cocktail-party problem.

All approaches with spatial clues require using a microphone array to capture the direction information of the sources in the mixture [2], [3], [16], [24], [36]. Some TSE approaches that exploit audio or visual clues also assume multi-channel recordings, such as the extensions of ICA/IVA approaches [5], [12].

Multi-channel approaches generally generate extracted signals with better quality and are thus preferable when recordings from a microphone array are available. However, sometimes they might fail when the sources are located in the same direction from the viewpoint of the recording device. Moreover, adopting a microphone array is not always an option when developing applications due to cost restrictions. In such cases, single-channel approaches are requested. They rely on spectral models of speech mixture using either F-HMM or recently NNs and exploit audio [10], [11] or visual clues [7], [8] to identify the target speech.

Recent single-channel neural TSE systems have achieved remarkable performance. Interestingly, such approaches can also be easily extended to multi-channel processing by augmenting the input with spatial features [3] or combining the processing

with beamforming [24], [30], as discussed in Section IV-C. For example, using a beamformer usually extracts a higher quality signal due to employing a spatial linear filter to perform extraction, which can benefit ASR applications [10].

C. Speaker-open vs speaker-close methods

We usually understand the clues used by TSE as short evidence about the target speaker obtained at the time of executing the method, e.g., one utterance spoken by the target speaker, a video of him/her speaking, or their current location. There are, however, also methods that use a more significant amount of data from the target speaker (e.g., several hours of their speech) to build a model specific to that person. These methods can also be seen as TSE except that the clues involve much more data.

We refer to these two categories as the speaker-open and speaker-close methods¹³. In speaker-open methods, the data of the target speaker are available only during the test time, i.e., the model is trained on the data of different speakers. In contrast, the target speaker is part of the training data in speaker-close methods. Many methods in the past were speaker-close, e.g., [4] or [19], where the models were trained on the clean utterances of the target speaker. Also, the first neural models for TSE used a speaker-specific network [22]. Most recent works on neural methods, which use a clue as an additional input, are speaker-open methods [3], [7], [8], [10], [11]. Recent IVE methods [12] are also speaker-open, i.e., they guide the inference of IVE using the embedding of a previously unseen speaker.

¹¹Since the first works that proposed beamforming were not model-based, we consider them neither generative nor discriminative.

¹²In speaker-close cases, the models are trained on target speaker's audio. We consider this an audio clue in Table I.

¹³Speaker-open and speaker-close categories are sometimes referred to as speaker-independent and speaker-dependent, respectively. We avoid this terminology, as in TSE, all systems are informed about the target speaker, and therefore the term speaker-independent might be misleading.

D. Generative vs discriminative

We can classify TSE into approaches using generative or discriminative models.

Generative approaches model the joint distribution of the observations, target signals, and clues. The estimated target speech is obtained by maximizing the likelihood. In contrast, discriminative approaches directly estimate the target speech signal given observations and clues.

In the TSE literature, generative models were the dominant choice in the pioneering works, including one [4] that used HMMs to jointly model audio and visual modalities. IVE [12] is also based on a generative model of the mixtures.

The popularity of discriminative models, in particular NNs, has increased since mid-2010's, and such models today are the choice for many problems, including TSE. With discriminative models, TSE is treated as a supervised problem, where the parameters of a TSE model are learned using artificially generated training data. The modeling power of NNs enables us to exploit large amounts of such data to build strong speech models. Moreover, the versatility of NNs enables complex dependencies to be learned between different types of observations (e.g., speech mixture and video/speaker embeddings), which allows the successful conditioning of the extraction process on various clues. However, NNs also bring new challenges, such as generalization to unseen conditions or high computational requirements [38].

Some recent works have also explored using generative NNs, such as variational autoencoders (VAEs) [29], which might represent a middle-ground between the traditional generative approaches and those using discriminative NNs.

E. Scope of overview paper

In the remainder of our paper, we focus on the neural methods for TSE emphasized in Table I. Recent neural TSE approaches opened the possibility of achieving high-performance extraction with various clues. They can be operated with a single microphone and applied for speaker-open conditions, which are very challenging constraints for other schemes. Consequently, these approaches have received increased attention from both academia and industry.

In the next section, we introduce a general framework to provide a uniformized view of the various NN-based TSE approaches, for both single- and multi-channel approaches, and independently of the type of clues. We then respectively review the approaches relying on audio, visual, and spatial clues in Sections V, VI, and VII.

IV. GENERAL FRAMEWORK FOR NEURAL TSE

In the previous section, we introduced a taxonomy that described the diversity of approaches to tackle the TSE problem. However, recent neural TSE systems have much in common. In this section, we introduce a general framework that provides a unified view of a neural TSE system, which shares the same processing flow independently of the type of clue used. By organizing the existing approaches into a common framework, we hope to illuminate their similarities and differences and establish a firm foundation for future research.

A neural TSE system consists of an NN that estimates the target speech conditioned on a clue. Fig. 3 is a schematic diagram of a generic neural TSE system that consists of two main modules: a clue encoder and a speech extraction module, described in more detail below.

A. Clue encoder

The clue encoder pulls out (from the clue, \mathbf{C}_s) information that allows the speech extraction module to identify and extract the target speech in the mixture. We can express the processing as

$$\mathbf{E}_s = \text{ClueEncoder}(\mathbf{C}_s; \theta^{\text{Clue}}), \quad (5)$$

where $\text{ClueEncoder}(\cdot; \theta^{\text{Clue}})$ represents the clue encoder, which can be an NN with learnable parameters θ^{Clue} , and \mathbf{E}_s are the clue embeddings. Naturally, the specific implementation of the clue encoder and the information carried within \mathbf{E}_s largely depend on the type of clues. For example, when the clue is an enrollment utterance, $\mathbf{E}_s = \mathbf{E}_s^{(a)} \in \mathbb{R}^{D^{\text{Emb}}}$ will be a speaker embedding vector of dimension D^{Emb} that represents the voice characteristics of the target speaker. When dealing with visual clues, $\mathbf{E}_s = \mathbf{E}_s^{(v)} \in \mathbb{R}^{D^{\text{Emb}} \times N}$ can be a sequence of the embeddings of length N , representing, e.g., the lip movements of the target speaker. Here N represents the number of time frames of the mixture signal.

Interestingly, the implementation of the speech extraction module does not depend on the type of clues used. To provide a description that is independent of the type of clues, hereafter, we consider that $\mathbf{E}_s \in \mathbb{R}^{D^{\text{Emb}} \times N}$ consists of a sequence of embedding vectors of dimension D^{Emb} of length N . Note that we can generate a sequence of embedding vectors for audio clue-based TSE systems by repeating the speaker embedding vector for each time frame.

B. Speech extraction module

The speech extraction module estimates the target speech from the mixture, given the target speaker embeddings. We can use the same configuration independently of the type of clue. Its process can be decomposed into three main parts: a mixture encoder, a fusion layer, and a target extractor:

$$\mathbf{Z}_y = \text{MixEncoder}(\mathbf{y}; \theta^{\text{Mix}}), \quad (6)$$

$$\mathbf{Z}_s = \text{Fusion}(\mathbf{Z}_y, \mathbf{E}_s; \theta^{\text{Fusion}}), \quad (7)$$

$$\hat{\mathbf{x}}_s = \text{TgtExtractor}(\mathbf{Z}_s, \mathbf{y}; \theta^{\text{TgtExtractor}}), \quad (8)$$

where $\text{MixEncoder}(\cdot; \theta^{\text{Mix}})$, $\text{Fusion}(\cdot; \theta^{\text{Fusion}})$, and $\text{TgtExtractor}(\cdot; \theta^{\text{TgtExtractor}})$ respectively represent the mixture encoder, the fusion layer, and the target extractor with parameters θ^{Mix} , θ^{Fusion} , and $\theta^{\text{TgtExtractor}}$. $\mathbf{Z}_y \in \mathbb{R}^{D^y \times N}$ and $\mathbf{Z}_s \in \mathbb{R}^{D^s \times N}$ are the internal representations of the mixture before and after conditioning on embedding \mathbf{E}_s .

The mixture encoder performs the following:

$$\mathbf{Y} = \text{FE}(\mathbf{y}; \theta^{\text{FE}}), \quad (9)$$

$$\mathbf{Z}_y = \text{MixNet}(\mathbf{Y}; \theta^{\text{MixNet}}), \quad (10)$$

where $\text{FE}(\cdot)$ and $\text{MixNet}(\cdot)$ respectively represent the feature extraction process and an NN with parameters θ^{FE} and θ^{MixNet} .

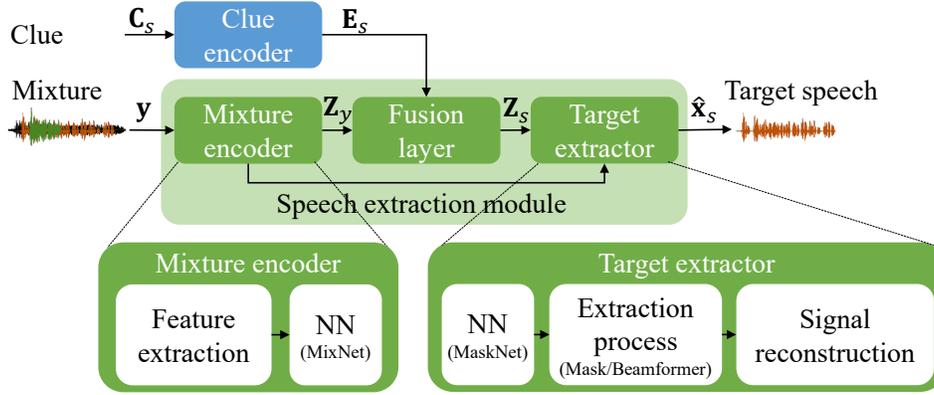


Fig. 3. General framework for neural TSE

TABLE II

TYPE OF FUSION LAYERS: \mathbf{L} , \mathbf{L}_1 , AND \mathbf{L}_2 ARE LINEAR TRANSFORMATIONS FOR MAPPING THE DIMENSION OF THE CLUE EMBEDDINGS, D^{Emb} , TO THE DIMENSION OF \mathbf{Z}_y , D^Z . \odot REPRESENTS THE ELEMENT-WISE HADAMARD MULTIPLICATION OPERATION OF MATRICES. \mathbf{e}_i IS A VECTOR CONTAINING THE ELEMENTS OF THE i -TH ROW OF \mathbf{E}_s AND $\text{diag}(\cdot)$ IS AN OPERATOR THAT CONVERTS A VECTOR INTO A DIAGONAL MATRIX.

Fusion type	Equation	Parameters (θ^{Fusion})
Concatenation	$\mathbf{Z}_s = [\mathbf{Z}_y, \mathbf{E}_s]$	-
Addition	$\mathbf{Z}_s = \mathbf{Z}_y + \mathbf{L}\mathbf{E}_s$	$\mathbf{L} \in \mathbb{R}^{D^Z \times D^{\text{Emb}}}$
Multiplication	$\mathbf{Z}_s = \mathbf{Z}_y \odot (\mathbf{L}\mathbf{E}_s)$	$\mathbf{L} \in \mathbb{R}^{D^Z \times D^{\text{Emb}}}$
Feature-wise Linear Modulation (FiLM)	$\mathbf{Z}_s = \mathbf{Z}_y \odot (\mathbf{L}_1\mathbf{E}_s) + \mathbf{L}_2\mathbf{E}_s$	$\mathbf{L}_1 \in \mathbb{R}^{D^Z \times D^{\text{Emb}}}$, $\mathbf{L}_2 \in \mathbb{R}^{D^Z \times D^{\text{Emb}}}$
Factorized layer	$\mathbf{Z}_s = \sum_{i=1}^{D^{\text{Emb}}} \mathbf{L}_i \mathbf{Z}_y \text{diag}(\mathbf{e}_i)$	$\mathbf{L}_i \in \mathbb{R}^{D^Z \times D^Z}$

The feature extractor computes the features from the observed mixture signal, $\mathbf{Y} \in \mathbb{R}^{D \times N}$. These can be such spectral features as magnitude spectrum coefficients derived from the short-time Fourier transform (STFT) of the input mixture [7], [8], [10], [11]. When using a microphone array, spatial features like interaural phase difference (IPD) defined in Eq. (21) in Section VII can also be appended. Alternatively, the feature extraction process can be implemented by an NN such as a 1-D convolutional layer that operates directly on the raw input waveform of the microphone signal [23], [39]. This enables learning of a feature representation optimized for TSE tasks.

The features are then processed with an NN, MixNet(\cdot), which performs a non-linear transformation and captures the time context, i.e., several past and future frames of the signal. The resulting representation, \mathbf{Z}_y , of the mixture is (at this point) agnostic of the target.

The fusion layer, sometimes denoted as an adaptation layer, is a key component of a TSE system and allows conditioning of the process on the clue. It combines \mathbf{Z}_y with the clue embeddings, \mathbf{E}_s . Conditioning an NN on auxiliary information is a general problem that has been studied for multi-modal processing or the speaker adaptation of ASR systems. TSE systems have borrowed fusion layers from these fields. Table II lists several options for the fusion layer. Some widely used fusion layers include: (1) the concatenation of \mathbf{Z}_y with the clue embeddings \mathbf{E}_s [7], [8]; (2) addition¹⁴ after transforming the embeddings with linear transformation \mathbf{L} to match the dimension of \mathbf{Z}_y ; (3) multiplication [10]; (4) a combination of

addition and multiplication denoted as FiLM; (5) a factorized layer [10], [30], i.e., the combination of different transformations of the mixture representation weighted by the clue embedding values. Other alternatives have also been proposed, including attention-based fusion [40]. Note that the fusion operations described here assume just one clue. It is also possible to use multiple clues, as discussed in Section VI-B. Some works also employ the fusion repeatedly at multiple positions in the model [31].

The last part of the speech extraction module is the target extractor, which estimates the target signal. We explain below the time-frequency masking-based extractor, which has been widely used [3], [7], [8], [41]. Recent approaches also perform a similar masking operation in the learned feature domain [23], [39].

The time-frequency masking approach was inspired by early BSS studies that relied on the sparseness assumption of speech signals, an idea based on the observation that the energy of a speech signal is concentrated in a few time-frequency bins of a speech spectrum. Accordingly, the speech signals of different speakers rarely overlap in the time-frequency domain in a speech mixture. We can thus extract the target speech by applying a time-frequency mask on the observed speech mixture, where the mask indicates the time-frequency bins where the target speech is dominant over other signals. Fig. 4 shows an example of an ideal binary mask for extracting a target speech in a mixture of two speakers. Such an ideal binary mask assumes that all the energy in each TF bin belongs to one speaker. In recent mask-based approaches that use real-

¹⁴Concatenation is similar to addition if a linear transformation follows it.

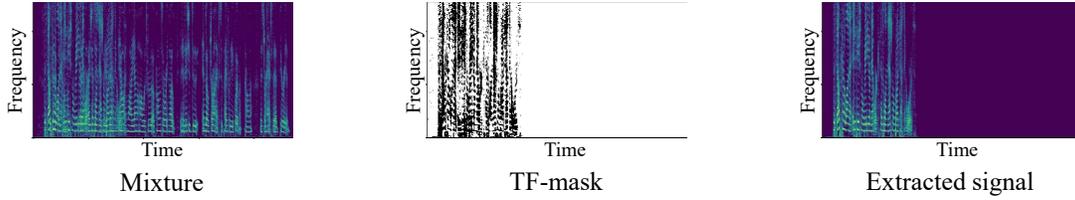


Fig. 4. Example of time-frequency mask for speech extraction: Time-frequency mask shows spectrogram regions where target source is dominant. By applying this mask to the mixture, we obtain an extracted speech signal that estimates the target speech.

valued (or complex) masks, this assumption or observation is not needed.

The processing of the masking-based extractor can be summarized as

$$\mathbf{M}_s = \text{MaskNet}(\mathbf{Z}_s; \theta^{\text{Mask}}), \quad (11)$$

$$\hat{\mathbf{X}}_s = \mathbf{M}_s \odot \mathbf{Y}, \quad (12)$$

$$\hat{\mathbf{x}}_s = \text{Reconstruct}(\hat{\mathbf{X}}_s; \theta^{\text{Reconst}}), \quad (13)$$

where $\text{MaskNet}(\cdot)$ is an NN that estimates the time-frequency mask for the target speech, $\mathbf{M}_s \in \mathbb{R}^{D \times N}$, θ^{Mask} are the network parameters, and \odot denotes the element-wise Hadamard multiplication. \mathbf{Y} and $\hat{\mathbf{X}}_s$ are the mixture and the estimated target speech signals in the feature domain. Eq. (12) shows the actual extraction process. $\text{Reconstruct}(\cdot)$ is an operation to reconstruct the time-domain signal by performing the inverse operation of the feature extraction of the mixture encoder, i.e., either inverse STFT (iSTFT) or a transpose convolution if using a learnable feature extraction. In the latter case, the reconstruction layer has learnable parameters, θ^{Reconst} .

There are other possibilities to perform the extraction process. For example, we can modify the $\text{MaskNet}(\cdot)$ NN to directly infer the target speech signal in the feature domain. Alternatively, as discussed in Section IV-C, we can replace the mask-based extraction process with beamforming when a microphone array is available.

C. Integration with microphone array processing

If we have access to a microphone array to record the speech mixture, we can exploit the spatial information to extract the target speech. One approach is to use spatial clues to identify the speaker in the mixture by informing the system about the target speaker's direction, as discussed in Section VII. Another approach combines TSE with beamforming and uses the latter to perform the extraction process instead of Eq. (12). For example, we can use the output of a TSE system to estimate the spatial statistics needed to compute the coefficients of a beamformer steering in the direction of the target speaker. This approach can also be used with audio or visual clue-based TSE systems and requires no explicit use of spatial clues to identify the target speaker in the mixture.

We briefly review the mask-based beamforming approach, which was introduced initially for noise reduction and BSS [42], [43]. A beamformer performs the linear spatial filtering of the observed microphone signals:

$$\hat{X}_s[n, f] = \mathbf{W}^H[f] \mathbf{Y}[n, f], \quad (14)$$

where $\hat{X}_s[n, f] \in \mathbb{C}$ is the STFT coefficient of the estimated target signal at time frame n and frequency bin f , $\mathbf{W}[f] \in \mathbb{C}^M$ is a vector of the beamformer coefficients, $\mathbf{Y}[n, f] = [Y^1[n, f], \dots, Y^M[n, f]]^T \in \mathbb{C}^M$ is a vector of the STFT coefficients of the microphone signals, M is the number of microphones, and H is the conjugate transpose. We can derive the beamformer coefficients from the spatial correlation matrices of the target speech and the interference. These correlation matrices can be computed from the observed signal and the time-frequency mask estimated by the TSE system [30].

This way of combining a TSE system with beamforming replaces the time-frequency masking operation of Eq. (12) with the spatial linear filtering operation of Eq. (14). It allows distortionless extraction, which is often advantageous when using TSE as a front-end for ASR [10].

D. Training a TSE system

Before using a TSE model, we first need to learn its parameters: $\theta^{\text{TSE}} = \{\theta^{\text{Mix}}, \theta^{\text{Clue}}, \theta^{\text{Fusion}}, \theta^{\text{TgtExtractor}}\}$. Most existing studies use fully supervised training, which requires a large amount of training data consisting of the triplets of speech mixture \mathbf{y} , target speech signal \mathbf{x}_s , and corresponding clue \mathbf{C}_s to learn parameters θ^{TSE} . Since this requires access to a clean target speech signal, such training data are usually simulated by artificially mixing clean speech signals and noise following the signal model of Eq. (1).

Figure 5 illustrates the data generation process using a multi-speaker audio-visual speech corpus containing multiple videos for each speaker. First, we generate a mixture using randomly selected speech signals from the target speaker, the interference speaker, and the background noise. We obtain an audio clue by selecting another speech signal from the target speaker as well as a visual clue from the video signal associated with the target speech.

The training of a neural TSE framework follows the training scheme of NNs with error back-propagation. The parameters are estimated by minimizing a training loss function:

$$\theta^{\text{TSE}} = \arg \min_{\theta} \mathcal{L}(\mathbf{x}_s, \hat{\mathbf{x}}_s), \quad (15)$$

where $\mathcal{L}(\cdot)$ is a training loss, which measures how close estimated target speech $\hat{\mathbf{x}}_s = \text{TSE}(\mathbf{y}, \mathbf{C}_s; \theta)$ is to the target source signal \mathbf{x}_s . We can use a similar loss as that employed for training noise reduction or BSS systems [14], [39].

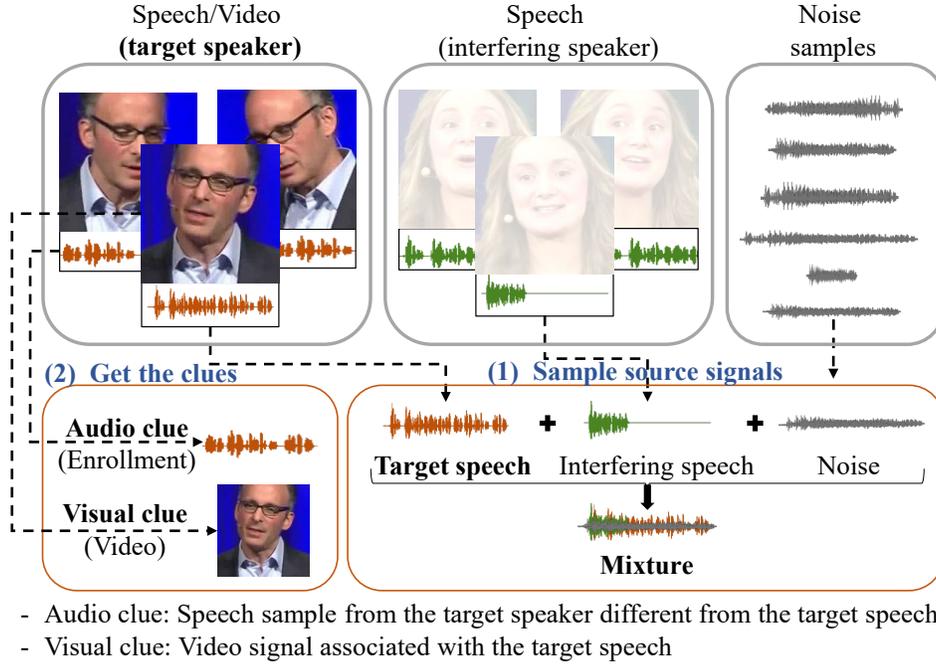


Fig. 5. Example of generating simulation data for training or testing: This example assumes videos are available so that audio and visual clues can be generated. No video is needed for audio clue-based TSE. For visual clue-based TSE, we do not necessarily need multiple videos from the same speaker.

Several variants of the losses operating on different domains exist, such as the cross-entropy between the oracle and the estimated time-frequency masks and the mean squared error (MSE) loss between the magnitude spectra of the source and the estimated target speech. Recently, a negative signal-to-noise ratio (SNR) measured in the time-domain has been widely used [6], [23], [39]:

$$\mathcal{L}^{\text{SNR}}(\mathbf{x}_s, \hat{\mathbf{x}}_s) = -10 \log_{10} \left(\frac{\|\mathbf{x}_s\|^2}{\|\mathbf{x}_s - \hat{\mathbf{x}}_s\|^2} \right). \quad (16)$$

The SNR loss is computed directly in the time-domain, which forces the TSE system to learn to correctly estimate the magnitude and the phase of the target speech signal. This loss has improved extraction performance [23]. Many works also employ versions of the loss which are invariant to arbitrary scaling, i.e., scale-invariant SNR (SI-SNR) [39] or linear filtering of the estimated signal, often called signal-to-distortion ratio (SDR) [44]. Besides training losses operating on the signal or mask levels, it is also possible to train a TSE system end-to-end with a loss defined on the output of an ASR system [45]. Such a loss can be particularly effective when targeting ASR applications, as discussed in Section VIII.

The clue encoder can be an NN trained jointly with a speech extraction module [10] or pre-trained on a different task, such as speaker identification for audio clue-based TSE [11] or lip-reading for visual clue-based TSE [7]. Using a pre-trained clue encoder enables the leveraging of large amounts of data to learn robust and highly discriminative embeddings. On the other hand, jointly optimizing the clue encoder allows learning embeddings to be optimized directly for TSE. These two trends can also be combined by fine-tuning the pre-trained encoder

or using multi-task training schemes, which add a loss to the output of the clue embeddings [46].

E. Considerations when designing a TSE system

We conclude this section with some considerations about the different options for designing a TSE system. In the above description, we intentionally ignored the details of the NN architecture used in the speech extraction module, such as the type of layers. Indeed, novel architectures have been and will probably continue to be proposed regularly, leading to gradual performance improvement. For concrete examples, we refer to some public implementations of TSE frameworks presented in Section X.

Most TSE approaches borrow a network configuration from architectures proven effective for BSS or noise reduction. One important aspect is that an NN must be able to see enough context in the mixture to identify the target speaker. This has been achieved using such recurrent neural network (RNN)-based architectures as a stack of bidirectional long short-term memory (BLSTM) layers [10], convolutional neural network (CNN)-based architectures with a stack of convolutional layers that gradually increases the receptive field over the time axis to cover a large context [7], [23] or attention-based architectures [47].

The networks in the mixture encoder and the extraction process generally use a similar architecture. The best performance was reported when using a shallow mixture encoder (typically a single layer/block) and a much deeper extraction network, i.e., where a fusion layer is placed on the lower part of the extraction module. Furthermore, we found in our experiments that the multiplication or FiLM layers usually perform well.

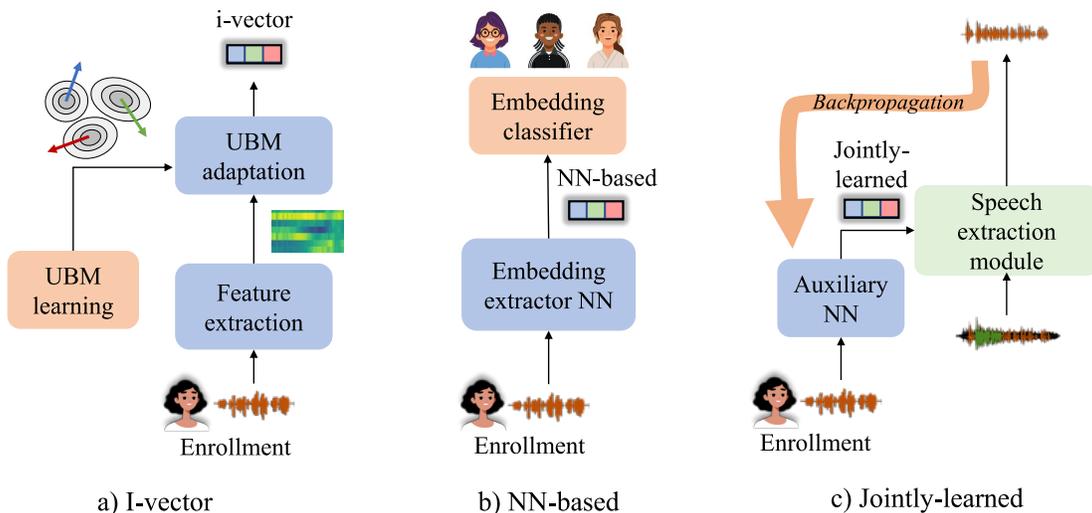


Fig. 6. Illustration of i-vector, NN-based vector, and jointly-trained embeddings: Orange parts are included only in training stage.

However, the impact of the choice of the fusion layer seems rather insignificant.

For the feature extraction, early studies used spectral features computed with STFT [7], [8], [10]. However, most recent approaches employ a learned feature extraction module following its success for separation [23], [39]. This approach allows direct optimization of the features for the given task. However, the choice of input features may depend on the acoustic conditions, and some have reported superior performance using STFT under challenging reverberant conditions [48] or using handcrafted filterbanks [49].

Except for such general considerations, it is difficult to make solid arguments for a specific network configuration since performance may depend on many factors, such as the task, the type of clue, the training data generation, and the network and training hyper-parameters.

V. AUDIO-BASED TSE

In this section, we explain how the general framework introduced in Section IV can be applied in the case of audio clues. In particular, we discuss different options to implement the clue encoder, summarize the development of the audio-based TSE, and present some representative experimental results.

A. Audio clue encoder

An audio clue is an utterance spoken by the target speaker from which we derive the characteristics of their voice, allowing identification in a mixture. This enrollment utterance can be obtained by pre-recording the user of a personal device or with a part of a recording where a wake-up keyword was uttered. The clue encoder is usually used to extract a single vector that summarizes the entire enrollment utterance.

Since the clue encoder’s goal is to extract information that defines the voice characteristics of the target speaker, embeddings from the speaker verification field are often used, such as i-vectors or NN-based embeddings (e.g., d-vectors or

x-vectors). Clue encoders trained directly for TSE tasks are also used. Fig. 6 describes these three options.

1) *I-vectors*: From their introduction around 2010, i-vectors [50] were the ruling speaker verification paradigm until the rise of NN speaker embeddings. The main idea behind i-vectors is modeling the features of an utterance using a Gaussian mixture model (GMM), whose means are constrained to a subspace and depend on the speaker and the channel effects. The subspace is defined by the Universal Background model (UBM), i.e., GMM trained on a large amount of data from many speakers, and a total variability subspace matrix. The super-vector of the means of utterance GMM μ is decomposed:

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

where \mathbf{m} is a super-vector of the means of the UBM, \mathbf{T} is a low-rank rectangular matrix representing the bases spanning the subspace, and \mathbf{w} is a random variable with standard normal prior distribution. Since an i-vector is the maximum a posteriori estimate of \mathbf{w} , it thus consists of values that enable the adaptation of the parameters of the generic UBM speaker model (\mathbf{m}) to a specific recording. As a result, it captures the speaker’s voice characteristics in the recording.

An important characteristic of i-vectors is that they capture both the speaker and channel variability. This case may be desired in some TSE applications, where we obtain enrollment utterances in identical conditions as the mixed speech. In such a situation, the channel information might also help distinguish the speakers. I-vectors have also been used in several TSE works [10].

2) *Neural network-based embeddings*: The state-of-the-art speaker verification systems predominantly use NN-based speaker embeddings, which were adopted later for TSE. The common idea is to train an NN for the task of speaker classification. Such an NN contains a “pooling layer” which converts a sequence of features into one vector. The pooling layer computes the mean and optionally the standard deviation of the sequence of features over the time dimension. The pooled vector is then classified into speaker classes or used

in other loss functions that encourage speaker discrimination. For TSE, the speaker embedding is then the vector of the activation coefficients of one of the last network layers. The most common of such NN-based speaker embeddings are d-vectors and x-vectors [51]. Many TSE works employ d-vectors [11].

Since NNs are trained for speaker classification or a related task, embeddings are usually highly speaker-discriminative. Most other sources of variability are discarded, such as the channel or content. Another advantage of this class of embeddings is that they are usually trained on large corpora with many speakers, noises, and other variations, resulting in very robust embedding extractors. Trained models are often publicly available, and the embeddings can be readily used for TSE tasks.

3) *Jointly-learned embeddings*: NN-based embeddings, such as x-vectors, are designed and trained for the task of speaker classification. Although this causes them to contain speaker information, it is questionable whether the same representation is optimal for TSE tasks. An alternative is to train the neural embedding extractor jointly with a speech extraction module. The resulting embeddings are thus directly optimized for TSE tasks. This approach has been used for TSE in several works [10], [31].

The NN performing the speaker embedding extraction takes an enrollment utterance $C_s^{(a)}$ as input and generally contains a pooling layer converting the frame-level features into one vector, similar to the embedding extractors discussed above. This NN is trained with the main NN using a common objective function. A second objective function can also be used on the embeddings to improve their speaker discriminability [46].

As mentioned above, the advantage of such embeddings is that they are trained directly for TSE and thus collect essential information for this task. On the other hand, the pre-trained embedding extractors are often trained on larger corpora and may be more robust. A possible middle ground might take a pre-trained embedding extractor and fine-tune it jointly with the TSE task. However, this has, to the best of our knowledge, not been done yet.

B. Existing approaches

The first neural TSE methods were developed around 2017. One of the first published works, SpeakerBeam [10], explored both the single-channel approach, where the target extractor was implemented by time-frequency masking, and the multi-channel approach using beamforming. This work also compared different variants of fusion layers and clue encoders. This was followed by VoiceFilter [11], which put more emphasis on ASR applications using TSE as a front-end and also investigated streaming variants with minimal latency. A slightly modified variant of the task was presented in works on speaker inventory [40], where not one but multiple speakers can be enrolled. Such a setting might be suitable for meeting scenarios. Recently, many works, such as SpEx [31], have started to use time-domain approaches, following their success in BSS [39].

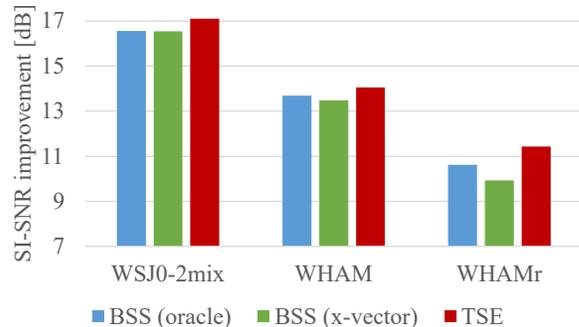


Fig. 7. Comparison of TSE and cascade BSS systems when using an audio clue in terms of SI-SNR improvement (higher is better) [52].

C. Experiments

An audio clue is a simple way to condition the system for extracting the target speaker. Many works have shown that the speaker information extracted from audio clues is sufficient for satisfactory performance. Demonstrations of many works are available online¹⁵. We present here some results to demonstrate the potential of audio clue-based approaches. The experiments were done with time-domain SpeakerBeam¹⁶, which uses a convolutional architecture, a multiplicative fusion layer, and a jointly-learned clue encoder.

The experiments were done on three different datasets (WSJ0-2mix, WHAM!, and WHAMR!) to show the performance in different conditions (clean, noisy, and reverberant, respectively). We describe these datasets in more detail in Section X. All the experiments were evaluated with the SI-SNR metric and measured the improvements over the SI-SNR of the observed mixture. More details about the experiments can be found in [52].

Figure 7 compares the TSE results with a cascade system, first doing BSS and then independent speaker identification. Speaker identification is done either in an oracle way (selecting the output closest to the reference) or with x-vectors (extracting the x-vectors from all the outputs and the enrollment utterances and selecting the output with the smallest cosine distance as the target). The BSS system uses the same convolutional architecture as TSE, differing only in that it does not have a clue encoder and the output layer is twice larger as it outputs two separated speech signals. The direct TSE scheme outperformed the cascade system, especially in more difficult conditions such as WHAMR!. This difference reflects a couple of causes: 1) the TSE model is directly optimized for the TSE task and does not spend any capacity on extracting other speakers or 2) the TSE model has additional speaker information.

Figure 8 shows an example of spectrograms obtained using TSE on a recording of two speakers from the WHAMR! database, including noise and reverberation. TSE correctly

¹⁵Demonstrations of audio clues approaches: VoiceFilter [11] <https://google.github.io/speaker-id/publications/VoiceFilter/>, SpeakerBeam [10] <https://www.youtube.com/watch?v=7FSHgKip6vI>.

¹⁶<https://github.com/butspeechfit/speakerbeam>

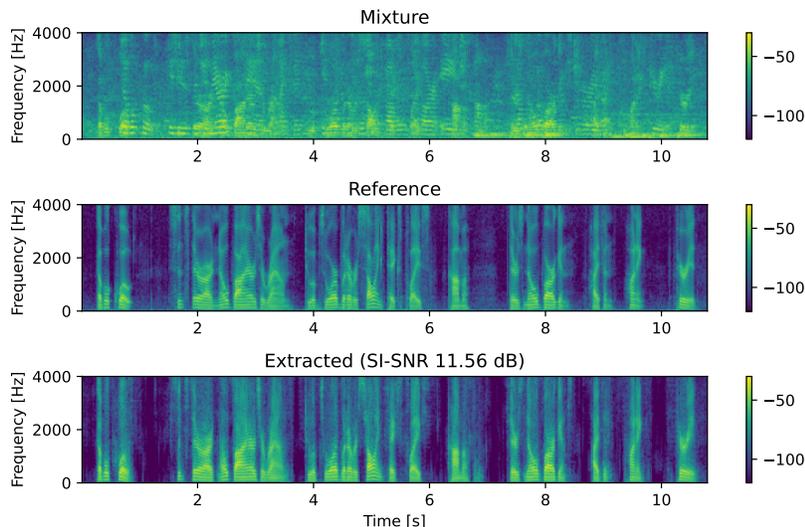


Fig. 8. Example of spectrograms of mixed, reference, and extracted speech: Example is taken from WHAMR! database.

identifies the target speaker and removes all the interference, including the second speaker, noise, and reverberation.

D. Limitations and outlook

Using TSE systems conditioned on audio clues is particularly practical due to the simplicity of obtaining the clues, i.e., no additional hardware is needed, such as cameras or multiple microphones. Considering the good performance demonstrated in the literature, these systems are widely applicable. Nowadays, the methods are rapidly evolving and achieving increasingly higher accuracy.

The main challenge in audio-clue-based systems is correct identification of the target speaker. The speech signal of the same speaker might have highly different characteristics in different conditions due to such factors as emotional state, channel effects, or the Lombard effect. TSE systems must be robust enough to such intra-speaker variability. On the other hand, different speakers might have very similar voices, leading to erroneous identification if the TSE system lacks sufficient accuracy.

Resolving both issues requires precise speaker modeling. In this regard, the TSE methods may draw inspiration from the latest advances in the speaker verification field, including advanced model architectures, realistic datasets with a huge number of speakers for training, or using pre-trained features from self-supervised models.

VI. VISUAL/MULTI-MODAL CLUE-BASED TSE

Visual clue-based TSE assumes that a video camera captures the face of the target speaker who is talking in the mixture [7], [8]. Using visual clues is motivated by psycho-acoustic studies (see the references in a previous work [6]) that revealed that humans look at lip movements to understand speech better. Similarly, the visual clues of TSE systems derive hints about the state of the target speech from the lip movements, such

as whether the target speaker is speaking or silent or more refined information about the phoneme being uttered.

A visual clue, which presents different characteristics than audio clues because it captures information from another modality, is time-synchronized with the target speech in the mixture without being corrupted by the interference speakers. Therefore, a visual clue-based TSE can better handle mixtures of speakers with similar voices, such as same-gender mixtures, than audio clue-based systems because the extraction process is not based on the speaker’s voice characteristics¹⁷. Another potential advantage is that the users may not need to pre-enroll their voice. Video signals are also readily available for many applications such as video-conferencing.

Figure 9 shows a diagram of a visual TSE system that follows the same structure as the general TSE framework introduced in Section IV. Only the visual clue encoder part is specific to the task. We describe it in more detail below and then introduce a multi-modal clue extension. We conclude this section with some experimental results and discussions.

A. Visual clue encoder

The visual clue encoder computes from the video signal a representation that allows the speech extraction module to identify and extract the target speech in the mixture. This processing involves the steps described below:

$$\mathbf{E}_s^{(v)} = \text{Upsample}(\text{NN}(\text{VFE}(\mathbf{C}_s^{(v)}), \theta^{v\text{-clue}})), \quad (18)$$

where $\mathbf{E}_s^{(v)} \in \mathbb{R}^{D^{\text{Emb}} \times N}$ represents the sequence of the visual embedding vectors, $\mathbf{C}_s^{(v)}$ is the video signal obtained after pre-processing, $\text{VFE}(\cdot)$ is the visual feature extraction module, $\text{NN}(\cdot, \theta^{v\text{-clue}})$ is an NN with parameters $\theta^{v\text{-clue}}$, and $\text{Upsample}(\cdot)$ represents the up-sampling operation. The latter up-sampling step is required because the sampling rates of the audio and video devices are usually different. Up-sampling

¹⁷Some works can even perform extraction from a mixture of the same speaker’s speech [8].

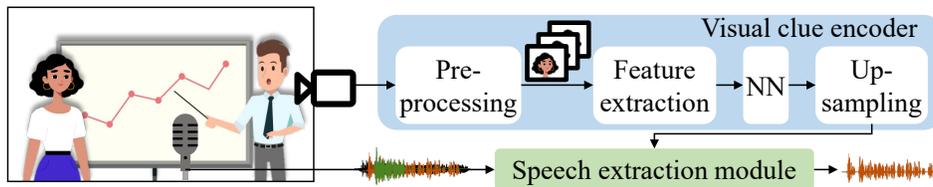


Fig. 9. Visual clue-based TSE system.

matches the number of frames of the mixture and visual clue encoders.

1) *Pre-processing*: First, the video signal captured by the camera requires pre-processing to isolate the face of the target speaker. Depending on the application, this may require detecting and tracking the target speaker’s face and cropping the video. These pre-processing steps can be performed using previously well-established video processing algorithms [6].

2) *Visual feature extraction*: Similar to an audio-clue-based TSE, the visual clue encoder can directly extract embeddings from raw video data or visual features. With the first option, the raw video is processed with a CNN whose parameters are jointly-learned with the speech extraction module to enable direct optimization of the features for the extraction task without any loss of information. However, since the video signals are high-dimensional data, achieving joint optimization can be complex. This approach has been used successfully with speaker-close conditions [53]. Extending it to speaker-open conditions might require a considerable amount of data or careful design of the training loss using, e.g., multi-task training to help the visual encoder capture relevant information.

Most visual TSE works use instead a visual feature extractor pre-trained on another task to reduce the dimensionality of the data. Such feature extractors can leverage a large amount of image or video data (that do not need to be speech mixtures) to learn representation robust to variations, such as resolution, luminosity, or head orientation. The first option is to use facial landmark points as features. Facial landmarks are the key points on a face that indicate the mouth, eyes, or nose positions and offer a very low-dimension representation of a face, which is interpretable. Moreover, face landmarks can be easily computed with efficient off-the-shelf algorithms [32].

The other option is to use neural embeddings derived from an image/video processing NN trained on a different task, which was proposed in three concurrent works [7]–[9]. Ephrat et al. [8] used visual embeddings obtained from an intermediate layer of a face recognition system called FaceNet. This face recognition system is trained so that embeddings derived from photographs of the same person are close and embeddings from different persons are far from each other. It thus requires only a corpus of still images with person identity labels for training the system. However, the embeddings do not capture the lip movement dynamics and are not explicitly related to the acoustic content.

Alternatively, Afouras et al. [7] proposed using embeddings obtained from a network trained to perform lip-reading, i.e., where a network is trained to estimate the phoneme or word

uttered from the video of the speaker’s lips. The resulting embeddings are thus directly related to the acoustic content. However, the training requires video with the associated phoneme or word transcriptions, which are more demanding and costly to obtain.

The third option introduced by Owens et al. [9] exploits embeddings derived from an NN trained to predict whether the audio and visual tracks of a video are synchronized. This approach enables self-supervised training, where the training data are simply created by randomly shifting the audio track by a few seconds. The embeddings capture information on the association between the lip motions and the timing of the sounds in the audio. All three options [7]–[9] can successfully perform a visual TSE.

3) *Transformation and up-sampling*: Except with joint-training approaches, the visual features are (pre-)trained on different tasks and thus do not provide a representation optimal for TSE. Besides, since some of the visual features are extracted from the individual frames of a video, the dynamics of lip movements are not captured. Therefore, the visual features are further transformed with an NN, which is jointly trained with the speech extraction module. The NN, which allows learning a representation optimal for TSE, can be implemented with long short-term memory (LSTM) or convolutional layers across the time dimension to model the time series of the visual features, enabling the lip movement dynamics to be captured. Finally, the visual embeddings are up-sampled to match the sampling rate of audio features \mathbf{Z}_y .

B. Audio-visual clue-based TSE

Audio and visual clue-based TSE systems have complementary properties. An audio clue-based TSE is not affected by speaker movements and visual occlusions. In contrast, a visual clue-based TSE is less affected by the voice characteristics of the speakers in the mixture. By combining these approaches, we can build TSE systems that exploit the strengths of both clues for improving the robustness to various conditions [33], [36].

Figure 10 shows a diagram of an audio-visual TSE system, which assumes access to the pre-recorded enrollment of the target speaker to provide an audio clue and a video camera for a visual clue. The system uses the audio and visual clue encoders described in Sections V-A and VI-A and combines these clues into an audio-visual embedding, which is given to the speech extraction module. Audio-visual embeddings can be simply the concatenation [35] or the summation of the audio and visual embeddings, or obtained as a weighted sum [33],

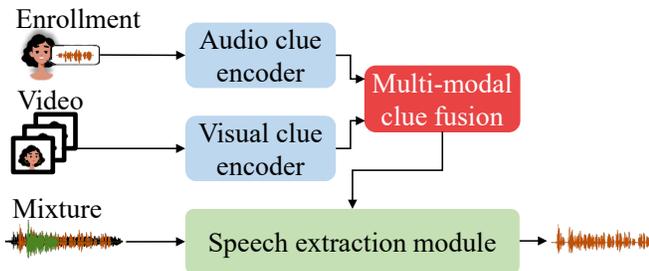


Fig. 10. Audio-visual clue-based TSE system

[34], where the weights can vary depending on the reliability of each clue. The weighted sum approach can be implemented with an attention layer widely used in machine learning, which enables dynamic weighting of the contribution of each clue.

C. Experimental results and discussion

Several visual TSE systems have been proposed, which differ mostly by the type of visual features used and the network configuration. These systems have demonstrated astonishing results, which can be attested by the demonstrations available online¹⁸. Here we briefly describe experiments using the audio, visual, and audio-visual time-domain SpeakerBeam systems [34], which use a similar configuration as the system in Section V-C. The speech extraction module employs a stack of time-convolutional blocks and a multiplicative fusion layer. The audio clue encoder consists of the jointly-learned embeddings described in Section V-A3. The visual clue encoder uses visual features derived from face recognition like a previous work [8]. The audio-visual system combines the visual and audio clues with an attention layer [34].

The experiments used mixtures of utterances from the LRS3-TED corpus¹⁹, which consists of single speaker utterances with associated videos. We analyzed the behavior under various conditions by looking at results from same and different gender mixtures and two examples of clue corruptions (enrollment corrupted with white noise at SNR of 0 dB and video with a mask on the speaker’s mouth). The details of the experimental setup are available in [34].

Figure 11 compares the extraction performance measured in terms of SDR improvement for audio, visual, and audio-visual TSE under various mixture and clue conditions. We confirmed that a visual clue-based TSE is less sensitive to the characteristics of the speakers in the mixture since the performance gap between different- and same-gender mixtures is smaller than with an audio clue-based TSE. When using a single clue, performance can be degraded when this clue is corrupted. However, the audio-visual system that exploits both clues can achieve superior extraction performance and is more robust to clue corruption.

¹⁸Demo samples for several approaches are available, e.g., for [9]: <https://andrewowens.com/multisensory>, for [8]: <https://looking-to-listen.github.io>, for [7]: <https://www.robots.ox.ac.uk/~vgg/demo/theconversation>, and for [34]: http://www.kecl.ntt.co.jp/icl/signal/member/demo/audio_visual_speakerbeam.html

¹⁹https://www.robots.ox.ac.uk/~vgg/data/lip_reading/lrs3.html

D. Discussions and outlook

Visual clue-based TSE approaches offer an alternative to audio-clue-based ones when a camera is available. The idea of using visual clues for TSE is not new [4], [5], although recent neural systems have achieved an impressive level of performance. This is probably because NNs can effectively model the relationship between the different modalities learned from a large amount of training data.

Issues and research opportunities remain with the current visual clue-based TSE systems. First, most approaches do not consider the speaker tracking problem and assume that the audio and video signals are synchronized. These aspects must be considered when designing and evaluating future TSE systems. Second, video processing involves high computational costs, and more research is needed to develop efficient online systems.

VII. SPATIAL CLUE-BASED TSE

When using a microphone array to record a signal, spatial information can be used to discriminate among sources. In particular, access to multi-channel recordings opens the way to extract target speakers based on their location, i.e., using spatial clues (as indicated in Fig. 1). This section explains how such spatial clues can be obtained and used in TSE systems. While enhancing speakers from a given direction has a long research history [2], we focus here on neural methods that follow the scope of our overview paper.

Note that multi-channel signals can also be utilized in the extraction process using beamforming. Such an extraction process can be used in systems with any type of clue, only requiring that the mixed speech be recorded with multiple microphones. This beamforming process was reviewed in Section IV-C. In this section, we focus specifically on the processing of spatial clues.

A. Obtaining spatial clues

In some situations, the target speaker’s location is approximately known in advance. For example, for an in-car ASR, the driver’s position is limited to a certain region in a car. In other scenarios, we might have access to a multi-channel enrollment utterance of the speaker recorded in the same position as the final mixed speech. In such a case, audio source localization methods can be applied. Conventionally, this can be done by methods based on generalized cross-correlation or steered-response power, but recently, deep learning methods have also shown success in this task. An alternative is to skip the explicit estimation of the location and directly extract features in which the location is encoded when a multi-channel enrollment is available. We will detail this approach further in the next section.

Spatial clues can also be obtained from a video using face detection and tracking systems. A previous work [36] demonstrated this possibility with a 180-degree wide-angle camera positioned parallel to a linear microphone array²⁰. By identifying the target speaker in the video, the azimuth with

²⁰<https://yongxuustc.github.io/grnbf>

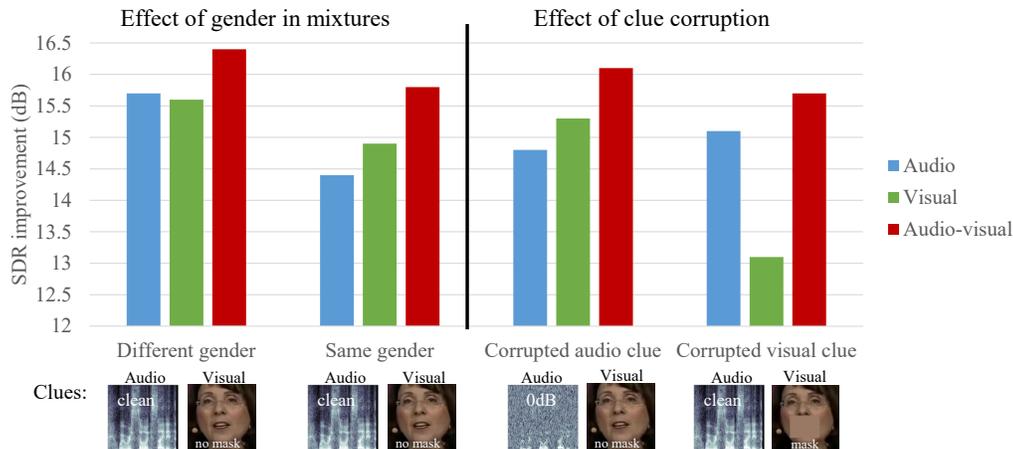


Fig. 11. SDR Improvement of TSE with audio, visual, and audio-visual clues for mixtures of same/different gender and for corruptions of audio and visual clues: Audio clues were corrupted by adding white noise at SNR of 0 dB to enrollment utterance. Video clues were corrupted by masking mouth region in video.

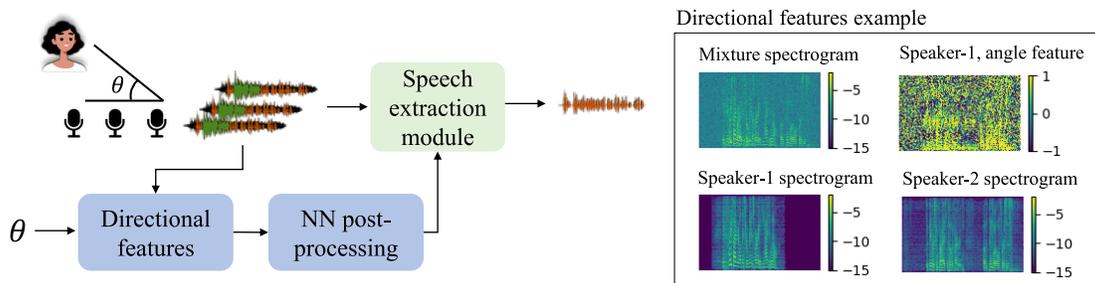


Fig. 12. Illustration of usage of spatial clue encoder and directional features

respect to the microphone array was roughly approximated. Depth cameras can also be used to estimate not only the azimuth but also the elevation and distance of the speaker.

B. Spatial clue encoder

The left part of Fig. 12 shows the overall structure and the usage of a spatial clue encoder, which usually consists of two parts: the extraction of directional features and an NN post-processing of them. Two possible forms of spatial clues are dominant in the literature: the angle of the target speaker with respect to the microphone array or a multi-channel enrollment utterance recorded in the target location. Both can be encoded into directional features.

When the spatial clue is DOA, the most commonly used directional features are the *angle features*, which are computed as the cosine of the difference between the IPD and the target phase difference (TPD):

$$AF[n, f] = \sum_{m_1, m_2 \in \mathcal{M}} \cos \left(\text{TPD}(m_1, m_2, \phi_s, f) - \text{IPD}(m_1, m_2, n, f) \right) \quad (19)$$

$$\text{TPD}(m_1, m_2, \phi_s, f) = \frac{2\pi f F_s}{F} \frac{\cos \phi_s \Delta_{m_1, m_2}}{c} \quad (20)$$

$$\text{IPD}(m_1, m_2, n, f) = \angle Y^{m_2}[n, f] - \angle Y^{m_1}[n, f], \quad (21)$$

where \mathcal{M} is a set of pairs of microphones used to compute the feature, F_s is the sampling frequency, ϕ_s is the target direction, c is the sound's velocity, and Δ_{m_1, m_2} is the distance from microphone m_1 to microphone m_2 . An example of angle features is shown on the right of Fig. 12. For time-frequency bins dominated by the source from direction ϕ_s , the value of the angle feature should be close to 1 or -1. Other directional features have been proposed that exploit a grid of fixed beamformers. A directional power ratio measures the ratio between the power of the response of a beamformer steered into the target direction and the power of the beamformer responses steered into all the directions in the grid. In a similar fashion, a directional signal-to-noise ratio can also be computed, which compares the response of a beamformer in the target direction with the response of a beamformer in the direction with the strongest interference.

If the spatial clue consists of a multi-channel enrollment utterance, the directional feature can be formed as a vector of IPDs computed from the enrollment. Alternatively, the DOA can be estimated from the enrollment, and the spatial features derived from it can be used.

Note that when using a spatial clue to determine the target speaker, the multi-channel input of the speech extraction module must also be used. This enables the identification of the speaker coming from the target location in the mixture. Furthermore, a target extractor is often implemented as beamforming, as explained in Section IV-C.

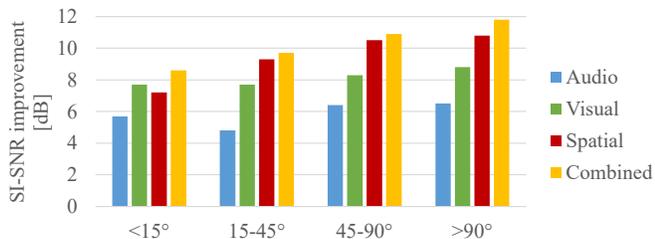


Fig. 13. SI-SNR improvement of TSE with audio, visual, and spatial clues in four conditions based on angle separation between speakers [36]

C. Combination with other clues

Although a spatial clue is very informative and generally can lead the TSE system to a correct extraction of the target, it does fail in some instances. Estimation errors of DOA are harmful to proper extraction. Furthermore, if the spatial separation of the speakers with respect to the microphone array is not significant enough, the spatial clue may not discriminate between them. Combining a spatial clue with audio or visual clues is an option to combat such failure cases.

D. Experimental results

We next report the results from an experiment with spatial clues [36] that compared the effectiveness of using audio, visual, and spatial clues. The audio-clue encoder was trained jointly with the extraction module, and the visual encoder was a pre-trained lip-reading network. The target speaker’s direction was encoded in the angle feature. The spatial and visual embeddings were fused with the extraction network by concatenation and the audio embedding with a factorized layer. The extraction module employed a neural network consisting of temporal convolutional layers.

The experiments were performed on a Mandarin audio-visual dataset containing mixtures of two and three speakers. The results in Fig. 13 were divided into several conditions, based on the angle separation between the closest speakers. The spatial clue is very effective, although the performance declines when speakers are near each other ($< 15^\circ$). A combination with other modalities outperformed any individual type of clue in all the conditions.

E. Discussion

Using spatial clues is a powerful way of conditioning a TSE system to extract the target speaker. It relies on the availability of signals from a microphone array and a way to determine the location of the target speaker. Unfortunately, these restrictions limit the applications to some extent. Neural TSE methods with spatial clues follow a long history of research on the topic, such as beamforming techniques, and extend them with non-linear processing. This approach unifies the methods with those using other clues and allows a straightforward combination of different clues into one system. Such combinations can alleviate the shortcomings of spatial clues, including the failures when the speakers are located in the same direction from the point of view of the microphones.

In most current neural TSE works, the target speaker’s location is assumed to be fixed. Although the methods should be easily extended to a dynamic case, investigations of such settings remain relatively rare [24].

VIII. EXTENSION TO OTHER TASKS

The ideas of TSE can be applied to other speech processing tasks, such as ASR and diarization.

A. Target-speaker ASR

An important application of TSE is TS-ASR, where the goal is to transcribe the target speaker’s speech and ignore all the interference speakers. The TSE approaches we described can be naturally used as a front-end to an ASR system to achieve TS-ASR. Such a cascade combination allows for a modular system, which offers ease of development and interpretability. However, the TSE system is often optimized with a signal loss, as in Eq. (16). Such a TSE system inevitably introduces artifacts caused by the remaining interferences, over-suppression, and other non-linear processing distortions. These artifacts limit the expected performance improvement from a TSE front-end.

One approach to mitigate the effect of such artifacts is to optimize the TSE front-end with an ASR criterion [10]. The TSE front-end and the ASR back-end are NNs and can be interconnected with differentiable operations, such as beamforming and feature extraction. Therefore, a cascade system can be represented with a single computational graph, allowing all parameters to be jointly trained. Such joint-training can significantly improve the TS-ASR performance.

Another approach inserts a fusion layer into an ASR system [26], [45] to directly perform clue conditioning. These integrated TS-ASR systems avoid any explicit signal extraction step, a decision that reduces the computational cost, although such systems may be less interpretable than cascade systems.

TS-ASR can use the audio clues provided by pre-recorded enrollment utterances [10], [26], [45] or from a keyword (anchor) for a smart-device scenario [54], for example. Some works have also exploited visual clues, which can be used for the extraction process and to implement an audio-visual ASR back-end, since lip-reading also improves ASR performance [55].

B. Target-speaker VAD and diarization

The problem of speech diarization consists of detecting who spoke when in a multi-speaker recording. This technology is essential for achieving, e.g., meeting recognition and analysis systems that can transcribe a discussion between multiple participants. Several works have explored using speaker clues to perform this task [27], [28].

For example, a personalized VAD [27] exploits a speaker embedding vector derived from an enrollment utterance of the target speaker to predict its activity, i.e., whether they are speaking at a given time. In principle, this can be done with a system like that presented in Section IV, where the output layer performs the binary classification of the speaker activity

instead of estimating the target speech signal. Similar systems have also been proposed using visual clues, called audio-visual VAD [56]. Predicting the target speaker’s activity is arguably a more straightforward task than estimating its speech signal. Consequently, TS-VAD can use simpler network architectures, leading to more lightweight processing.

The above TS-VAD systems, which estimate the speech activity of a single target speaker, have been extended to simultaneously output the activity of multiple target speakers [28]. The resulting system achieved the top diarization performance in the CHiME 6 evaluation campaign²¹.

IX. REMAINING ISSUES AND OUTLOOK

Research toward computational selective hearing has been a long endeavor. Recent developments in TSE have enabled identifying and extracting a target speaker’s voice in a mixture by exploiting audio, visual, or spatial clues, which is one step closer to solving the cocktail-party problem. Progress in speech processing (speech enhancement, speaker recognition) and image processing (face recognition, lip-reading), combined with deep learning technologies to learn models that can effectively condition processing on auxiliary clues, triggered the progress in the TSE field. Some of the works we presented have achieved levels of performance that seemed out-of-reach just a few years ago and are already being deployed in products²².

Despite substantial achievements, many opportunities remain for further research, some of which we list below.

A. Deployment of TSE systems

Most of the systems we described operate offline and are computationally expensive. They are also evaluated under controlled (mostly simulated mixture) settings. Deploying such systems introduces engineering and research challenges to reduce computational costs while maintaining high performance under less controlled recording conditions. We next discuss some of these aspects.

1) *Inactive target speaker*: Most TSE systems have been evaluated assuming that the target speaker is actively speaking in the mixture. In practice, we may not know beforehand whether the target speaker will be active. We expect that a TSE system can output no signal when the target speaker is inactive, which may not actually be the case with most current systems that are not explicitly trained to do so. The inactive target speaker problem is specific to TSE. The type of clue used may also greatly impact the difficulty of tackling this problem. For instance, visual voice activity detection [5] might alleviate this issue. However, it is more challenging with audio clues [57], and further research may be required.

²¹The results of the CHiME 6 challenge can be found at: <https://chimechallenge.github.io/chime6/results.html>. The top system used TS-VAD among other technologies. DiHARD 3 performed a diarization evaluation on the CHiME 6 challenge data. Here the top system also used TS-VAD: <https://dihardchallenge.github.io/dihard3/results>

²²The following blog details the effort for deploying a visual clue-based TSE system for on-device processing: <https://ai.googleblog.com/2020/10/audiovisual-speech-enhancement-in.html>.

2) *Training and evaluation criteria*: Most TSE systems are trained and evaluated using such signal level metrics as SNR or SDR. Although these metrics are indicative of the extraction performance, their use presents two issues.

First, they may not always be correlated with human perception and intelligibility or with ASR performance. This issue is not specific to TSE; it is common to BSS and noise reduction methods. For ASR we can train a system end-to-end, as discussed in Section VIII-A. When targeting applications for human listeners, the problem can be partly addressed using other metrics for training or evaluation that correlate better with human perception, such as short-time objective intelligibility (STOI) or perceptual evaluation of speech quality (PESQ) [6]. However, controlled listening tests must be conducted to confirm the impact of a TSE on human listeners [6].

Second, unlike BSS and noise reduction, a TSE system needs to identify the target speech, implying other sources of errors. Indeed, failing to identify the target may lead to incorrectly estimating an interference speaker or inaccurately outputting the mixture. Although these errors directly impact the SDR scores, it would be fruitful to agree on the evaluation metrics that separate extraction and identification performance to better reveal the behavior of TSE systems. Signal level metrics might not satisfactorily represent the extraction performance for inactive speaker cases. A better understanding of the failures might help develop TSE systems that can recognize when they cannot identify the target speech, which is appealing for practical applications.

Consequently, developing better training and evaluation criteria are critical research directions.

3) *Robustness to recording conditions*: Training neural TSE systems requires simulated mixtures, as discussed in Section IV-D. Applying these systems to real conditions (multi-speaker mixtures recorded directly with a microphone) requires that the training data match the application scenario relatively well. For example, the type of noise and reverberation may vary significantly depending on where a system is deployed. This raises questions about the robustness of TSE systems to various recording conditions.

Neural TSE systems trained with a large amount of simulated data have been shown to generalize to real recording conditions [8]. However, exploiting real recordings where no reference target speech signal is available could further improve performance. Real recordings might augment the training data or be used to adapt a TSE system to a new environment. The issue is defining unsupervised training losses correlated with the extraction performance of the target speech without requiring access to the reference target signal.

Another interesting research direction is combining neural TSE systems, which are powerful under matched conditions, with such generative-based approaches as IVE [12], which are adaptive to recording conditions.

4) *Lightweight and low-latency systems*: Research on lightweight and low-latency TSE systems is gaining momentum as the use of teleconferencing systems in noisy environments has risen in response to the Covid pandemic. Other important use cases for TSE are hearing aids and

hearables, both of which impose very severe constraints in terms of computation costs and latency. The recent DNS²³ and Clarity²⁴ challenges that target teleconferencing and hearing aid application scenarios include tracks where target speaker clues (enrollment data) can be exploited. This demonstrates the growing interest in practical solutions for TSE.

Since TSE is related to BSS and noise reduction, the development of online and low-latency TSE systems can be inspired from the progress of BSS/noise reduction in that direction. However, TSE must also identify the target speech, which may need specific solutions that exploit the long-context of the mixture to reliably and efficiently capture a speaker's identity.

5) *Spatial rendering*: For applications of TSE to hearing aids or hearables, sounds must be localized in space after the TSE processing. Therefore, a TSE system must not only extract the target speech but also estimate its direction to allow rendering it so that a listener feels the correct direction of the source.

B. Self-supervised and cross-modal learning

A TSE system identifies the target speech in a mixture based on the intermediate representation of the mixture and the clue. Naturally, TSE benefits from better intermediate representations. For example, speech models learned with self-supervised learning criteria have gained attention as a way to obtain robust speech representations. They have shown potential for pre-training many speech processing downstream tasks, such as ASR, speaker identification, and BSS. Such self-supervised models could also reveal advantages for TSE since they could improve robustness by allowing efficient pre-training on various acoustic conditions. Moreover, for audio-based TSE, using the same self-supervised pre-trained model for the audio clue encoder and the speech extraction module will help to learn the common embedding space between the enrollment and speech signals in the mixture. Similarly, the progress in cross-modal learning, which aims to learn the joint representation of data across modalities, could benefit such multi-modal approaches as visual clue-based TSE.

C. Exploring other clues

We presented three types of clues that have been widely used for TSE. However, other clues can also be considered. For example, recent works have explored other types of spatial clues such as the distance [58]. Moreover, humans do not only rely on physical clues to perform selective hearing. We also use more abstract clues, such as semantic ones. Indeed, we can rapidly focus our attention on a speaker when we hear our name or a topic we are interested in. Reproducing a similar mechanism would require TSE systems that operate with semantic clues, which introduces novel challenges concerning how to represent semantic information and exploit it within a TSE system. Some works have started to explore this direction,

such as conditioning on languages [59] or more abstract concepts [60].

Other interesting clues consist of signals that measure a listener's brain activity to guide the extraction process. Indeed, the electroencephalogram (EEG) signal of a listener focusing on a speaker correlates with the envelope of that speaker's speech signal. Ceolini et al. identified the possibility of using EEG as clues for TSE with a system similar to the one described in Section IV [61]. An EEG-guided TSE might open the door for futuristic hearing aids controlled by the user's brain activity, which might automatically emphasize the speaker a user wants to hear. However, research is still needed because developing a system that requires marginal tuning to the listener is especially challenging. Moreover, collecting a large amount of training data is very complicated since it is more difficult to control the quality of such clues. Compared to audio and visual TSE clues, EEG signals are very noisy and affected by changes in the attention of the listener, body movements, and other factors.

D. Beyond speech

Human selective listening abilities go beyond speech signals. For example, we can focus on listening to the part of an instrument in an orchestra or switch our attention to a siren or a barking dog. In this paper, we focused on TSE, but similar extraction problems have also been explored for other audio-processing tasks. For example, much research has been performed on extracting the track of an instrument in a piece of music conditioned on, e.g., the type of instrument [62], video of the musician playing [63], or EEG signal of the listener [64]. These approaches may be important to realize, e.g., audio-visual music analysis [65].

Recently, the problem was extended to the extraction of arbitrary sounds from a mixture [66], [67], e.g., extracting the sound of a siren or a klaxon from a recording of a mixture of street sounds. We can use such systems as that introduced in Section IV to tackle these problems, where the clue can be a class label indicating the type of target sound [66], the enrollment audio of a similar target sound [67], a video of the sound source [9] or a text description of the target sound [68]. Target sound extraction may become an important technology to design, e.g., hearables or hearing aids that could filter out nuisances and emphasize important sounds in our surroundings, or audio visual scene analysis [9].

Psycho-acoustic studies suggest that humans process speech and music partly using shared auditory mechanisms and that exposure to music can lead to better discrimination of speech sounds [69]. It would be interesting to explore whether, similarly to humans, TSE systems could benefit from exposure to other acoustic signals by training a system to extract target speech, music, or arbitrary sounds.

X. RESOURCES

We conclude by providing pointers to selected datasets and toolkits available for those motivated to experiment with TSE.

TSE works mostly use datasets designed for BSS. These datasets consist generally of artificial mixtures generated from

²³<https://www.microsoft.com/en-us/research/academic-program/deep-noise-suppression-challenge-icassp-2022/>

²⁴https://claritychallenge.github.io/clarity_CC_doc/

TABLE III
SOME DATASETS AND TOOLKITS

	Name	Description	Link
Dataset	WSJ0-mix	Mixtures of two or three speakers	www.merl.com/demos/deep-clustering
	WHAM(R)!	Noisy and reverberant versions of WSJ0-mix	wham.whisper.ai
	Librimix	Larger dataset of mixtures of two or three speakers	github.com/JorisCos/LibriMix
	LibriCSS	Meeting-like mixtures recorded in a room	github.com/chenzhuo1011/libri_css
	MC-WSJ0-mix	Spatialized version of WSJ0-2mix	www.merl.com/demos/deep-clustering
	SMS-WSJ	Multi-channel corpus based on WSJ	github.com/fgnt/sms_wsj
	LRS	Audio-visual corpus from TED or BBC videos	www.robots.ox.ac.uk/~vgg/data/lip_reading
	AVSpeech	Very large audio-visual corpus from YouTube videos	looking-to-listen.github.io/avspeech
Tools	SpeakerBeam	Time-domain audio-based TSE system	github.com/butspeechfit/speakerbeam
	SpEx+	Time-domain audio-based TSE system [31]	github.com/xuchenglin28/speaker_extraction_SpEx
	VoiceFilter	Time-domain audio-based TSE system (Unofficial) [11]	github.com/mindsfab-ai/voicefilter
	Multisensory	Visual clue-based TSE [9]	github.com/andrewowens/multisensory
	AV Speech enh.	Face landmark-based visual clue-based TSE [32]	github.com/dr-pato/audio_visual_speech_enhancement
	FaceNet	Visual feature extractor used in [8], [33], [34]	github.com/davidsandberg/facenet

the isolated signals of the individual speakers and background. This allows evaluation of the performance by comparing the estimated signals to the original references. Additionally, TSE methods also require a clue, i.e., an enrollment utterance for the target speaker or video signal. We can obtain enrollment utterances by choosing a random utterance of the target speaker from the same database, provided that the utterance is different from the one in the mixture. For a video clue, it requires using an audio-visual dataset. The top of Table III lists some of the most commonly used datasets for audio and visual TSE.

Several implementations of TSE systems are openly available and listed in the lower part of Table III. Although there are no public implementations for some of the visual TSE systems, they can be re-implemented following the audio TSE toolkits and using openly available visual feature extractors such as FaceNet, which was used in some previous works [8], [33], [34].

XI. ACKNOWLEDGMENTS

This work was partly supported by the Czech Ministry of Education, Youth and Sports from project no. LTAIN19087 "Multi-linguality in speech technologies." Computing on IT4I supercomputer was supported by the Czech Ministry of Education, Youth and Sports from the Large Infrastructures for Research, Experimental Development and Innovations project "e-Infrastructure CZ – LM2018140". The figures contain elements designed by pikisuperstar/Freepik.

REFERENCES

- [1] A. W. Bronkhorst, "The cocktail-party problem revisited: early processing and selection of multi-talker speech," *Attention, Perception, & Psychophysics*, vol. 77, no. 5, pp. 1465–1487, 2015.
- [2] J. L. Flanagan, J. D. Johnston, R. Zahn, and G. W. Elko, "Computer-steered microphone arrays for sound transduction in large rooms," *The Journal of the Acoustical Society of America*, vol. 78, no. 5, pp. 1508–1518, 1985.
- [3] R. Gu, L. Chen, S.-X. Zhang, J. Zheng, Y. Xu, M. Yu, D. Su, Y. Zou, and D. Yu, "Neural spatial filter: Target speaker speech separation assisted with directional information," in *Interspeech*, pp. 4290–4294, 2019.
- [4] J. Hershey and M. Casey, "Audio-visual sound separation via Hidden Markov Models," *Advances in Neural Information Processing Systems*, vol. 14, 2001.
- [5] B. Rivet, W. Wang, S. M. Naqvi, and J. A. Chambers, "Audiovisual speech source separation: An overview of key methodologies," *IEEE Signal Processing Magazine*, vol. 31, no. 3, pp. 125–134, 2014.
- [6] D. Michelsanti, Z.-H. Tan, S.-X. Zhang, Y. Xu, M. Yu, D. Yu, and J. Jensen, "An overview of deep-learning-based audio-visual speech enhancement and separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 1368–1396, 2021.
- [7] T. Afouras, J. S. Chung, and A. Zisserman, "The Conversation: Deep Audio-Visual Speech Enhancement," in *Interspeech*, pp. 3244–3248, 2018.
- [8] A. Ephrat, I. Mosseri, O. Lang, T. Dekel, K. Wilson, A. Hassidim, W. T. Freeman, and M. Rubinstein, "Looking to listen at the cocktail party: A speaker-independent audio-visual model for speech separation," *ACM Trans. Graph.*, vol. 37, jul 2018.
- [9] A. Owens and A. A. Efros, "Audio-visual scene analysis with self-supervised multisensory features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 631–648, 2018.
- [10] K. Žmolíková, M. Delcroix, K. Kinoshita, T. Ochiai, T. Nakatani, L. Burget, and J. Černocký, "SpeakerBeam: Speaker aware neural network for target speaker extraction in speech mixtures," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 4, pp. 800–814, 2019.
- [11] Q. Wang, H. Muckenhirn, K. Wilson, P. Sridhar, Z. Wu, J. R. Hershey, R. A. Saurous, R. J. Weiss, Y. Jia, and I. L. Moreno, "VoiceFilter: Targeted Voice Separation by Speaker-Conditioned Spectrogram Masking," in *Interspeech*, pp. 2728–2732, 2019.
- [12] J. Janský, J. Málek, J. Čmejla, T. Kounovský, Z. Koldovský, and J. Žďánský, "Adaptive blind audio source extraction supervised by dominant speaker identification using x-vectors," in *ICASSP*, pp. 676–680, 2020.
- [13] H. Sawada, N. Ono, H. Kameoka, D. Kitamura, and H. Saruwatari, "A review of blind source separation methods: two converging routes to ILRMA originating from ICA and NMF," *APSIPA Transactions on Signal and Information Processing*, vol. 8, p. e12, 2019.
- [14] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 26, no. 10, pp. 1702–1726, 2018.
- [15] S. E. Eskimez, T. Yoshioka, H. Wang, X. Wang, Z. Chen, and X. Huang, "Personalized speech enhancement: new models and comprehensive evaluation," in *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 356–360, 2022.
- [16] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, "A consolidated perspective on multimicrophone speech enhancement and source separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 4, pp. 692–730, 2017.
- [17] J. R. Hershey, S. J. Rennie, P. A. Olsen, and T. T. Kristjansson, "Superhuman multi-talker speech recognition: A graphical modeling approach," *Computer Speech and Language*, vol. 24, no. 1, pp. 45 – 66, 2010.
- [18] T. Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 3, pp. 1066–1074, 2007.
- [19] A. Ozerov and E. Vincent, "Using the fasst source separation toolbox for

- noise robust speech recognition,” in *International Workshop on Machine Listening in Multisource Environments (CHiME 2011)*, 2011.
- [20] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *ICASSP*, pp. 31–35, 2016.
- [21] D. Yu, M. Kolbaek, Z.-H. Tan, and J. Jensen, “Permutation invariant training of deep models for speaker-independent multi-talker speech separation,” in *ICASSP*, pp. 241–245, 2017.
- [22] J. Du, Y. Tu, Y. Xu, L. Dai, and C.-H. Lee, “Speech separation of a target speaker based on deep neural networks,” in *2014 12th International Conference on Signal Processing (ICSP)*, pp. 473–477, IEEE, 2014.
- [23] C. Xu, W. Rao, E. S. Chng, and H. Li, “Time-domain speaker extraction network,” in *ASRU*, pp. 327–334, IEEE, 2019.
- [24] J. Heitkaemper, T. Fehér, M. Freitag, and R. Haeb-Umbach, “A study on online source extraction in the presence of changing speaker positions,” in *International Conference on Statistical Language and Speech Processing*, pp. 198–209, Springer, 2019.
- [25] M. Delcroix, K. Zmolikova, K. Kinoshita, A. Ogawa, and T. Nakatani, “Single channel target speaker extraction and recognition with Speaker-Beam,” in *ICASSP*, pp. 5554–5558, IEEE, 2018.
- [26] P. Denisov and N. T. Vu, “End-to-End Multi-Speaker Speech Recognition Using Speaker Embeddings and Transfer Learning,” in *Interspeech*, pp. 4425–4429, 2019.
- [27] S. Ding, Q. Wang, S.-Y. Chang, L. Wan, and I. Lopez Moreno, “Personal VAD: Speaker-Conditioned Voice Activity Detection,” in *Odyssey*, pp. 433–439, 2020.
- [28] I. Medennikov, M. Korenevsky, T. Prisyach, Y. Khokhlov, M. Korenevskaya, I. Sorokin, T. Timofeeva, A. Mitrofanov, A. Andrusenko, I. Podluzhny, A. Laptev, and A. Romanenko, “Target-Speaker Voice Activity Detection: A Novel Approach for Multi-Speaker Diarization in a Dinner Party Scenario,” in *Interspeech*, pp. 274–278, 2020.
- [29] M. Sadeghi, S. Leglaive, X. Alameda-Pineda, L. Girin, and R. Horaud, “Audio-visual speech enhancement using conditional variational auto-encoders,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1788–1800, 2020.
- [30] K. Žmolíková, M. Delcroix, K. Kinoshita, T. Higuchi, A. Ogawa, and T. Nakatani, “Speaker-aware neural network based beamformer for speaker extraction in speech mixtures,” in *Interspeech*, pp. 2655–2659, 2017.
- [31] M. Ge, C. Xu, L. Wang, E. S. Chng, J. Dang, and H. Li, “SpEx+: A Complete Time Domain Speaker Extraction Network,” in *Proc. Interspeech 2020*, pp. 1406–1410, 2020.
- [32] G. Morrone, S. Bergamaschi, L. Pasa, L. Fadiga, V. Tikhanoff, and L. Badino, “Face landmark-based speaker-independent audio-visual speech enhancement in multi-talker environments,” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6900–6904, IEEE, 2019.
- [33] T. Ochiai, M. Delcroix, K. Kinoshita, A. Ogawa, and T. Nakatani, “Multimodal SpeakerBeam: Single channel target speech extraction with audio-visual speaker clues,” in *Interspeech*, pp. 2718–2722, 2019.
- [34] H. Sato, T. Ochiai, K. Kinoshita, M. Delcroix, T. Nakatani, and S. Araki, “Multimodal attention fusion for target speaker extraction,” in *2021 IEEE Spoken Language Technology Workshop (SLT)*, pp. 778–784, IEEE, 2021.
- [35] T. Afouras, J. S. Chung, and A. Zisserman, “My lips are concealed: Audio-visual speech enhancement through obstructions,” *arXiv preprint arXiv:1907.04975*, 2019.
- [36] R. Gu, S.-X. Zhang, Y. Xu, L. Chen, Y. Zou, and D. Yu, “Multi-modal multi-channel target speech separation,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 3, pp. 530–541, 2020.
- [37] S.-W. Chung, S. Choe, J. S. Chung, and H.-G. Kang, “FaceFilter: Audio-Visual Speech Separation Using Still Images,” in *Proc. Interspeech 2020*, pp. 3481–3485, 2020.
- [38] M. Maciejewski, G. Sell, Y. Fujita, L. P. Garcia-Perera, S. Watanabe, and S. Khudanpur, “Analysis of robustness of deep single-channel speech separation using corpora constructed from multiple domains,” in *WASPAA*, pp. 165–169, 2019.
- [39] Y. Luo and N. Mesgarani, “TasNet: Surpassing ideal time-frequency masking for speech separation,” in *ICASSP*, 2018.
- [40] X. Xiao, Z. Chen, T. Yoshioka, H. Erdogan, C. Liu, D. Dimitriadis, J. Droppo, and Y. Gong, “Single-channel speech extraction using speaker inventory and attention network,” in *ICASSP*, pp. 86–90, 2019.
- [41] K. Žmolíková, M. Delcroix, K. Kinoshita, T. Higuchi, A. Ogawa, and T. Nakatani, “Learning speaker representation for neural network based multichannel speaker extraction,” in *ASRU*, pp. 8–15, IEEE, 2017.
- [42] J. Heymann, L. Drude, and R. Haeb-Umbach, “Neural network based spectral mask estimation for acoustic beamforming,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 196–200, IEEE, 2016.
- [43] H. Erdogan, J. R. Hershey, S. Watanabe, M. I. Mandel, and J. Le Roux, “Improved MVDR beamforming using single-channel mask prediction networks,” in *Interspeech*, pp. 1981–1985, 2016.
- [44] E. Vincent, R. Gribonval, and C. Févotte, “Performance measurement in blind audio source separation,” *IEEE transactions on audio, speech, and language processing*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [45] M. Delcroix, S. Watanabe, T. Ochiai, K. Kinoshita, S. Karita, A. Ogawa, and T. Nakatani, “End-to-end speakerbeam for single channel target speech recognition,” in *Interspeech*, 2019.
- [46] S. Mun, S. Choe, J. Huh, and J. S. Chung, “The sound of my voice: Speaker representation loss for target voice separation,” in *ICASSP*, pp. 7289–7293, IEEE, 2020.
- [47] X. Li, Y. Xu, M. Yu, S.-X. Zhang, J. Xu, B. Xu, and D. Yu, “MIMO Self-Attentive RNN Beamformer for Multi-Speaker Speech Separation,” in *Proc. Interspeech 2021*, pp. 1119–1123, 2021.
- [48] T. Cord-Landwehr, C. Boeddeker, T. Von Neumann, C. Zorilă, R. Doddipatla, and R. Haeb-Umbach, “Monaural source separation: From anechoic to reverberant environments,” in *IWAENC*, pp. 1–5, IEEE, 2022.
- [49] D. Ditter and T. Gerkmann, “A multi-phase gammatone filterbank for speech separation via Tasnet,” in *ICASSP*, pp. 36–40, IEEE, 2020.
- [50] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, “Front-end factor analysis for speaker verification,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 788–798, 2010.
- [51] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, “X-vectors: Robust DNN embeddings for speaker recognition,” in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 5329–5333, 2018.
- [52] K. Žmolíková, *Neural target speech extraction*. Ph.d. thesis, Brno University of Technology, Faculty of Information Technology, 2022.
- [53] A. Gabbay, A. Shamir, and S. Peleg, “Visual Speech Enhancement,” in *Proc. Interspeech 2018*, pp. 1170–1174, 2018.
- [54] B. King, I.-F. Chen, Y. Vaizman, Y. Liu, R. Maas, S. H. K. Parthasarathi, and B. Hoffmeister, “Robust speech recognition via anchor word representations,” in *Interspeech*, pp. 2471–2475, 2017.
- [55] J. Yu, S.-X. Zhang, B. Wu, S. Liu, S. Hu, M. Geng, X. Liu, H. Meng, and D. Yu, “Audio-visual multi-channel integration and recognition of overlapped speech,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2067–2082, 2021.
- [56] D. Soderoy, B. Rivet, L. Girin, J.-L. Schwartz, and C. Jutten, “An analysis of visual speech information applied to voice activity detection,” in *2006 IEEE International Conference on Acoustics Speech and Signal Processing Proceedings*, vol. 1, pp. I–I, 2006.
- [57] C. Zhang, M. Yu, C. Weng, and D. Yu, “Towards robust speaker verification with target speaker enhancement,” in *Proc. of ICASSP’21*, pp. 6693–6697, 2021.
- [58] E. Tzinis, G. Wichern, A. S. Subramanian, P. Smaragdis, and J. Le Roux, “Heterogeneous Target Speech Separation,” in *Proc. Interspeech 2022*, pp. 1796–1800, 2022.
- [59] M. Borsdorf, H. Li, and T. Schultz, “Target language extraction at multilingual cocktail parties,” in *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pp. 717–724, 2021.
- [60] Y. Ohishi, M. Delcroix, T. Ochiai, S. Araki, D. Takeuchi, D. Nizumi, A. Kimura, N. Harada, and K. Kashino, “Conceptbeam: Concept driven target speech extraction,” *MM ’22*, (New York, NY, USA), p. 4252–4260, Association for Computing Machinery, 2022.
- [61] E. Ceolini, J. Hjortkjær, D. D. Wong, J. O’Sullivan, V. S. Raghavan, J. Herrero, A. D. Mehta, S.-C. Liu, and N. Mesgarani, “Brain-informed speech separation (BISS) for enhancement of target speaker in multi-talker speech perception,” *NeuroImage*, vol. 223, p. 117282, 2020.
- [62] P. Seetharaman, G. Wichern, S. Venkataramani, and J. L. Roux, “Class-conditional embeddings for music source separation,” in *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 301–305, 2019.
- [63] H. Zhao, C. Gan, A. Rouditchenko, C. Vondrick, J. McDermott, and A. Torralba, “The sound of pixels,” in *The European Conference on Computer Vision (ECCV)*, September 2018.
- [64] G. Cantisani, M. Essid, and G. Richard, “Neuro-steered music source separation with eeg-based auditory attention decoding and contrastive-nmf,” in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 36–40, 2021.

- [65] Z. Duan, S. Essid, C. C. Liem, G. Richard, and G. Sharma, "Audiovisual analysis of music performances: Overview of an emerging field," IEEE Signal Processing Magazine, vol. 36, no. 1, pp. 63–73, 2019.
- [66] T. Ochiai, M. Delcroix, Y. Koizumi, H. Ito, K. Kinoshita, and S. Araki, "Listen to What You Want: Neural Network-Based Universal Sound Selector," in Interspeech, pp. 1441–1445, 2020.
- [67] B. Gfeller, D. Roblek, and M. Tagliasacchi, "One-shot conditional audio filtering of arbitrary sounds," in ICASSP, pp. 501–505, 2021.
- [68] X. Liu, H. Liu, Q. Kong, X. Mei, J. Zhao, Q. Huang, M. D. Plumbley, and W. Wang, "Separate What You Describe: Language-Queried Audio Source Separation," in Proc. Interspeech 2022, pp. 1801–1805, 2022.
- [69] S. S. Asaridou and J. M. McQueen, "Speech and music shape the listening brain: evidence for shared domain-general mechanisms," Frontiers in psychology, vol. 4, p. 321, 2013.