# SoundBeam: Target Sound Extraction Conditioned on Sound-Class Labels and Enrollment Clues for Increased Performance and Continuous Learning

Marc Delcroix *Senior Member, IEEE*, Jorge Bennasar Vázquez, Tsubasa Ochiai *Member, IEEE*,
Keisuke Kinoshita *Senior Member, IEEE*, Yasunori Ohishi *Member, IEEE*, Shoko Araki *Fellow, IEEE*

*Abstract*—In many situations, we would like to hear desired sound events (SEs) while being able to ignore interference. Target sound extraction (TSE) tackles this problem by estimating the audio signal of the sounds of target SE classes in a mixture of sounds while suppressing all other sounds. We can achieve this with a neural network that extracts the target SEs by conditioning it on clues representing the target SE classes. Two types of clues have been proposed, i.e., target *SE class labels* and *enrollment audio samples* (or audio queries), which are pre-recorded audio samples of sounds from the target SE classes. Systems based on SE class labels can directly optimize embedding vectors representing the SE classes, resulting in high extraction performance. However, extending these systems to extract new SE classes not encountered during training is not easy. Enrollment-based approaches extract SEs by finding sounds in the mixtures that share similar characteristics to the enrollment audio samples. These approaches do not explicitly rely on SE class definitions and can thus handle new SE classes. In this paper, we introduce a TSE framework, SoundBeam, that combines the advantages of both approaches. We also perform an extensive evaluation of the different TSE schemes using synthesized and real mixtures, which shows the potential of SoundBeam.

*Index Terms*—Target Sound Extraction, Sound Event, Sound-Beam, Few-shot adaptation, Deep Learning.

## I. Introduction

**H**uman beings can listen to a desired sound within a complex acoustic scene consisting of a mixture of various sound events (SEs). This phenomenon is called the cocktail party effect or selective hearing [1]. For example, it enables us to listen to an interlocutor in a noisy cafe, focus on a particular instrument in a song, or notice a siren on the road. One of the long-term goals of speech and audio processing research is to reproduce the selective hearing ability of humans computationally.

Target sound extraction (TSE) is one approach toward achieving this goal. We define the TSE problem as the extraction of one or multiple desired sounds from a mixture of various SEs, given user-specified clues characterizing the target SE classes. When multiple desired SE classes are selected, we output a signal that consists of the sum of all the SEs from these classes. TSE exploits auxiliary clues such
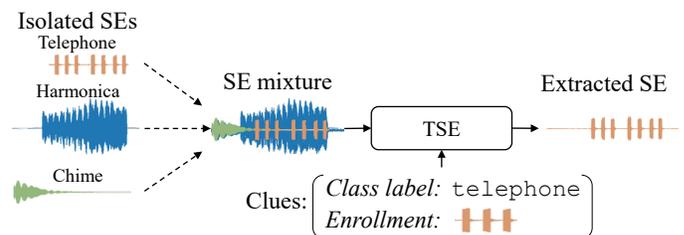
Fig. 1. Schematic diagram of a TSE system based on target SE class labels or enrollment clues.

as class labels to inform the system about the target class [2]. Figure 1 is a schematic diagram of a TSE system.

In general, sound recordings are often polyphonic [3] and contain different SEs that overlap in time, which makes the problem particularly challenging. Solving the TSE problem has many direct practical implications. For example, it would allow flexible personal hearables or hearing aids that could filter sounds depending on the situation to provide important information about our surroundings (e.g., a klaxon when we cross a street) while removing nuisances that disrupt our concentration (e.g., the same klaxon when we work at home with a window open). TSE could also be used for sound post-production to emphasize or remove [2] specific sounds in, e.g., a video recording.

Research on the processing of SEs has focused mostly on sound event classification (SEC), audio tagging (AT), and sound event detection (SED) problems [3]–[9]. SEC/AT assign SE class labels to audio signals, while SED also localizes in time the SEs in the audio. SED does not explicitly extract individual sound signals from the mixture, and thus it cannot solve the TSE problem when SEs overlap.

We can use source separation to separate a sound mixture into each source. For example, universal sound separation (USS) [10] proposes to use a neural network (NN) to separate a mixture of arbitrary SEs. However, one issue with USS is that it requires knowing or estimating the number (or maximum number) of sources in the mixture. This requirement may be challenging when dealing with a mixture of arbitrary SEs, since the number of SEs can vary greatly depending on the environment. Moreover, USS provides access to the different sources but does not identify them. Consequently, USS is also unable to solve the TSE problem.

Recently, TSE systems that can directly extract sounds from

the target SE classes have received increased interest [2], [11]–[14]. Such TSE systems use a sound extraction NN that accepts a sound mixture and an auxiliary target SE embedding vector [1] that represents the sound characteristics of the target SE . The sound extraction NN extracts the sound signals in the mixture that match the characteristics of the target SE embedding. The system thus implicitly performs "separation" and "sound identification" internally.

There are several advantages to this approach. First, the sound extraction NN extracts only the target SEs and thus has a single output, which makes the NN architecture independent of the number of sources in the mixture. Naturally, the computational complexity is also unaffected by the number of sources. Finally, the "separation" and "identification" processes are performed simultaneously, allowing us to exploit the information on the sound characteristics to help the internal "separation" process and optimize the TSE system globally.

A key component of TSE systems are the target SE embeddings. They can be derived from *class labels* or *enrollment audio samples*[2]. With the first approach, class labels represented as 1-hot vectors are mapped to an embedding space representing the different SE classes using an embedding layer, which is learned jointly with the sound extraction NN [2], [11], [17]. Such class label-based TSE systems can directly learn the target SE embeddings, which can thus optimally represent the characteristics of the SE classes for the TSE task. However, class label-based approaches assume a predetermined set of target SE classes, and thus they cannot perform TSE of *new SE classes* that were not encountered during training.

With the second approach, the target SE embeddings are derived from an enrollment audio sample, which is a short recording of an isolated sound similar to the target SE [12], [13], [16], [18], [19]. Enrollment-based TSE systems learn to extract sounds that share similar characteristics to the enrollment sample without explicitly relying on SE class labels. They can thus perform TSE of new SE classes. However, unlike class label-based approaches, the embeddings are derived from the enrollment samples. Therefore, they are not directly optimized for the SE classes, which may lead to suboptimal performance.

In this paper, we provide a unified description of the class label and enrollment-based TSE approaches and introduce *SoundBeam*, which is a TSE model that combines both approaches. SoundBeam generalizes the target *speech* extraction framework called SpeakerBeam [18] to arbitrary *sounds*. "Beam" refers to focusing on a particular sound in a mixture but does not imply using a microphone array or a beamformer [20]. Indeed, this work focuses on single-channel processing.

SoundBeam learns an embedding space common to the class label- and enrollment-based target SE embeddings by sharing the extraction NN. Consequently, we can use optimal SE embeddings for known SE classes with the class label embeddings and enrollment embeddings for new SE classes.

[1] Some works call the "target SE embedding vector" a "conditioning vector" [11], a "conditioning embedding" [12], or a "query embedding" [15].

[2] Some works use the term "audio query" [15], [16] or conditioning audio [12] instead of "enrollment audio samples."

Furthermore, we show that SoundBeam has potential for *continuous learning*, since it can learn class label embeddings for new classes with few-shot adaptation.

We presented the initial ideas of SoundBeam in our previous work [13]. This paper provides a more detailed explanation of the approach and an extensive evaluation of the TSE frameworks, covering the following four practical aspects.

**Extraction of new SE classes with few-shot adaptation**: The number of potential SE classes is huge, if not infinite, and therefore it may not be realistic to develop a system that can handle all possible SE classes. Moreover, the target SE classes may depend on the application scenarios. Therefore, we believe TSE systems should be able to learn new SE classes after their deployment, i.e., continuous learning. We propose a few-shot adaptation approach, which allows SoundBeam to learn optimal SE embeddings from new SE classes using a few enrollment samples. We introduced this adaption scheme earlier [13]. Here, we provide a deeper experimental evaluation, including a class-wise analysis of the effect of adaptation and experiments on the influence of the number of enrollment samples on performance.

**Handling of inactive target SE classes**: For many practical applications, a TSE system should output a silent or zero signal if the target SE is *inactive* in the mixture. However, most prior works have not rigorously evaluated this scenario. We discuss how to learn to handle inactive SE classes and perform a thorough experimental analysis with the different TSE frameworks.

**Simultaneous extraction of multiple SEs:** In some applications, there may be several target SE classes. The target would thus consist of a mixture of target SE sounds. We show that we can extract several SE sounds simultaneously simply by adding the SE embeddings of the different target classes. We first introduced this idea for class label-based TSE [2]. Here, we develop it for the enrollment and SoundBeam models.

**Challenges faced with real recordings:** Finally, we explore the challenges of processing real mixtures, which include evaluation without access to isolated target SE sources and adaptation to real recording conditions without access to strong labels consisting of the isolated SE signals. We perform preliminary experiments in these directions using real mixtures taken from the Freesound dataset (FSD) [21].

The remainder of the paper is as follows. We review related works in Section II. The TSE problem is formalized while introducing the notations in Section III. We present details of the TSE frameworks by introducing the class label, enrollment, and SoundBeam in Section IV. In Section V, we discuss how to handle the four practical aspects discussed above. We then provide experimental results comparing the different TSE frameworks in Section VI. Finally, Section VII concludes the paper and previews possible future works.

## II. RELATED WORKS

### A. Source separation

TSE is related to source separation since both problems are regression problems from a mixture signal. Source sepa-

ration, including speech [22]–[24], music [25]–[29], and universal sound separation [10], [30], [31], has made tremendous progress with the advent of deep NNs (DNNs). We can borrow some of the ideas proposed for source separation directly to design TSE systems, such as the model architectures and the training losses. For example, we derive our implementation of the different TSE frameworks from the fully-convolutional time-domain audio separation network (Conv-TasNet) [24], which was originally proposed for speech separation.

There is a vast diversity of sounds occurring in our everyday environments. Therefore, universal sound separation systems need to handle a large number of possible SE classes, $N$. Here, the number of SE classes, $N$, includes all the SE classes that can occur in any recording we wish to process. Each mixture would typically consist of a smaller number of sound sources, $M$, taken from a subset of the possible SE classes. In this paper, we consider all sounds that are from the same SE class as a single source in the mixture. A system that could handle recordings in various situations, such as street, home, orchestra, train stations, etc., would have a large number of possible SE classes, $N$. However, since all situations do not contain sounds from all SE classes, the number of sources in a mixture, $M$, is generally smaller, i.e., $M < N$. In our experiments, we use datasets with a number of possible SE classes, $N$, of 41 or 61, but a number of sound sources in a mixture, $M$, from two to nine, as described in Section VI and in Table I.

There are two categories of sound separation systems. The first type of system uses *as many outputs as the total number of possible SE classes*, $N$ [25], [26], [28], [32]. The separation systems in this category perform both separation and identification and could thus be used directly for TSE. These approaches have mainly been investigated for music processing, where the number of possible SE classes, $N$, is relatively small, e.g., four instruments ($N = 4$) in many studies [25]. It may be challenging to scale these systems to a large number of SE classes, $N$, that a TSE system would require because the computational complexity would increase significantly [29].

As the second type of sound separation system, USS proposes instead to use *as many outputs as the maximum number of sources in a mixture, $M$*, which is typically much smaller than the number of SE classes, $N$ [10]. Such separation systems can be trained using permutation invariant training (PIT) [22]. USS systems can handle an arbitrarily large number of SE classes, $N$, and potentially new classes not encountered during training. However, they are limited by the number of sources in the mixtures, $M$, they can handle. Arguably, this may be a significant issue because the number of SEs in our everyday surroundings vary significantly and can be relatively large. Moreover, they separate the mixtures without identifying the sounds.

TSE could be achieved by combining USS with SEC, i.e., first separating the mixture into all its source signals and then identifying the target SE class using SEC. Note that some works proposed to combine source separation and SEC/SED to improve SED performance but not to achieve TSE [33]–[37]. To the best of our knowledge, combining USS with SEC

has not been proposed explicitly for TSE, but we still consider it a natural baseline in our experiments of Section VI-D. Although intuitive, this approach has several drawbacks. First, it inherits from USS the requirement of knowing or estimating the number of sources in the mixture, $M$. Second, the SEC process must be performed on all separated signals, which may become computationally expensive when dealing with many potential sources. Third, such a cascade combination may not be optimal since source separation may introduce errors or processing artifacts that can impact SEC performance. Besides, the separation process is independent of the target SE class, although knowing the sound characteristics of the target SE class could help improve separation.

### B. Target sound extraction

Recently, there has been increased interest in approaches that aim at extracting a target signal in an audio mixture [2], [11]–[14], [16]–[19], [38]–[40]. These works cover various domains of applications and are implemented with various network architectures. However, they share the same idea of conditioning the extraction process on auxiliary clues to identify a target signal in a mixture.

The domain of application of TSE includes speech [18], [38], music [16], [17], [19], [39], [40], and universal sounds [2], [11]–[14]. Various types of auxiliary clues have been proposed, including enrollment audio samples [12], [16], [18], [19], [38], class labels [2], [11], [40], video signals of the target source [39], and recently even onomatopoeia [14]. We deal with universal sounds, which is challenging because it implies dealing with a much larger number of SE classes, $N$, than for music (i.e., in this paper, we use up to 61 SE classes) and with a greater variety of sounds than speech signals.

This work focuses on approaches based on enrollment and class labels. Enrollment-based TSE was introduced first for target speech and music extraction [16], [41], where the enrollment sample consists of an utterance from the target speaker or an audio query of a target instrument. Similar ideas have recently been applied to sound extraction [12]. Meanwhile, class label-based TSE was introduced in concurrent works [2], [11].

We reproduce systems based on the concepts of prior enrollment- and class label-based TSE approaches [2], [11], [12] using the same formalization to allow a fair comparison between them. The main difference between this paper and previous works [11], [12] is that we base our experiments on a dataset that contains isolated SEs annotated with class labels, whereas they did not assume that class labels were available. However, we could also use similar schemes to previous works [11], [12] to train SoundBeam on datasets without class labels annotations. Furthermore, the model architecture also differs, i.e., our study is based on Conv-TasNet (as in e.g., [2]), while the above works [11], [12] are based on U-Net.

### C. Training on real mixtures

Training a system with real mixtures requires modifying the classical supervised learning schemes since isolated reference sources are unavailable. There have been several proposals to

train separation systems from real mixtures, including training on mixtures-of-mixtures [42] and using losses that do not require clean sources, such as a generative adversarial network (GAN) [43], [44] or an SEC-based loss [45]. Several TSE systems have been trained with a similar principle as mixtures-of-mixtures [11], [12], [39], [46]. For example, one work [11] used SED to detect isolated SE in real mixtures, while another [12] showed that the enrollment-based TSE could also be trained even when the audio query consists of sound mixtures.

Although mixtures-of-mixtures approaches permit training with real mixtures, they still generate artificial recordings since unrelated audio signals are mixed. An alternative consists of retraining directly on real mixtures using a weakly supervised loss based on SEC [45]. We adopt this scheme originally proposed for training separation systems with a small number of SE classes (i.e., $N = 5$ in [45]) and apply it to adapt a TSE system on real data with a large number of classes (i.e., $N = 61$, which corresponds to the number of SE class encountered during training in our experiments, as described in Section VI-A).

## III. TARGET SOUND EXTRACTION PROBLEM FORMULATION

We consider the problem of extracting sounds of one or multiple target SE classes from a mixture of SEs captured with a single microphone. The observed mixture signal is given as,

$$\mathbf{y} = \sum_{n=1}^{N} \mathbf{x}^n + \mathbf{b}, \tag{1}$$

where $\mathbf{y} \in \mathbb{R}^T$, $\mathbf{x}^n \in \mathbb{R}^T$, and $\mathbf{b} \in \mathbb{R}^T$ are the observed single-channel mixture, the source signal from the $n$-th SE class, and the background noise, respectively. The background noise, $\mathbf{b}$, includes ambient noise and could also potentially include sounds from SEs that are not defined in the SE classes of the system. $T$ is the signal duration.

We assume $\mathbf{x}^n = \mathbf{0}$ when no source from the $n$-th SE class is active, where $\mathbf{0}$ denotes a vector of all zeros. Later, we refer to this case as an *inactive SE class*. Note that $\mathbf{x}^n$ may consist of multiple SE sources from the same SE class. For example, if the mixture includes barking sounds from several dogs, the source signal $\mathbf{x}^n$ for the barking SE class will consist of a mixture of all barking sounds. In Eq. (1), the summation is over all possible SE classes, $N$, including target and non-target ones. Since, in practice, many SE classes are often inactive, the actual number of sources in a mixture, $M$, is usually smaller than the total number of SE classes, $N$.

### A. Single target extraction

In this paper, we mostly deal with TSE of a single target SE class. In this case, the goal of TSE is to estimate the target signal, $\mathbf{x}^s$, where $s$ is the index of the target SE class. If the target SE class is inactive, TSE should estimate a zero signal, i.e., $\mathbf{x}^s = \mathbf{0}$.

TSE requires clues to indicate the target SE classe. In the following, we consider two types of clues, i.e., class labels and enrollment audio samples. Target class labels can be represented as 1-hot vectors, $\mathbf{o}^s = [o_1^s, \dots, o_N^s]^T$, where

$$o_n^s = \begin{cases} 1 & \text{if } n = s, \\ 0 & \text{if } n \neq s. \end{cases} \tag{2}$$

A class label-based TSE system estimates thus the target source given the 1-hot vector as,

$$\hat{\mathbf{x}}^s = \text{TSE}(\mathbf{y}, \mathbf{o}^s), \tag{3}$$

where $\hat{\mathbf{x}}^s \in \mathbb{R}^T$ is an estimate of the target source and $\text{TSE}(\cdot)$ represents a TSE system.

Alternatively, we can also use an enrollment audio sample of the target SE class, $\mathbf{a}^s \in \mathbb{R}^{T_a}$, as a clue, where $T_a$ is the duration of the enrollment sample. An enrollment-based TSE system performs thus the following,

$$\hat{\mathbf{x}}^s = \text{TSE}(\mathbf{y}, \mathbf{a}^s). \tag{4}$$

### B. Multi-target extraction

We can generalize TSE to the simultaneous extraction of sounds from multiple target SE classes [2]. We denote by $S = \{s_j\}_{j=1}^J$ the set of indexes of the target SE classes that we want to extract, where $J$ is the number of target SE classes, which can differ for each processed mixture. For example, if we want to extract sounds from two classes simultaneously ($J = 2$), such as both piano and guitar, $S = \{s_1 = n_{\text{piano}}, s_2 = n_{\text{guitar}}\}$, where $n_{\text{piano}}$ and $n_{\text{guitar}}$ are the indexes of the piano and guitar classes, respectively. The goal of TSE becomes estimating the target signal $\mathbf{x}^S$, which consists of the sum of the signals from all target SE classes:

$$\mathbf{x}^S = \sum_{n \in S} \mathbf{x}^n. \tag{5}$$

Note that if all target SEs are inactive, TSE should estimate a zero signal, i.e., $\mathbf{x}^S = \mathbf{0}$.

When using class labels clues, we can generalize the 1-hot vector of Eq. (2) to a n-hot vector, $\mathbf{o}^S = [o_1^S, \dots, o_N^S]^T$, where

$$o_n^S = \begin{cases} 1 & \text{if } n \in S; \\ 0 & \text{if } n \notin S. \end{cases} \tag{6}$$

A class label-based multi-target TSE system estimates the mixture of the target sources as,

$$\hat{\mathbf{x}}^S = \text{TSE}(\mathbf{y}, \mathbf{o}^S), \tag{7}$$

where $\hat{\mathbf{x}}^S$ is an estimate of the target signal $\mathbf{x}^S$.

When using enrollment audio samples as clues, we consider here that we have separate enrollment samples for each of the target SE classes. An enrollment-based multi-target TSE system performs thus the following,

$$\hat{\mathbf{x}}^S = \text{TSE}(\mathbf{y}, \{\mathbf{a}^s\}_{s \in S}), \tag{8}$$

where $\{\mathbf{a}^s\}_{s \in S} = \{\mathbf{a}_{s_1}, \dots, \mathbf{a}_{s_J}\}$ is the set of separate enrollment audio samples from each target SE class.

In most parts of the paper, we use the notations for single target extraction to simplify the explanations. We discuss the multi-target extraction in Section V-C.

## IV. SoundBeam framework for TSE

In this section, we introduce the generic framework for NN-based TSE from which we derive SoundBeam. The TSE models share a common structure and are composed of two modules, i.e., (1) a sound extraction NN and (2) a clue encoder that computes the target SE embeddings. Figure 2 schematically overviews the TSE frameworks for (a) class label-based TSE, (b) enrollment-based TSE, and (c) SoundBeam.

### A. Sound extraction NN

The sound extraction NN estimates the target signal, $\mathbf{x}^s$, from the mixture, $\mathbf{y}$, and a $D$-dimensional target embedding vector, $\mathbf{e}^s \in \mathbb{R}^D$, as,

$$\hat{\mathbf{x}}^s = f(\mathbf{y}, \mathbf{e}^s), \tag{9}$$

where $f(\cdot)$ is the sound extraction NN. The embedding vector, $\mathbf{e}^s$, represents the characteristics of the target SE class. It is computed with the clue encoder derived from the class label or enrollment clues described in section IV-B.

There are many possible ways to implement the extraction NN. We borrow the implementation from Conv-TasNet because it has been widely used for speech separation [24] and target speech extraction [47]. However, our derivation is general enough to be extended to other network architectures such as U-Net or to other systems using the short-time Fourier transform (STFT)/inverse STFT (iSTFT) as encoder/decoder layers. The detailed processing is as follows.

The time domain input signal, $\mathbf{y}$, is processed by a trainable encoder layer as,

$$\mathbf{Y} = \mathrm{Encoder}(\mathbf{y}), \tag{10}$$

where $\mathbf{Y} \in \mathbb{R}^{D \times T'}$ is the encoded mixture representation, $T'$ is the number of time frames, and $\mathrm{Encoder}(\cdot)$ is the encoder layer, which consists of a 1-D convolution layer [24]. The estimated target signal, $\hat{\mathbf{x}}^s$, is obtained as,

$$\hat{\mathbf{x}}^s = \mathrm{Decoder}(\mathbf{M}^s \odot \mathbf{Y}), \tag{11}$$

where $\mathbf{M}^s$ is a target mask, $\odot$ is the element-wise product, and $\mathrm{Decoder}(\cdot)$ is a decoder layer that maps the output representation of the extracted signal back to the time domain.

Compared to Conv-TasNet, TSE only requires computing the mask for the target SE class instead of one for each source. The target mask, $\mathbf{M}^s$, is computed as follows,

$$\mathbf{Z} = \mathrm{MixBlock}(\mathbf{Y}), \tag{12}$$
$$\mathbf{Z}^s = \mathrm{Adapt}(\mathbf{Z}, \mathbf{e}^s), \tag{13}$$
$$\mathbf{M}^s = \mathrm{TgtBlock}(\mathbf{Z}^s), \tag{14}$$

where $\mathbf{Z} \in \mathbb{R}^{D \times T'}$ and $\mathbf{Z}^s \in \mathbb{R}^{D \times T'}$ are the internal representation of the mixture and the target, respectively. $\mathrm{MixBlock}(\cdot)$ is the lower part of the sound extraction NN, which transforms $\mathbf{Y}$ into a general internal representation of the mixture independent of the target SE class. $\mathrm{TgtBlock}(\cdot)$ is the upper part of the sound extraction NN, which computes the target mask, $\mathbf{M}^s$, that specifically extracts the sounds from the target SE classes (note that in [13], $\mathrm{MixBlock}(\cdot)$ and $\mathrm{TgtBlock}(\cdot)$ are referred to as $\mathrm{ExtractBlock1}(\cdot)$ and

$\mathrm{ExtractBlock2}(\cdot)$, respectively). $\mathrm{Adapt}(\cdot)$ is an adaptation layer. $\mathrm{MixBlock}(\cdot)$ and $\mathrm{TgtBlock}(\cdot)$ are implemented as stacks of 1-D convolutional blocks [24].

One of the key differences from the original Conv-TasNet is the adaptation layer, which combines the internal representation of the mixture, $\mathbf{Z}$, with the target SE embedding, $\mathbf{e}^s$. It plays the crucial role of conditioning the extraction process on the target SE classes. Various types of adaptation layers have been proposed in the context of target speech extraction, including factorized [18], addition/concatenation [38], multiplication [47], attention [48], and Feature-wise Linear Modulation (FiLM) [12] layers. Here, we use the simple yet sufficiently powerful [18] element-wise multiplication layer as,

$$\mathbf{z}^s{}_t = \mathbf{z}_t \odot \mathbf{e}^s, \tag{15}$$

where $\mathbf{z}^s{}_t \in \mathbb{R}^D$ and $\mathbf{z}_t \in \mathbb{R}^D$ are the $t$-th frame of $\mathbf{Z}$ and $\mathbf{Z}^s$, respectively.

### B. Clue encoder

The different TSE frameworks differ by the clue encoders they use. The clue encoder computes the SE embedding vectors from the class label or enrollment clues. Below, we describe these two configurations and a mixed encoder, which SoundBeam employs.

*1) Class label encoder:* The class label encoder consists of an embedding layer that converts the 1-hot vector, $\mathbf{o}^s$, defined in Eq. (6) into a continuous representation as,

$$\mathbf{e}^{\mathrm{class},s} = \mathbf{W}\mathbf{o}^s, \tag{16}$$

where $\mathbf{W} = [\mathbf{e}^{\mathrm{class},1}, \ldots, \mathbf{e}^{\mathrm{class},N}] \in \mathbb{R}^{D \times N}$ is an embedding matrix whose columns contain the embedding vectors of each SE class.

Figure 2(a) shows a diagram of the class label-based TSE model similar to previous works [2], [11]. By learning the class label encoder and the sound extraction NN jointly, we can optimize the SE embedding vectors directly for the TSE task. However, the SE classes that the method can handle are fixed by the embedding matrix, $W$, and limited to the $N$ *"seen" SE classes* encountered during training. Consequently, we cannot directly use it for the extraction of new SE classes.

*2) Enrollment encoder:* The enrollment encoder computes the target SE embedding vector from an enrollment audio sample of the target SE class, $\mathbf{a}^s$. Here, we use a sequence summary NN [18], [49] to convert the variable-length enrollment, $\mathbf{a}^s$, to a vector of fixed dimension as,

$$\mathbf{e}^{\mathrm{enrl},s} = g(\mathbf{a}^s). \tag{17}$$

$g(\cdot)$ is an enrollment encoder NN implemented as,

$$\mathbf{A}^s = \mathrm{Encoder}(\mathbf{a}^s) \tag{18}$$
$$\mathbf{Z}^a = \mathrm{EnrlBlock}(\mathbf{A}^s) \tag{19}$$
$$\mathbf{e}^{\mathrm{enrl},s} = \frac{1}{T'_a} \sum_t \mathbf{z}^a_t, \tag{20}$$

where $\mathrm{Encoder}(\cdot)$ is an encoder layer similar to that in Eq. (10) but with different parameters, $\mathrm{EnrlBlock}(\cdot)$ consists of a stack of 1-D convolution blocks, $\mathbf{A}^s \in \mathbb{R}^{D \times T'_a}$ and $\mathbf{Z}^a \in$
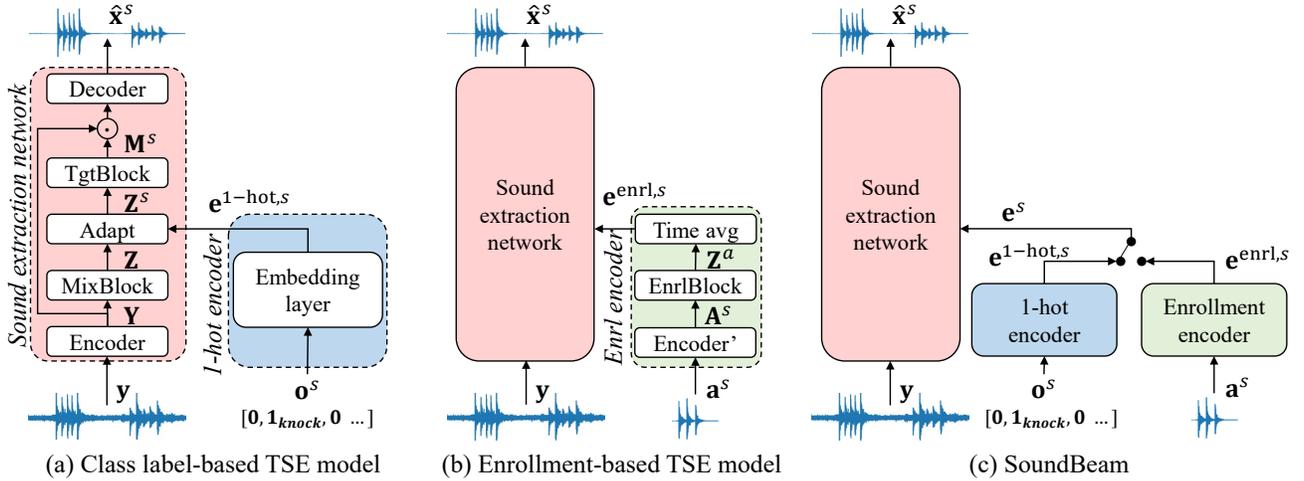
Fig. 2. Diagrams of three different TSE frameworks.

$\mathbb{R}^{D \times T'_a}$ are internal representations of the enrollment, $\mathbf{z}^a_t \in \mathbb{R}^D$ is the $t$-th time frame of $\mathbf{Z}^a$, and $T'_a$ is the number of time frames. The last layer performs average pooling over the time dimension.

Figure 2(b) shows a diagram of enrollment-based TSE similar to previous works [12], [19]. Compared to the class label-based TSE, the enrollment-based approach does not explicitly rely on well-defined class labels but instead identifies and extracts sounds in the mixture based on their similarity to the enrollment. Consequently, the enrollment-based model can generalize to a new target SE class as long as an enrollment audio sample of the new class can be collected and the model can be trained with a sufficient variety of SE classes. On the other hand, the embedding vector, $\mathbf{e}^{\text{enrl},s}$, is not directly optimized for the target SE class, which may result in lower performance compared with the class label-based model when extracting sounds from seen SE classes.

*3) Mixed encoder of SoundBeam:* Both class label and enrollment models use an embedding vector to characterize the target SE class. When trained independently, the embedding vectors of both models are mapped to different embedding spaces. However, we can enforce mapping them to a joint embedding space by designing a model that shares the sound extraction NN and alternates between the class label and enrollment encoders during training. The sound extraction network would thus be trained to output the same extracted signal when conditioning on embedding vectors derived from either the class labels or enrollment samples. One way to achieve this is if the class label and enrollment encoders provide similar representations for SEs of the same class, i.e., the embedding vectors derived from the class labels and enrollment samples are similar and mapped to the same region of the embedding space. Note that we could also include explicit constraints to further enforce learning a joint embedding space [13], but in the experiments of this paper, such an additional constraint was not necessary.

At test time, we can use either the class label or the enrollment encoder to perform the extraction. Figure 2(c) is a diagram of SoundBeam, which uses the mixed encoder. The

switch indicates that we can use either the class label or the enrollment encoder.

SoundBeam offers several advantages. First, the alternate training scheme, described in Section IV-C, acts as multi-task training, which may lead to better performance. Moreover, we can employ at test time the class label encoder to use optimal target SE embeddings for seen SE classes or the enrollment encoder to handle new classes. Finally, we can derive the few-shot adaptation scheme described in section V-A, which enables us to learn optimal embeddings for new SE classes with a few enrollments.

### C. Training TSE systems

This section introduces the fully supervised training scheme of the TSE systems, where we assume that the sound mixture, $\mathbf{y}$, the target sound signal, $\mathbf{x}^s$, and the class label or enrollment clues are available. We learn the sound extraction NN and the clue encoder jointly to ensure that the embedding vectors capture the information needed for the TSE task.

We use the negative thresholded signal-to-distortion ratio (SDR) [42] as extraction loss,

$$\mathcal{L}^{\text{ext}}(\hat{\mathbf{x}}^s, \mathbf{x}^s) = -10 \log_{10}\left(\frac{\|\mathbf{x}^s\|^2}{\|\mathbf{x}^s - \hat{\mathbf{x}}^s\|^2 + \tau\|\mathbf{x}^s\|^2}\right), \quad (21)$$

$$\Leftrightarrow 10 \log_{10}\left(\|\mathbf{x}^s - \hat{\mathbf{x}}^s\|^2 + \tau\|\mathbf{x}^s\|^2\right) \quad (22)$$

where $\tau$ is a threshold that limits the upper value of the loss, thus preventing the low error/distortion training samples from dominating the gradient. Here, we follow the prior work [42] and fix in all experiments $\tau = 10^{-\frac{\eta}{10}}$, where $\eta = 30$ dB is a soft threshold value following.

To train SoundBeam, we sum the extraction loss obtained with the class label and enrollment encoders as,

$$\mathcal{L}^{\text{SoundBeam}} = \alpha \mathcal{L}^{\text{ext}}\left(\hat{\mathbf{x}}^s = f\left(\mathbf{y}, \mathbf{e}^{\text{class},s}\right), \mathbf{x}^s\right) + (1-\alpha)\mathcal{L}^{\text{ext}}\left(\hat{\mathbf{x}}^s = f\left(\mathbf{y}, \mathbf{e}^{\text{enrl},s}\right), \mathbf{x}^s\right), \quad (23)$$

where $\alpha$ is a multi-task weight that we fix to $\alpha = 0.5$ in all experiments. This loss is computed by propagating each training sample twice by alternating the clue encoders. In

other words, for each mini-batch, we first propagate with the class labels encoder ($\hat{\mathbf{x}}^s = f\left(\mathbf{y}, \mathbf{e}^{\text{class},s}\right)$) and then with the enrollment encoder ($\hat{\mathbf{x}}^s = f\left(\mathbf{y}, \mathbf{e}^{\text{enrl},s}\right)$) while accumulating the gradients, so that each training sample is extracted with both clues. This enables multi-task learning, which can have a beneficial effect on performance. More importantly, SoundBeam learns a common embedding space for class label-based and enrollment-based TSE by sharing the same sound extraction NN and training simultaneously with both the class label and enrollment encoders.

## V. PRACTICAL ISSUES

### A. Extraction of new SE classes with few-shot adaptation

It is a challenging task to create a universal TSE system that can cover all SE classes at its initial deployment since there is a wide variety of possible SEs that occur in our daily environments. In addition, the target SE classes may depend on the environments and applications. We believe that an essential property of a TSE system is to extend its capabilities flexibly and be able to continuously learn how to extract new target SE classes given a few audio samples of the new classes.

The class label model cannot handle new SE classes as discussed in Section IV-B1. The enrollment model performs extraction based on sound similarities between the characteristics of the sounds in the mixture and the enrollment. Consequently, it can naturally generalize to new SE classes if we can record enrollment audio samples for these classes. However, the embedding vectors are not optimized directly for extracting the new SE classes. Using SoundBeam, we can combine the advantages of both approaches and learn to extract sounds from new SE classes using a few enrollments (few-shot adaptation). We describe the procedure for this below.

Let $\tilde{s} = N + 1$ be the index of a new SE class. Here, we describe the addition of a single new SE class, but we can equivalently add multiple new SE classes simultaneously. We assume that we can collect $K$ enrollments, $\{\mathbf{a}_1^{\tilde{s}}, \ldots, \mathbf{a}_K^{\tilde{s}}\}$, from that new class. First, we compute an average embedding, $\mathbf{e}^{\tilde{s}}$, for the new SE class using the enrollment encoder as,

$$\mathbf{e}^{\tilde{s}} = \frac{1}{K} \sum_{k=1}^{K} g(\mathbf{a}_k^{\tilde{s}}). \tag{24}$$

Because SoundBeam learns a common embedding space for the embedding vectors obtained with the class label and enrollment encoders, this average embedding can be used as a new entry to the embedding matrix of the class label encoder as,

$$\tilde{\mathbf{W}} = [\mathbf{W}, \mathbf{e}^{\tilde{s}}] \in \mathbb{R}^{D \times (N+1)}, \tag{25}$$

where $\tilde{\mathbf{W}}$ is a new embedding matrix. This process would already enable extracting sounds from the new SE class with the class label encoder, but the new embedding vector is not yet optimized for the TSE purpose.

Consequently, we propose to fine-tune the new embedding vector, $\mathbf{e}^{\tilde{s}}$, which we call *few-shot adaptation to new SE classes*. First, we generate adaptation data by mixing the $K$ enrollments with SE signals from the original training data. We then retrain the new embedding vector, $\mathbf{e}^{\tilde{s}}$, using the adaptation data following a similar training procedure as in Section IV-C. Note that we can optimize the embedding vector of the new SE class while preserving the performance on the original SE classes by fixing all model parameters but the new embedding vector during the fine-tuning process. This simple approach allows learning new SE classes with few samples while avoiding catastrophic forgetting [50], making it suitable for continuous learning.

Note that we assume that we could record a few enrollment audio samples of the new SE classes. This implies being able to record isolated SEs of the new class, which may be challenging in real-life recordings where many sounds often overlap. Since our proposed few-shot adaptation requires only a small number of enrollment samples, we assume that the user could carefully record sounds of the target SE class or select portions of a recording where the target SE class is clearly dominant. Extension of the approach for noisy enrollments will be part of our future works.

### B. Handling of inactive target SE classes

Dealing with inactive target SE classes is another important issue of TSE systems. For example, when the target SE class is dog barking, but there are no barking sounds in the mixture, the output should be silent, i.e., $\mathbf{x}^s = \mathbf{0}$. However, it remains unclear how well TSE systems can learn this behavior since most studies focused only on active SE classes. For example, one proposed TSE system [11] was trained with data that included inactive classes, but the system was not evaluated for inactive class extraction. More recently, another approach [51] proposed evaluating the extraction for inactive speakers but not for SEs.

The problem of extraction of inactive classes is related to the problem of separation with an unknown number of speakers, where the number of sources in the mixture can be smaller than the number of outputs of the separation system [30]. The separation systems address this problem by outputting zero signals for the inactive sources, which is similar to the inactive class problem for TSE. The extraction loss of Eq. (21) is ill-defined when $\mathbf{x}^s = \mathbf{0}$. However, we can define an extraction loss for the inactive SE cases as [30],

$$\mathcal{L}^{\text{inactive}}(\hat{\mathbf{x}}^s, \mathbf{y}) = 10 \log_{10}\left(\|\hat{\mathbf{x}}^s\|^2 + \tau^{\text{inactive}}\|\mathbf{y}\|^2\right), \tag{26}$$

where $\tau^{\text{inactive}}$ is a soft threshold set at $\tau^{\text{inactive}} = 10^{-2}$. This loss consists of the denominator term of Eq. (21) as shown in Eq. (22), with a different setting for the soft threshold (i.e., $\mathbf{x}^s$ replaced by $\mathbf{y}$).

During training, following a previous study [11], we include 10 % of inactive samples (IS) to learn how to handle inactive SE classes. The training loss thus becomes,

$$\mathcal{L}(\hat{\mathbf{x}}^s, \mathbf{x}^s, \mathbf{y}) = \begin{cases} \mathcal{L}^{\text{active}}(\hat{\mathbf{x}}^s, \mathbf{x}^s) & \text{if } \mathbf{x}^s \neq \mathbf{0}, \\ \mathcal{L}^{\text{inactive}}(\hat{\mathbf{x}}^s, \mathbf{y}) & \text{if } \mathbf{x}^s = \mathbf{0}, \end{cases} \tag{27}$$

where $\mathcal{L}^{\text{active}}(\hat{\mathbf{x}}^s, \mathbf{x}^s)$ is the active loss defined in Eqs. (21) and (23). For SoundBeam, $\mathcal{L}^{\text{inactive}}(\hat{\mathbf{x}}^s, \mathbf{y})$ is computed with the class label and enrollment encoders, as for the active cases.

Note that we opted for a scale-dependent extraction loss for $\mathcal{L}^{\text{active}}$ as shown in Eq.(21) instead of the scale-independent

version widely used for source separation [24] because we believe that the scale of the output signal may matter in practical applications to detect inactive target SE classes. For example, we can evaluate how well the system could internally detect active/inactive target classes by looking at the attenuation from the mixture, which would not be possible when using a scale-independent loss as the system could choose to output signals with, e.g., very low energy.

If the TSE system could output zero signals for inactive SE classes, it would be possible to detect whether a class is active by considering the attenuation ratio between the mixture and the extracted signal, $\mathcal{A}^{\text{mixture}}(\hat{\mathbf{x}}^s, \mathbf{y})$, as,

$$\mathcal{A}^{\text{mixture}}(\hat{\mathbf{x}}^s, \mathbf{y}) = -10 \log_{10} \left( \frac{\|\mathbf{y}\|^2}{\|\hat{\mathbf{x}}^s\|^2} \right). \qquad (28)$$

Consequently, an SE class would be considered inactive if the ratio were smaller than a threshold. We use this approach to evaluate the different TSE systems in cases of inactive SE classes.

### C. Simultaneous multi-target extraction

For some applications, the user may want to hear sounds from multiple target SE classes, e.g., both sirens and klaxon sounds, when walking in the street. We can permit this by independently performing TSE for each target SE class and then summing up the extracted signals. However, this naive approach would increase the computational complexity by the number of target SE classes. We propose instead learning to simultaneously extract a mixture of target SE classes [2].

For simultaneous multi-target extraction, the embedding vector should characterize all target SE classes. For class label-based TSE, we can replace the 1-hot vector, $\mathbf{o}^s$, with an n-hot vector, $\mathbf{o}^S$, that represents multiple SE classes as defined in Eq. (6). The class label encoder of Section IV-B1 can naturally generalize to accept n-hot vectors as,

$$\mathbf{e}^{\text{class},S} = \mathbf{W}\mathbf{o}^S = \sum_{s \in S} \mathbf{W}\mathbf{o}^s. \qquad (29)$$

Equation (29) shows that the multi-target embedding vector, $\mathbf{e}^{\text{class},S}$, is simply the sum of the embedding vectors of all target SE classes in $S$.

Similarly, for the enrollment-based TSE, we define the multi-target embedding vector as the summation of the embedding vectors obtained from the individual enrollments of each class as,

$$\mathbf{e}^{\text{enrl},S} = \sum_{s \in S} \mathbf{e}^{\text{enrl},s} = \sum_{s \in S} g(\mathbf{a}^s), \qquad (30)$$

where $\mathbf{a}^s$ consists of isolated recordings from each target SE class.

We learn simultaneous extraction by including multi-target training samples, performing extraction with the above multi-target embeddings, and computing the loss with the multi-target signal, $\mathbf{x}^S$, instead of the single-target, $\mathbf{x}^s$, in Eq. (21). During training, we randomly select the number of target SE classes so that the same system can learn how to extract various numbers of classes.

TABLE I
DETAILS OF DIFFERENT DATASETS.

| Data set | Nb Classes $N$ | Audio Dur. $T$ | No. of sources $M$ | No. of mixtures Train | Valid. | Test |
|---|---|---|---|---|---|---|
| Single target (41) | 41 | 6s | 3 | 50k | 10k | 3k |
| Single target (61) | 61 | 6s | 3 | 50k | 10k | 3k |
| Multi-target | 61 | 6s | 3-5 | 50k | 10k | 3k |
| Few-shot adaptation | 41+20 | 6s | 3 | 3k | 500 | - |
| Real mixtures | 61 | 1-30s | 2-9 | 406 | 175 | 727 |
| SEC | 61 | 1-30s | 1 | 18.8k | - | 3.6k |

### D. Challenges faced with real recordings

Finally, we discuss the challenges of applying TSE to real recordings. Training a TSE system with the supervised loss of Eq. (21) requires access to the target source signal, $\mathbf{x}^s$. This requirement implies that we can only train with simulated mixtures since the source signals are usually unavailable for real recordings. However, it is challenging to create a large number of realistic mixtures of SE sounds by simulation due to, for example, the difficulty of correctly recording isolated SE sounds, the diversity of acoustic scenes, sound occurrence patterns, and duration. Therefore, there may be a significant mismatch between the training and testing conditions, which could impede extraction performance when processing real recordings.

There are various ways to tackle this problem, as discussed in subsection II-C. Here, we hypothesize that the classification performance of a SEC system should improve the better we extract SEs. Therefore, we explore retraining a TSE system on real recordings using a pre-trained SEC model to compute a weakly supervised training loss, similar to a previous work [45], instead of the supervised loss of Eq. (21). The SEC loss is computed on the output of the TSE system as,

$$\mathcal{L}^{\text{SEC}}(\hat{\mathbf{x}}^s, s) = \text{CE}\left( c(\hat{\mathbf{x}}^s), s \right), \qquad (31)$$

where $\text{CE}(\cdot)$ is the binary cross-entropy and $c(\cdot)$ is an SEC model with fixed parameters [45]. The SEC loss can be computed without clean reference signals as long as the SE class labels, $s$, are available (i.e., weak labels).

A similar loss has been proposed to train sound separation models [45]. However, compared with separation, TSE uses clues about the target SE class as an auxiliary input. Therefore, the system can easily exploit the information about the target SE to maximize the classification loss independently of the input mixture. We avoid this issue by pre-training the system using the fully supervised loss of Eq. (21) and only retraining the lower layers of the TSE system, which are not exposed to the target SE clues.

## VI. EXPERIMENTS

We evaluate the class label, enrollment, and SoundBeam models, and compare their performance for handling new SE classes, inactive classes, multi-target extraction, and real data.

### A. Data

To compare the effectiveness of the TSE frameworks, we created several datasets of simulated and real sound event

mixtures based on the FSD corpora, including FSD-Kaggle 2018 [52] and FSD50K [21]. In all experiments, we downsampled the sounds to 8 kHz to reduce the computational and memory costs. Table I summarizes the details of the different datasets.

*1) Single target dataset:* The single-target dataset consists of simulated mixtures. We created mixtures by mixing several SEs selected randomly from different SE classes ($N = 41$ or $N = 61$, depending on the experiment). The SEs include human sounds, object sounds, and musical instruments. We selected the 41 SE classes with the most training samples in the FSD corpora. On average, each SE class has 220 samples in the training set. We added 20 additional SE classes to generate the 61 SE class dataset. We chose the additional SE classes that were relatively well represented in the FSD corpora. However, they had fewer samples in the training data, i.e., on average, 47 samples per SE class. We generated six-second-long mixtures ($T = 6$ secs) using the Scaper toolkit [53] by inserting isolated SEs at random time positions on top of the background noise. The isolated SEs consisted of signals of 2 to 5 secs randomly selected from the FSD corpora. The background noise consisted of stationary noise from the REVERB challenge corpus [54] mixed at a signal-to-noise ratio (SNR) between 15 and 25 dB. The mixtures were composed of three SEs ($M = 3$). For the enrollment-based experiments, we randomly selected a sample from the target SE class that differs from the target sound in the mixture.

*2) Multi-target dataset:* For the multi-target TSE experiments, we created a dataset of mixtures created similarly to the single-target dataset but with three to five SEs ($M = 3 \sim 5$) per mixture.

*3) Few-shot adaptation dataset:* For the few-shot adaptation experiments, we consider the 41 classes of the single-target dataset as seen classes and the remaining 20 classes as new. We sampled $K$ enrollments from each new SE class, with $K = 1, \ldots, 15$. We used these enrollments to compute the average embedding of Eq. (24) and to generate adaptation data. We created the adaptation data by mixing one enrollment of the new SE classes with two SE training samples from the 41 seen classes, using a similar simulation procedure as above. The total number of adaptation data consisted of 3000 mixtures, covering all 20 new SE classes. Here, the adaptation data contains all 20 new SE classes because we perform the adaptation of multiple new SE classes simultaneously. However, we could also adapt the models for each new class independently. We generated six trials for the adaptation experiments by varying the random seed for the sampling enrollments. We used the test set of the single-target class datasets with 61 classes for evaluation.

*4) Real mixture dataset:* In addition to the simulated mixtures, we also used a small dataset of real mixtures that consists of recordings from the FSD50K that contained two or more labeled SE classes. The recording and mixing conditions vary considerably compared to the simulated mixtures. For example, most mixtures had two to three classes, but some contained up to nine, and the duration varied between 1 to 30 seconds. Moreover, although the recordings include multiple SE classes, there is no control over the actual temporal overlap.

Therefore, we can reasonably expect that some of the real recordings consist of partially or non-overlapping SEs.

*5) SEC dataset:* Finally, we used the original training data from both the FSD-Kaggle and the FSD50K corpora to retrain the SEC system.

### B. Experimental settings

We used the Asteroid toolkit for all experiments [55].

*1) TSE models:* The configuration of the different TSE models followed Conv-TasNet [24]. In the explanations below, we follow the notations of original Conv-TasNet [24], except for the encoded features dimension $D$. In particular, the encoder and decoder layers consisted of 1-D convolution and deconvolution layers that operated on segments of $L = 20$ taps with 50% overlap. The dimension of the encoded features was $D = 256$. MixBlock, EnrlBlock and TgtBlock in Figure 2 consisted of stacked dilated 1-D convolutional blocks with $H = 512$ channels, kernel size of $P = 3$, and $B = 256$ bottleneck channels. MixBlock and EnrlBlock consisted of a single stack ($R = 1$), while in TgtBlock the stack was repeated seven times ($R = 7$). We used $X = 8$ convolutional blocks per stack for MixBlock and TgtBlock, and $X = 4$ for EnrlBlock.

We trained all models using the Adam optimizer [56] with a learning rate of $10^{-4}$, a batch size of 8, and a maximum of 200 epochs. In all experiments, we used the models achieving the lowest cross-validation loss value.

For the few-shot adaptation experiments, we fixed all network parameters except the embedding layer and trained for 50 epochs with a learning rate of $10^{-3}$.

For retraining with the weakly supervised SEC loss, we fixed all parameters except for MixBlock (i.e., the part of the extraction NN not exposed to the auxiliary inputs). This scheme prevents the network from exploiting the embedding vectors to reduce the SEC loss artificially.

*2) USS model:* We compared the TSE models with a USS system based on Conv-TasNet with three outputs. We used a similar network configuration for the separation system to that of SoundBeam and trained it using PIT.

To realize TSE, we chose from among the outputs of the separation the one output with the highest posterior probability for the target SE class. We computed the posteriors with the SEC model described below.

*3) SEC model:* We used the publicly available pretrained audio NN (PANN) model [6] for our experiments with SEC. We retrained the model for SEC on the SEC dataset following the publicly available recipe[3]. The base model consisted of the CNN14 network trained on the AudioSet dataset [57], where we replaced the output layer to classify the 61 SE classes of our experiments. We retrained the model for 10,000 iterations using the Adam optimizer [56], with a learning rate of $10^{-4}$ and a batch size of 32.

### C. Evaluation metrics

We measured the extraction performance for active SE classes in terms of scale-invariant SDR (SI-SDR) computed

---

[3]https://github.com/qiuqiangkong/audioset_tagging_cnn

| | Model | No. classes (Training) | No. classes (Test) 41 | 61 |
|---|---|---|---|---|
| 1 | *USS (oracle selection)* | 41 | *11.1* | *9.9* |
| 2 | *USS (oracle selection)* | 61 | *10.6* | *9.8* |
| 3 | USS (SEC-based selection) | 41 | 7.7 | 6.4 |
| 4 | USS (SEC-based selection) | 61 | 7.3 | 6.2 |
| 5 | Enrl | 41 | 9.8 | 7.7 |
| 6 | Enrl | 61 | 10.1 | 7.9 |
| 7 | Class | 41 | 10.9 | - |
| 8 | Class | 61 | 10.3 | 8.3 |
| 9 | SoundBeam (Enrl) | 41 | 9.9 | 7.9 |
| 10 | SoundBeam (Enrl) | 61 | 9.9 | 7.6 |
| 11 | SoundBeam (Class) | 41 | **11.9** | - |
| 12 | SoundBeam (Class) | 61 | 11.1 | 9.2 |
| 13 | SoundBeam-adapt (Class), $K = 10$ | 41+20 | **11.9** | **9.5** |

with the BSSEval toolkit [58] and reported the SDR improvement (SDRi) compared to the mixtures. The SDRi values were obtained by averaging the values for each SE in the mixtures of the test set for the single-target extraction experiments. For the multi-target experiments, we extracted one combination of target SEs per mixture for a number of targets $J$ of 1, 2, and 3.

For inactive SE classes, we measured the attenuation relative to the mixture, $\mathcal{A}^{\mathrm{mix}}$, shown in Eq. (28) and relative to the minimum source, $\mathcal{A}^{\mathrm{src}}$, which is computed by replacing $\mathbf{y}$ with the signal of the source having minimum power, $\mathbf{x}^{min}$, in Eq. (28). The results show the average attenuation over inactive test samples (one randomly selected inactive SE class per mixture). Additionally, we evaluated the inactive class detection derived from the mixture attenuation ratio as explained in section V-B by plotting the receiver operating characteristic (ROC) curve (Recall vs. False positive rate) and measuring the area under the curve (AUC).

Finally, for the evaluation of real mixtures, we used the SEC model to predict the class labels after extraction and used the mean average precision (mAP) to measure the classification performance as it is widely used for AT [6]. mAP values were obtained by first computing the average precision (area under the recall-precision curve) per SE class and then taking the average over the classes. The mAP value serves as a proxy to measure the extraction performance when the target source signals are unavailable.

### D. Results of single-target experiments

We first compare the different TSE frameworks and a USS baseline on the single-target datasets. Table II shows the SDRi for the systems trained and tested on the single-target datasets with 41 and 61 SE classes.

The first two systems (rows 1-4) consist of USS-based approaches, which first separate the mixtures into three estimated source signals and then choose the target SE class

among the separated signals. The "oracle selection" (rows 1-2) chooses the separated signal with the highest SDR value. The "SEC-based selection" (rows 3-4) chooses the signal with the highest posterior probability of belonging to the target class according to the SEC model. The "oracle selection" represents an upper bound for USS-based approaches, while the "SEC-based selection" provides a more realistic performance level for cascade USS and SEC systems. We see that separating the sound mixtures with a high SDRi of about 10 dB is possible. However, identifying the target SE class from the separated output with an SEC is more challenging, which results in a large performance drop of more than 3 dB. We could improve the cascade system by retraining the SEC on the extracted signal, but it would not outperform the oracle selection.

"Enrl" (rows 5-6) and "Class" (rows 7-8) are the class label and enrollment-based TSE systems, respectively. "SoundBeam (Enrl)" (rows 9-10) refers to the SoundBeam model using the enrollment encoder at test time. "SoundBeam (Class)" (rows 11-12) is the same model using the class label encoder at test time, and "SoundBeam-adapt (Class)" (row 13) is the SoundBeam (Class) model adapted to the new SE classes. All TSE models outperform the "USS (SEC-based selection)" system, which demonstrates the potential of models optimized for TSE. In the remaining experiments, we focus on TSE and omit further comparison with USS-based systems.

The class label models outperform enrollment-based models because they can directly optimize the target SE embedding vectors for each SE class. "SoundBeam (Class)" (rows 11-12) performs the best and even outperforms the "USS (oracle selection)" for the test set with 41 SE classes. This result may be due to the multi-task training effect that seems to help to learn better class label embeddings.

Comparing the results with the 41- and 61-class test sets in Table II, we observe that the SDRi is consistently worse with 61 classes than with 41 classes. This is even the case for the "USS (oracle selection)" baseline (rows 1-2), indicating that it is more difficult to separate the mixtures as well. There are two possible reasons for this issue. First, the additional 20 SE classes may have sound characteristics that are more difficult to model. Second, the training data contains fewer samples from these SE classes since they are less represented in the original FSD datasets (i.e., on average 220 samples per SE class for the first 41 SE classes, but only 47 for the remaining 20 classes).

### E. Results of few-shot adaptation on new SE classes

The results of Table II confirm the ability of enrollment-based approaches trained with the dataset of 41 SE classes (rows 5 and 9) to extract new classes from the test set with 61 SE classes. However, SDRi is lower by between 0.6 dB and 1.3 dB than when using the class label approaches trained with 61 classes (rows 8 and 12).

Row 13 of Table II shows the results obtained using SoundBeam with few-shot adaptation to the 20 new SE classes as described in section V-A. The SoundBeam-adapt model was first trained with the dataset with 41 SE classes and then adapted to the 20 new SE classes using adaptation data generated from

$K = 10$ enrollments from each new SE class. To simplify the experiments, we performed adaptation simultaneously for 20 new SE classes. Since we only update the embedding matrix, this is equivalent to performing adaptation separately for each new SE class.

With the few-shot adaptation scheme, we can directly optimize the embedding vectors for the new SE classes while preserving performance on the seen SE classes. Interestingly, the SoundBeam-adapt model performs slightly better than SoundBeam trained on 61 classes from scratch. These results confirm the effectiveness of the proposed adaptation scheme even if it uses fewer samples to learn the new classes (i.e., on average, 47 samples when training from scratch compared with 10 for the few-shot adaptation).

Figure 3 plots the SDRi for the different target SE classes for the SoundBeam model with adaptation (row 13 of Table II) and without adaptation (row 12 of Table II). The first 41 classes are those in the 41 SE class dataset. "Sound-Beam(Class) 61 classes" was trained on the dataset with 61 SE classes from scratch. For the first 41 classes, the adapted model behaves identically to the model trained on the 41-class dataset (row 11 of Table II) since adaptation does not modify the embedding matrix for the seen SE classes.

We observe that, overall, "SoundBeam-Adapt" provides comparative performance to the "SoundBeam(Class) 61 classes" model for the last 20 SE classes. However, it performs slightly better for most of the first 41 SE classes, which translates to better average performance as seen in Table II. Note, this behavior may change depending on the data used and especially the type of SE classes and the amount of training samples from each class.

Figure 3 shows that SoundBeam can successfully extract the target sounds with SDRi of more than 5 dB for most target SE classes. We also confirm that extraction tends to be more challenging for the 20 new SE classes (especially for classes 53, 54, 56, 57, and 59), even when the training data included these classes.

Finally, we investigated the impact of the number of enrollment audio samples on the adaptation performance. Figure 4 shows the SDR improvement for the extraction of the new SE classes as a function of the number of enrollment samples, $K$, for the enrollment-based TSE model ("Enrl"), Sound-Beam ("SoundBeam (Enrl)"), and SoundBeam with adaptation ("SoundBeam-adapt (Class)"). "Enrl" and "SoundBeam (Enrl)" used the averaged embedding vectors of Eq. (24) to perform the extraction. Furthermore, SoundBeam with adaptation fine-tuned the averaged embeddings with the adaptation data. We performed experiments over six trials of randomly selected enrollments and here report results on only the new SE classes.

Figure 4 demonstrates that the proposed adaptation can improve performance by about 1 dB with only a few enrollments (i.e., $K \geq 3$). As a comparison, we also experimented with training the new embedding vectors from randomly initialized values, but this led to negligible SDRi.

The experiments with new SE classes reveal that Sound-Beam, with the adaptation scheme, can optimize the embedding vectors for the new SE classes, which leads to

TABLE III
RESULTS FOR INACTIVE CLASS EXTRACTION EXPERIMENT FOR TSE MODELS TRAINED WITH DATASETS HAVING 61 SE CLASSES.

| Model | IS | 41 test classes | | | | 61 test classes | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\mathcal{A}^{\mathrm{mix}} \downarrow$ | $\mathcal{A}^{\mathrm{src}} \downarrow$ | AUC $\uparrow$ | SDRi $\uparrow$ | $\mathcal{A}^{\mathrm{mix}} \downarrow$ | $\mathcal{A}^{\mathrm{src}} \downarrow$ | AUC $\uparrow$ | SDRi $\uparrow$ |
| Enrl | - | -10.1 | -0.2 | 0.62 | 10.1 | -10.3 | -0.5 | 0.60 | 7.9 |
| | ✓ | -31.5 | -21.6 | 0.79 | 9.4 | -31.7 | -21.9 | 0.74 | 6.6 |
| Class | - | -17.6 | -7.7 | 0.77 | 10.3 | -18.2 | -8.4 | 0.75 | 8.3 |
| | ✓ | -40.2 | -30.3 | 0.89 | 9.6 | -42.7 | -32.9 | **0.85** | 7.2 |
| SoundBeam (Enrl) | - | -11.1 | -1.1 | 0.65 | 9.9 | -11.1 | -1.3 | 0.62 | 7.6 |
| | ✓ | -30.2 | -20.3 | 0.79 | 9.4 | -30.3 | -20.5 | 0.75 | 6.8 |
| SoundBeam (Class) | - | -16.0 | -6.1 | 0.74 | **11.1** | -16.9 | -7.1 | 0.72 | **9.2** |
| | ✓ | **-49.6** | **-39.7** | **0.90** | 10.3 | **-52.2** | **-42.4** | 0.84 | 7.1 |

improved performance. This result demonstrates the potential of SoundBeam for continuous learning of new SE classes.

### F. Results for inactive SE classes

The next experiment investigates TSE of inactive SE classes as discussed in Section V-B. Table III compares the different TSE models trained with and without IS on the dataset with 61 SE classes and tested on the evaluation data with 41 and 61 SE classes. In this experiment, for each mixture in the test sets, we randomly chose one *inactive* SE class and computed the attenuation and AUC for that class. We report the averaged results in the table. $\mathcal{A}^{\mathrm{mix}}$ and $\mathcal{A}^{\mathrm{src}}$ are the attenuation with respect to the mixture and the lowest source. In the table, $\downarrow$ indicates that the lower the attenuation, the better the model outputs zero signals for inactive classes. The AUC measures the performance of inactive SE class detection using a simple classifier that considers a source inactive when $\mathcal{A}^{\mathrm{mix}}$ is lower than a threshold. $\uparrow$ indicates that the higher the AUC, the better the model can detect inactive classes. Figure 5 shows the ROC curves that were used to compute the AUC. Finally, we also report, in Table III, the SDRi for the extraction of the *active* SE classes in the mixture.

The results of Table III confirm that models trained without IS cannot detect inactive SE classes well and thus cannot output zero signals. In particular, the enrollment-based models perform extraction based on the similarities of the enrollments and sounds in the mixture. Consequently, if not explicitly taught to output zero signals for inactive sources, they tend to pick up the closest sounds in the mixture leading to a $\mathcal{A}^{\mathrm{src}}$ value of around 0 dB. When using IS, which adds 10 % of inactive SE classes to the training data as discussed in Section V-B, all models detect inactive SE classes better with $\mathcal{A}^{\mathrm{src}}$ values below -20 dB and AUC values above 0.75 in all cases.

Overall, the class label and SoundBeam models with IS can identify inactive SE classes well. Note that detecting inactive classes comes at the expense of lower extraction performance for the active SE classes, especially for the test set with 61 SE classes. As discussed previously, the 61 SE class test set contains classes that are more challenging to identify, and thus the model trained with IS tends to more often consider an active class as inactive. In future works, we will consider approaches to mitigate this issue by training better models for the 61 SE classes using, for example, training data
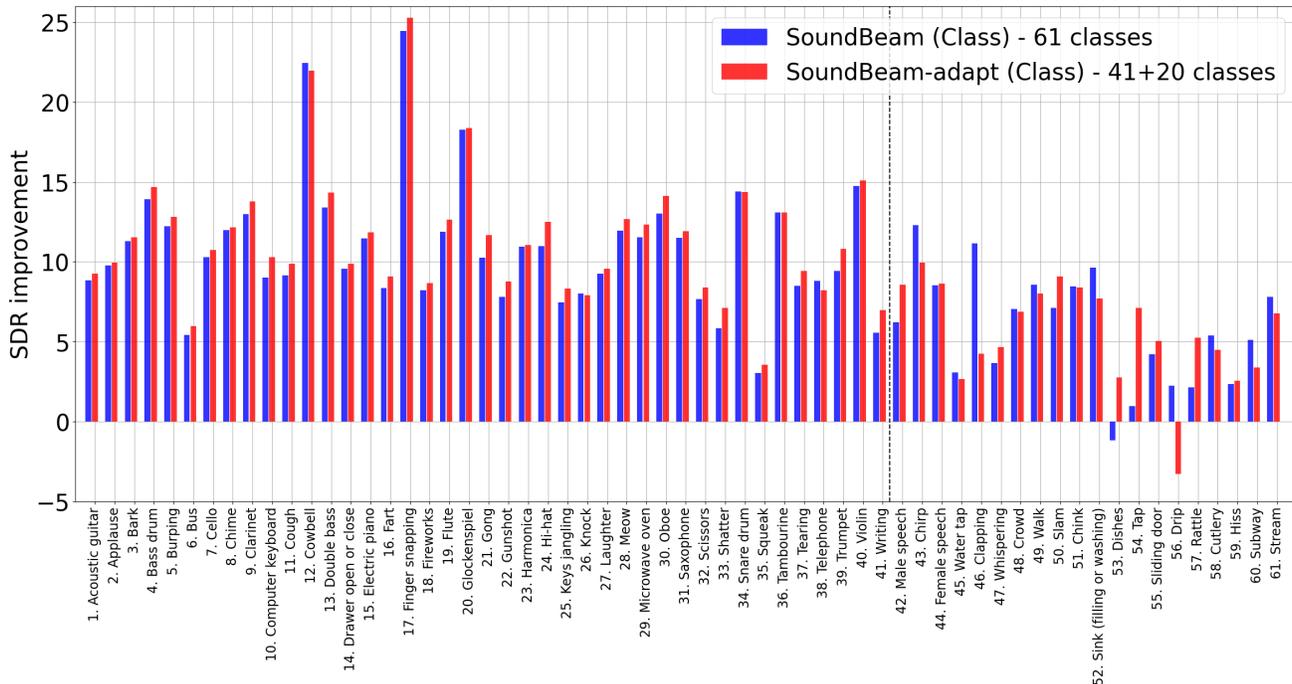
Fig. 3. SDRi [dB] for each SE class for SoundBeam models without and with few-shot adaptation, i.e. model trained with the single-target datasets with 61 classes (SoundBeam (Class)) and model trained with 41 classes then adapted to 20 new SE classes using the adaptation data (SoundBeam-adapt (Class).
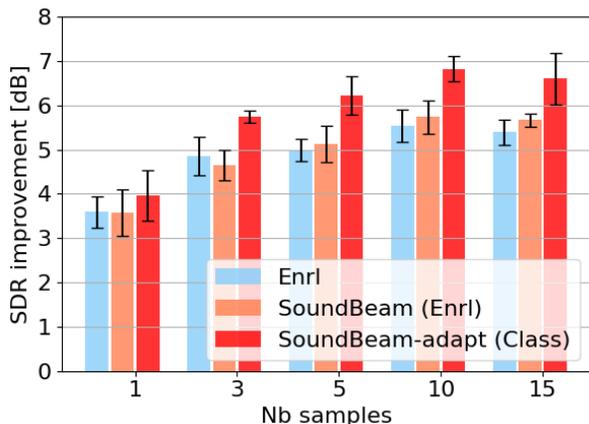


Fig. 4. SDRi as a function of the number of enrollments, $K$, for the Enrollment model, SoundBeam (Enrl) model trained with 41 SE classes, and SoundBeam-adapt (Class) model. The figure shows the mean and variance over six trials of randomly selected enrollments.



Fig. 5. ROC curves for inactive event detection using different TSE models. Results are shown for the test set with 41 classes.

augmentation for the challenging SE classes, as well as better tuning the amount of IS during training. Besides, some works dealing with target speech extraction [59], [60] have recently proposed to address the inactive speaker case by using a speaker verification module to confirm that the extracted signal was really from the target speaker and otherwise considering the target as inactive. In future works, we will explore similar verification-based ideas for the TSE problem.

*G. Results of multi-target experiments*

We tested the different models for the simultaneous extraction of multiple target SE classes as discussed in Section V-C.
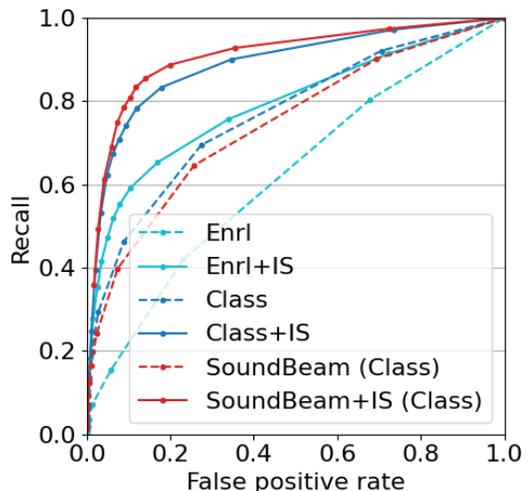
It is also possible to extract multiple-target SE classes using an iterative scheme, which extracts each class at a time. We showed in an earlier work [2] that both schemes achieved comparable performance, but the iterative scheme is more computationally expensive. Therefore, we omit here the results with the iterative scheme.

Table IV compares the SDR values for various numbers of sources in the mixtures ($M$ from 3 to 5) and targets ($J$ from 1 to 3) for the multi-target dataset described in Section VI-A2. The top three rows indicate the SDR of the mixtures, which, as expected, increases with the number of target SE classes. In the extreme case of extracting sounds from three

| Model | No. of targets $J$ | No. of sources in mixture $M$ | | |
|---|---|---|---|---|
| | | 3 | 4 | 5 |
| Mixture | 1 | -3.6 | -5.6 | -6.8 |
| | 2 | 4.0 | 0.2 | -2.0 |
| | 3 | 21.0 | 6.0 | 2.3 |
| Enrl | 1 | 3.2 (6.8) | 0.6 (6.2) | -1.0 (5.8) |
| | 2 | 7.9 (3.9) | 3.6 (3.5) | 1.4 (3.4) |
| | 3 | 21.3 (0.3) | 7.6 (1.5) | 3.9 (1.7) |
| Class | 1 | 4.4 (8.0) | 2.4 (7.9) | 0.8 (7.6) |
| | 2 | 8.3 (4.3) | **5.2 (5.0)** | 2.9 (4.9) |
| | 3 | 17.2 (-3.7) | 8.6 (2.5) | 5.4 (3.1) |
| SoundBeam (Enrl) | 1 | 2.9 (6.5) | 0.3 (5.8) | -1.3 (5.5) |
| | 2 | 7.8 (3.8) | 3.5 (3.3) | 1.2 (3.2) |
| | 3 | 21.6 (0.6) | 7.5 (1.4) | 3.8 (1.6) |
| SoundBeam (Class) | 1 | **4.9 (8.5)** | **2.8 (8.3)** | **1.1 (7.9)** |
| | 2 | **8.8 (4.8)** | **5.2 (5.0)** | **3.0 (5.0)** |
| | 3 | 18.2 (-2.8) | **8.9 (2.9)** | **5.5 (3.2)** |

| Model | 41 test classes | | 61 test classes | |
|---|---|---|---|---|
| | SDRi | mAP | SDRi | mAP |
| Enrl | 10.1 | 0.46 | 7.9 | 0.36 |
| Class | 10.3 | 0.53 | 8.3 | 0.45 |
| SoundBeam (Enrl) | 9.9 | 0.46 | 7.6 | 0.35 |
| SoundBeam (Class) | 11.1 | 0.55 | 9.2 | 0.47 |

SE classes, $J = 3$, in mixtures of three SE sounds, $M = 3$, the mixture and the reference, $\mathbf{x}^S$, differ only by the presence of background noise, and thus the SDR of the mixture attains 21 dB. It is thus naturally more challenging to improve SDR for multi-target extraction.

The results of Table IV confirm that simultaneous multi-target extraction is possible. Here again, SoundBeam outperforms the enrollment- or class label-based models.

Interestingly, we also observe that the extraction performance is not greatly affected by the number of sources in the mixture, i.e., we achieved similar SDRi values for the mixtures of 3 to 5 sources. This result confirms that TSE can operate independently of the number of sources in the mixture. Note that in our previous work [2], we also confirmed that TSE was possible when the number of sources in the mixtures was not encountered during training.

### H. Investigations with real mixtures

Finally, we perform preliminary experiments to explore the challenge of TSE when using real mixtures as discussed in Section V-D.

*1) Measuring performance with SEC:* For real recordings, we do not have access to the isolated target SE signals, $\mathbf{x}^s$. It is thus not possible to compute SDR values. Therefore, to evaluate the performance on real recordings, we propose measuring the classification performance of SEC applied to the
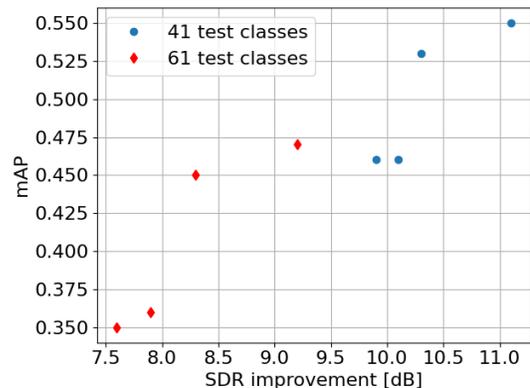


Fig. 6. mAP as a function of the SDRi for various models.

| Model | mAP |
|---|---|
| SoundBeam (Class) | 0.47 |
| + Weakly supervised retraining | 0.49 |

extracted signals. A similar approach was recently proposed for evaluating speech separation performance [61].

First, we evaluate how well we can measure TSE performance with SEC by comparing SDRi and mAP on *simulated mixtures*. Table V shows the SDRi and mAP obtained by classifying the output of TSE for the simulated single-target dataset. We should compare the mAP values with those obtained by classifying the mixture directly. When scoring the mixture against the target SE class (i.e., considering one target class as the reference), we obtain a mAP(tgt) score of 0.13. When scoring the mixture against all active SE classes in the mixture (i.e., considering all active classes as the references), we obtain a mAP(all) score of 0.39. Here, we only report mAP(tgt) when applying SEC after extraction.

We observe that the mAP values increase greatly after TSE, achieving values over 0.45 for most models. Figure 6 plots the mAP values as a function of the SDRi for the different systems of Table V. Although the number of systems is limited, the results suggest a relationship between SDRi and mAP values. This result justifies our use of mAP to evaluate extraction performance on real mixtures when it is impossible to compute the SDR values.

*2) Extraction results with real mixtures:* Table VI shows the mAP values on the real mixtures for the SoundBeam model trained with the single-target dataset with 61 SE classes. The mAP increases from 0.25 to 0.47, which indicates that SEC performance improved after extraction, suggesting that SoundBeam could extract the target.

There is a mismatch between the recording conditions of the real mixtures and the simulated training data. We experimented with retraining the model on training data consisting of real mixtures, using the weakly supervised SEC loss described in Section V-D. The last row of Table VI shows the results after

retraining, which improved mAP by 2 points. This difference is significant according to a Student's paired t-test [62] for a *p*-value of 0.059.

We should emphasize that these results are only suggestive and should be considered with precaution because good SEC does not always mean good extraction. For example, if only a portion of the sound of the target SE class is extracted, the mAP can be high, although the extraction would be imperfect.

Although they provide a likely imperfect evaluation, the results of Table VI combined with informal listening[4] indicate that SoundBeam can perform TSE on real mixtures. However, they also imply that future work is required to improve extraction performance and measure performance more accurately.

## VII. CONCLUSION

In this paper, we introduced a TSE framework, Sound-Beam, and performed extensive experiments comparing it to enrollment-based and class label-based schemes. We showed that the SoundBeam model combined the strengths of both enrollment- and class label-based schemes, which translated to better overall performance in various conditions, including inactive classes, new classes, and multi-target extraction. Furthermore, it offers the possibility of learning how to extract SE classes with few-shot adaptation. We also discussed the applicability of SoundBeam to processing recordings of real sound mixtures. Sound samples of the proposed SoundBeam are available on our demo webpage[4]

These experiments show the potential of TSE to tackle practical applications. However, there are many remaining issues. First, we focused on offline processing using computationally intensive network configurations, but many applications of TSE would require online processing with limited computational resources. The ideas presented could also be applied to a causal implementation of the models [24] or more efficient models [31], but more investigations would be required. Second, although we demonstrated promising initial results on real recordings, further research is still needed to improve extraction performance. For example, future works could include designing simulated training data that approximate better mixing conditions of real mixtures, exploiting larger datasets for semi-supervised retraining, or combining mixtures-of-mixtures and SEC-based training/adaptation strategies. Finally, we would like to extend our investigations to an even larger number of SE classes using larger datasets such as AudioSet [57].

## REFERENCES

[1] E. C. Cherry, "Some experiments on the recognition of speech, with one and with two ears," JASA, vol. 25, no. 5, pp. 975–979, 1953.

[2] 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 Proc. of Interspeech, 2020, pp. 1441–1445.

[3] E. Cakir, T. Heittola, H. Huttunen, and T. Virtanen, "Polyphonic sound event detection using multi label deep neural networks," in Proc. of IJCNN, 2015, pp. 1–7.

[4] A. Mesaros, T. Heittola, T. Virtanen, and M. D. Plumbley, "Sound event detection: A tutorial," IEEE Signal Processing Magazine, vol. 38, no. 5, pp. 67–83, 2021.

[5] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold et al., "CNN architectures for large-scale audio classification," in Proc. of ICASSP, 2017, pp. 131–135.

[6] Q. Kong, Y. Cao, T. Iqbal, Y. Wang, W. Wang, and M. D. Plumbley, "Panns: Large-scale pretrained audio neural networks for audio pattern recognition," IEEE/ACM Trans. ASLP, vol. 28, pp. 2880–2894, 2020.

[7] C. Clavel, T. Ehrette, and G. Richard, "Events detection for an audio-based surveillance system," in Proc. of ICME, 2005, pp. 1306–1309.

[8] P. Atrey, N. Maddage, and M. Kankanhalli, "Audio based event detection for multimedia surveillance," in Proc. of ICASSP, 2006, pp. 813–816.

[9] A. Mesaros, T. Heittola, A. Diment, B. Elizalde, A. Shah, E. Vincent, B. Raj, and T. Virtanen, "DCASE 2017 Challenge setup: Tasks, datasets and baseline system," in Proc. of DCASE, 2017.

[10] I. Kavalerov, S. Wisdom, H. Erdogan, B. Patton, K. Wilson, J. Le Roux, and J. R. Hershey, "Universal sound separation," in Proc. of WASPAA, 2019, pp. 175–179.

[11] Q. Kong, Y. Wang, X. Song, Y. Cao, W. Wang, and M. D. Plumbley, "Source separation with weakly labelled data: An approach to computational auditory scene analysis," in Proc. of ICASSP, 2020, pp. 101–105.

[12] B. Gfeller, D. Roblek, and M. Tagliasacchi, "One-shot conditional audio filtering of arbitrary sounds," in Proc. of ICASSP, 2021, pp. 501–505.

[13] M. Delcroix, J. B. Vázquez, T. Ochiai, K. Kinoshita, and S. Araki, "Few-shot learning of new sound classes for target sound extraction," Proc. of Interspeech, 2021.

[14] Y. Okamoto, S. Horiguchi, M. Yamamoto, K. Imoto, and Y. Kawaguchi, "Environmental sound extraction using onomatopoeia," arXiv preprint arXiv:2112.00209, 2021.

[15] K. Chen, X. Du, B. Zhu, Z. Ma, T. Berg-Kirkpatrick, and S. Dubnov, "Zero-shot audio source separation through query-based learning from weakly-labeled data," in Proc. of the AAAI, vol. 36, no. 4, 2022, pp. 4441–4449.

[16] J. H. Lee, H. Choi, and K. Lee, "Audio query-based music source separation," in Proc. of ISMIR, 2019, pp. 878–885.

[17] O. Slizovskaia, L. Kim, G. Haro, and E. Gomez, "End-to-end sound source separation conditioned on instrument labels," in Proc. of ICASSP, 2019, pp. 306–310.

[18] K. Zmolikova, M. Delcroix, K. Kinoshita, T. Ochiai, T. Nakatani, L. Burget, and J. Cernocky, "SpeakerBeam: Speaker aware neural network for target speaker extraction in speech mixtures," IEEE JSTSP, vol. 13, no. 4, pp. 800–814, 2019.

[19] P. Seetharaman, G. Wichern, S. Venkataramani, and J. Le Roux, "Class-conditional embeddings for music source separation," in Proc. of ICASSP, 2019, pp. 301–305.

[20] M. Delcroix, K. Zmolikova, K. Kinoshita, A. Ogawa, and T. Nakatani, "Single channel target speaker extraction and recognition with Speaker-Beam," in Proc. of ICASSP, 2018, pp. 5554–5558.

[21] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, "FSD50k: an open dataset of human-labeled sound events," arXiv preprint arXiv:2010.00475, 2020.

[22] M. Kolbæk, D. Yu, Z.-H. Tan, and J. Jensen, "Multitalker speech separation with utterance-level permutation invariant training of deep recurrent neural networks," IEEE/ACM Trans. ASLP, vol. 25, no. 10, pp. 1901–1913, 2017.

[23] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watanabe, "Deep clustering: Discriminative embeddings for segmentation and separation," in Proc. of ICASSP, 2016, pp. 31–35.

[24] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time–frequency magnitude masking for speech separation," IEEE/ACM Trans. ASLP, vol. 27, no. 8, pp. 1256–1266, 2019.

[25] E. Cano, D. FitzGerald, A. Liutkus, M. D. Plumbley, and F.-R. Stöter, "Musical source separation: An introduction," IEEE Signal Processing Magazine, vol. 36, no. 1, pp. 31–40, 2019.

[26] S. Uhlich, M. Porcu, F. Giron, M. Enenkl, T. Kemp, N. Takahashi, and Y. Mitsufuji, "Improving music source separation based on deep neural networks through data augmentation and network blending," in Proc. of ICASSP, 2017, pp. 261–265.

[27] Y. Luo, Z. Chen, J. R. Hershey, J. Le Roux, and N. Mesgarani, "Deep clustering and conventional networks for music separation: Stronger together," in Proc. of ICASSP, 2017, pp. 61–65.

[28] Z. Rafii, A. Liutkus, F.-R. Stöter, S. I. Mimilakis, D. FitzGerald, and B. Pardo, "An overview of lead and accompaniment separation in music," IEEE/ACM Trans. ASLP, vol. 26, no. 8, pp. 1307–1335, 2018.

[29] R. Kumar, Y. Luo, and N. Mesgarani, "Music Source Activity Detection and Separation Using Deep Attractor Network," in Proc. Interspeech, 2018, pp. 347–351.

---

[4] Sound samples can be found on our demo webpage www.kecl.ntt.co.jp/icl/signal/member/marcd/SoundBeamDemo

[30] S. Wisdom, H. Erdogan, D. P. Ellis, R. Serizel, N. Turpault, E. Fonseca, J. Salamon, P. Seetharaman, and J. R. Hershey, "What's all the FUSS about free universal sound separation data?" in Proc. of ICASSP, 2021, pp. 186–190.

[31] E. Tzinis, Z. Wang, and P. Smaragdis, "Sudo rm-rf: Efficient networks for universal audio source separation," in Proc. of MLSP, 2020, pp. 1–6.

[32] M. Olvera, E. Vincent, R. Serizel, and G. Gasso, "Foreground-background ambient sound scene separation," in Proc. of EUSIPCO, 2021, pp. 281–285.

[33] T. Heittola, A. Mesaros, T. Virtanen, and A. Eronen, "Sound event detection in multisource environments using source separation," in Machine Listening in Multisource Environments, 2011.

[34] J. F. Gemmeke, L. Vuegen, P. Karsmakers, B. Vanrumste et al., "An exemplar-based NMF approach to audio event detection," in Proc. of WASPAA, 2013, pp. 1–4.

[35] S. Cornell, G. Pepe, E. Principi, M. Pariente, M. Olvera, L. Gabrielli, and S. Squartini, "The UNIVPM-INRIA systems for the DCASE 2020 task 4," Proc. of DCASE, 2020.

[36] N. Turpault, S. Wisdom, H. Erdogan, J. Hershey, R. Serizel, E. Fonseca, P. Seetharaman, and J. Salamon, "Improving sound event detection in domestic environments using sound separation," Proc. of DCASE, 2020.

[37] Y. Huang, L. Lin, S. Ma, X. Wang, H. Liu, Y. Qian, M. Liu, and K. Ouch, "Guided multi-branch learning systems for sound event detection with sound separation," Proc. of DCASE, 2020.

[38] 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 Proc. of Interspeech, 2019, pp. 2728–2732.

[39] H. Zhao, C. Gan, A. Rouditchenko, C. Vondrick, J. McDermott, and A. Torralba, "The sound of pixels," in Proc. of ECCV, 2018, pp. 570–586.

[40] D. Samuel, A. Ganeshan, and J. Naradowsky, "Meta-learning extractors for music source separation," in Proc. of ICASSP, 2020, pp. 816–820.

[41] K. Zmolikova, M. Delcroix, K. Kinoshita, T. Higuchi, A. Ogawa, and T. Nakatani, "Speaker-aware neural network based beamformer for speaker extraction in speech mixtures," in Proc. of Interspeech, 2017, pp. 2655–2659.

[42] S. Wisdom, E. Tzinis, H. Erdogan, R. J. Weiss, K. Wilson, and J. R. Hershey, "Unsupervised sound separation using mixtures of mixtures," Proc. of NeurIPS, 2020.

[43] N. Zhang, J. Yan, and Y. Zhou, "Weakly supervised audio source separation via spectrum energy preserved wasserstein learning," Proc. of IJCAI-18, pp. 4574–4580, 7 2018.

[44] D. Stoller, S. Ewert, and S. Dixon, "Adversarial semi-supervised audio source separation applied to singing voice extraction," in Proc. of ICASSP, 2018, pp. 2391–2395.

[45] F. Pishdadian, G. Wichern, and J. Le Roux, "Finding strength in weakness: Learning to separate sounds with weak supervision," IEEE/ACM Trans. ASLP, vol. 28, pp. 2386–2399, 2020.

[46] R. Gao and K. Grauman, "Co-separating sounds of visual objects," in Proc. of ICCV, 2019, pp. 3879–3888.

[47] M. Delcroix, T. Ochiai, K. Zmolikova, K. Kinoshita, N. Tawara, T. Nakatani, and S. Araki, "Improving speaker discrimination of target speech extraction with time-domain SpeakerBeam," in Proc. of ICASSP, 2020, pp. 691–695.

[48] 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 Proc. of ICASSP, 2019, pp. 86–90.

[49] K. Vesely, S. Watanabe, K. Zmolikova, M. Karafiat, L. Burget, and J. H. Cernocky, "Sequence summarizing neural network for speaker adaptation," in Proc. of ICASSP, 2016, pp. 5315–5319.

[50] I. J. Goodfellow, M. Mirza, D. Xiao, A. Courville, and Y. Bengio, "An empirical investigation of catastrophic forgetting in gradient-based neural networks," Proc. of ICLR, 2014.

[51] M. Borsdorf, C. Xu, H. Li, and T. Schultz, "Universal Speaker Extraction in the Presence and Absence of Target Speakers for Speech of One and Two Talkers," in Proc. Interspeech, 2021, pp. 1469–1473.

[52] E. Fonseca, M. Plakal, F. Font, D. P. Ellis, X. Favory, J. Pons, and X. Serra, "General-purpose tagging of freesound audio with audioset labels: Task description, dataset, and baseline," in Proc. of DCASE, 2018.

[53] J. Salamon, D. MacConnell, M. Cartwright, P. Li, and J. P. Bello, "Scaper: A library for soundscape synthesis and augmentation," in Proc. of WASPAA, 2017, pp. 344–348.

[54] K. Kinoshita, M. Delcroix, S. Gannot, E. A. Habets, R. Haeb-Umbach, W. Kellermann, V. Leutnant, R. Maas, T. Nakatani, B. Raj et al., "A summary of the REVERB challenge: state-of-the-art and remaining challenges in reverberant speech processing research," EURASIP Journal on Advances in Signal Processing, vol. 7, 2016.

[55] M. Pariente, S. Cornell, J. Cosentino, S. Sivasankaran, E. Tzinis, J. Heitkaemper, M. Olvera, F.-R. Stöter, M. Hu, J. M. Martin-Doñas, D. Ditter, A. Frank, A. Deleforge, and E. Vincent, "Asteroid: The PyTorch-Based Audio Source Separation Toolkit for Researchers," in Proc. of Interspeech, 2020, pp. 2637–2641.

[56] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in Proc. of ICLR, 2015.

[57] J. F. Gemmeke, D. P. Ellis, D. Freedman, A. Jansen, W. Lawrence, R. C. Moore, M. Plakal, and M. Ritter, "Audio set: An ontology and human-labeled dataset for audio events," in Proc. of ICASSP, 2017, pp. 776–780.

[58] E. Vincent, R. Gribonval, and C. Fevotte, "Performance measurement in blind audio source separation," IEEE trans. ASLP, vol. 14, no. 4, pp. 1462–1469, 2006.

[59] C. Zhang, M. Yu, C. Weng, and D. Yu, "Towards robust speaker verification with target speaker enhancement," in Proc. of ICASSP'21, 2021, pp. 6693–6697.

[60] M. Delcroix, K. Kinoshita, T. Ochiai, K. Zmolikova, H. Sato, and T. Nakatani, "Listen only to me! how well can target speech extraction handle false alarms?" in Proc. of Interspeech'22, 2022.

[61] M. Maciejewski, S. Watanabe, and S. Khudanpur, "Speaker Verification-Based Evaluation of Single-Channel Speech Separation," in Proc. Interspeech, 2021, pp. 3520–3524.

[62] M. D. Smucker, J. Allan, and B. Carterette, "A comparison of statistical significance tests for information retrieval evaluation," in Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management. New York, NY, USA: Association for Computing Machinery, 2007, p. 623–632. [Online]. Available: https://doi.org/10.1145/1321440.1321528