

Tackling the Cocktail Fork Problem for Separation and Transcription of Real-World Soundtracks

Darius Petermann, Gordon Wichern, Aswin Shanmugam Subramanian, Zhong-Qiu Wang, and Jonathan Le Roux

Abstract—Emulating the human ability to solve the cocktail party problem, i.e., focus on a source of interest in a complex acoustic scene, is a long standing goal of audio source separation research. Much of this research investigates separating speech from noise, speech from speech, musical instruments from each other, or sound events from each other. In this paper, we focus on the cocktail fork problem, which takes a three-pronged approach to source separation by separating an audio mixture such as a movie soundtrack or podcast into the three broad categories of speech, music, and sound effects (SFX - understood to include ambient noise and natural sound events). We benchmark the performance of several deep learning-based source separation models on this task and evaluate them with respect to simple objective measures such as signal-to-distortion ratio (SDR) as well as objective metrics that better correlate with human perception. Furthermore, we thoroughly evaluate how source separation can influence downstream transcription tasks. First, we investigate the task of activity detection on the three sources as a way to both further improve source separation and perform transcription. We formulate the transcription tasks as speech recognition for speech and audio tagging for music and SFX. We observe that, while the use of source separation estimates improves transcription performance in comparison to the original soundtrack, performance is still sub-optimal due to artifacts introduced by the separation process. Therefore, we thoroughly investigate how remixing of the three separated source stems at various relative levels can reduce artifacts and consequently improve the transcription performance. We find that remixing music and SFX interferences at a target SNR of 17.5 dB reduces speech recognition word error rate, and similar impact from remixing is observed for tagging music and SFX content.

Index Terms—audio source separation, remixing, speech, music, sound effects, soundtrack, speech recognition, audio tagging, sound event detection

I. INTRODUCTION

OVER the last decade and especially with the recent advent of data-driven approaches, many studies have been investigating the separation of audio sources found in media content; whether addressing the separation of speech from non-speech in speech enhancement [1], [2], speech from other speech in speech separation [3]–[5], individual

D. Petermann is with the Department of Intelligent Systems Engineering, Indiana University, Bloomington, IN 47408, USA (e-mail: daripete@iu.edu).

G. Wichern and J. Le Roux are with Mitsubishi Electric Research Laboratories (MERL), Cambridge, MA 02139, USA (e-mail: {wichern,leroux}@merl.com).

A. S. Subramanian was with Mitsubishi Electric Research Laboratories (MERL), Cambridge, MA 02139, USA, and is now with Microsoft Corporation, One Microsoft Way, Redmond, WA 98052, USA (e-mail: asubra13@alumni.jh.edu).

Z.-Q. Wang is with the Language Technologies Institute, Carnegie Mellon University, Pittsburgh, PA 15213, USA (e-mail: wang.zhongqiu41@gmail.com).

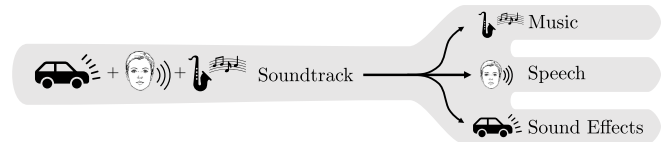


Fig. 1. Illustration of the cocktail fork problem: given a soundtrack consisting of an audio mixture of speech, music, and sound effects, the goal is to separate it into the three corresponding stems.

musical instruments in music source separation [6]–[8] or non-speech sound events (or sound effects) in universal sound separation [9]–[12], source separation finds its fit in many fields of application. Separating an audio mixture (e.g., movie soundtrack) into the three broad categories of speech, music, and sound effects (understood to include ambient noise and natural sound events) has however been left largely unexplored despite a wide domain of practical applications. A system properly trained on this task could indeed offer many potential benefits from the consumer standpoint, such as enhancing the listening experience by means of independent volume control over the sources (i.e., remixing), enabling the automatic captioning of sound events, or improving speech transcription accuracy, to name a few.

In our preliminary work [13], we formalized this new task of music-sound effects (SFX)-speech separation for real-world soundtracks as the cocktail fork problem (CFP), as illustrated in Fig. 1, and introduced a dataset specifically tailored towards this task, the Divide and Remaster (DnR) dataset, to foster research on this topic. We benchmarked multiple existing popular source separation models and proposed a multi-resolution short-time Fourier transform (STFT) architecture to better address the variety of acoustic characteristics of the three source types. We reported model performance in terms of scale-invariant signal-to-distortion ratio (SI-SDR) [14] and showcased the benefit of our proposed model, which produced SI-SDR improvements over the mixture of 11.0 dB for music, 11.2 dB for speech, and 10.8 dB for sound effects. We additionally included an analysis on system performance under different overlapping source conditions. Finally, we investigated separation performance at different sampling rates.

In this work, we explore the natural extension of the CFP towards transcription, and explore multiple techniques for integrating source separation and transcription such that the performance of both can be improved. The DnR dataset includes time-stamped annotations for all three source types, in the form of speech transcription for speech, music genre labels for music, and sound-event tags for SFX. Specifically, we investigate transcription for the CFP with the following

novel contributions:

- We present an activity detection system for the three parent classes (music, speech, SFX). By first detecting boundaries, we can then use sophisticated speech recognition and audio tagging models that expect pre-segmented chunks of audio as input.
- We showcase the benefit of integrating the output of our activity detection system with our source separation model by means of various conditioning mechanisms.
- We investigate the three classification tasks using well-established pre-trained models. For music genre and sound-event tagging, we evaluate YAMNet, a deep net that predicts 521 audio event classes from the AudioSet-YouTube corpus [15]. For speech transcription, we evaluate a conformer-based end-to-end automatic speech recognition (ASR) model provided by ESPnet [16]. The model is pre-trained on LibriSpeech [17], the same dataset DnR is based on.
- We explore the idea of source remixing, that is, the act of weighting and adding the separated sources back together, with the goal of increasing classification performance. We show that, in some cases, due to the imperfect nature of the separation, transcription can benefit from source remixing compared to using the raw separation output.
- In addition to SDR, we propose to evaluate our source separation models using two other metrics, that are arguably closer to human perception for audio quality assessment: PESQ [18] for speech and the 2f-model [19] for all three sources.

In this paper, we aim at further investigating the CFP task, not only from the source separation angle as previously achieved, but from that of transcription as well. We hope this paper serves as an important step to a system which would not only be capable of enhancing the listening experience but provide additional semantic understanding as well.

II. RELATED WORK ON THE INTERACTION BETWEEN SEPARATION AND DOWNSTREAM TASKS

The three broad categories of music, speech, and sound effects in audio signals can be found in many different types of settings, from podcasts to radio broadcasts, movies and TV-shows, their presence and overlap are ubiquitous. In more elementary scenarios, for example involving only one or two of the three classes, the task of source separation has been well investigated and has shown great promises towards various downstream tasks. For instance, separating individual instruments in musical signals has been found to be beneficial towards music transcription [20], source remixing in the context of music production [21], or even towards audio compression [22]. In the context of noisy speech, speech enhancement (or speech denoising) has shown great promise towards ASR tasks [23], [24]. More recent work proposed a joint audio-tagging and ASR system capable of fulfilling both tasks simultaneously [25]. Conversely, several works have considered how classification can improve performance in the context of general sound separation [10] or music separation [26], [27]. The interaction between voice activity detection and speech enhancement has also been explored [28], [29].

While the quality of source separation output has no doubt increased dramatically in the deep learning era, artifacts in the separated outputs remain a big problem. For speech enhancement in broadcast applications in particular, recent studies [30], [31] have shown that remixing the separated speech with some amount of the separated noise can substantially reduce artifacts and improve the listening experience compared to the noisy original. In this paper, we explore how similar remixing ideas can benefit the downstream transcription tasks of ASR, sound event detection, and music genre classification.

Particularly in complex mixtures, where many sounds may interfere with the source of interest, transcription becomes a very challenging task. One way to address this obstacle and potentially improve on the classification output would be to somehow reduce the amount of interfering sources in the mixture and consequently benefit the transcription task. For speech, the idea of using speech enhancement as a front end for recognition has been widely explored [16], [32]. However, these front-end systems are not perfect and may lead to artifacts and unwanted residual noises in the clean speech estimates, which ultimately may negatively affect the downstream ASR task. In [33], the authors propose to overcome this limitation by introducing a mechanism controlling an optimal level of noise reduction for the ASR task. An approach for learning whether to use the enhanced speech or the noisy mixture for ASR was presented in [34], where the authors ultimately found a soft combination between the enhanced signal and the noisy mixture signal performed best. A related study [35] also demonstrated the benefit of remixing the enhanced speech and the noisy mixture for ASR. In this work, we not only evaluate remixing for ASR in the presence of difficult music and sound effect background signals, but also study the benefit of remixing separated signals for sound event tagging and music genre recognition.

In [36], the authors evaluated source separation as a pre-processing step to improve the sound-event detection (SED) task by first breaking down mixtures into their constituent sounds. SED was then applied by combining the separated sources and the input mixture at different stages in the architecture. Although the separator is not trained jointly in this case (and consequently no active remixing is performed during training), one could argue that the remixing may be done implicitly within the SED network. Source separation in music applications often focuses on remixing separated musical stems [37]. Separation has also been widely used for music transcription (i.e., scoring), either as a pre-processing front-end [38], [39] or within a joint approach [20], [40]. In [41], the authors extensively explore the idea of source separation specifically applied to choir ensemble mixtures, allowing for a set of potential downstream applications such as F_0 contour analysis, synthesis, transposition, unison analysis, as well as singing group remixing. In [42], the authors investigate existing source separation algorithms and their perceptual impact on songs given various remixing scenarios. The claim suggests that existing separation approaches may suffer from imperfect separation, resulting in perceptible artifacts on individual source estimates, which can then jeopardize downstream task performance. To the best of our knowledge, there has not

TABLE I

SOME EXAMPLES OF EVENTS FOUND IN THE DNR METADATA AND HOW WE FORMULATE THEM IN THE CONTEXT OF THE CFP TRANSCRIPTION. “# CLASSES” INDICATES THE VOCABULARY SIZE FOR EACH TRANSCRIPTION TASK (E.G., THE NUMBER OF TOKENS FOR ASR IS 5,000).

Task	# classes	⟨ tags, ⟩	start_time,	end_time ⟩
Music	11	⟨ rock	2.3 s	19.5 s ⟩
SFX-Fg.	85	⟨ dog,bark,animal	10.1 s	13.8 s ⟩
SFX-Bg.	35	⟨ rain,thunder	1.0 s	30.6 s ⟩
Speech	5,000	⟨ [‘the man walks ...’]	5.9 s	24.2 s ⟩

been any work investigating the task of music source remixing specifically targeted towards downstream labeling tasks such as genre recognition in a manner similar to what we explore in this work.

III. METHODS

A. Problem Setup

In this work, we assume that we observe a single-channel mixture $x \in \mathbb{R}^T$ composed of three submixes:

$$x = y^{(s)} + y^{(m)} + y^{(e)}, \quad (1)$$

where $y^{(s)}$ is the submix containing all speech signals, $y^{(m)}$ that of all music signals, and $y^{(e)}$ that of all sound effects. We use the term sound effects (SFX) to broadly cover all sources not categorized as speech or music, and choose it over alternatives such as sound events or noise, as the term is especially relevant to our target application where x is a soundtrack. We here define the cocktail fork problem as that of recovering, from the audio soundtrack x , its music, speech, and sound effect submixes, as opposed to extracting individual musical instruments, speakers, or sound effects.

Additionally, we consider the case where the submixes have associated collections $l_{1:N_s}^{(s)}$, $l_{1:N_m}^{(m)}$, $l_{1:N_e}^{(e)}$ of metadata labels describing the content for each source type. Specifically, as illustrated by the examples in Table I, for speech we consider the speech recognition task where label $l_{i_s}^{(s)}$ represents the transcription of the i_s -th utterance and associated time boundaries, with index $i_s \in \{1, \dots, N_s\}$, where N_s is the number of utterances in $y^{(s)}$. Similarly, $l_{i_m}^{(m)}$ represents a music genre label and associated time boundaries for the i_m -th music excerpt. For sound effects, $l_{i_e}^{(e)}$ represents a list of audio tags describing the i_e -th sound event, along with the associated time boundaries; these tags are further split into two sub-categories of foreground events (SFX-Fg) such as “dog barking” and background events (SFX-Bg) such as “traffic noise,” which are treated separately for transcription.

In this work, we focus both on separation applications where the goal is to recover estimates $\hat{y}^{(s)}$, $\hat{y}^{(m)}$, $\hat{y}^{(e)}$ of the submixes, and on transcription applications where the goal is to estimate metadata label collections $\hat{l}_{1:N_s}^{(s)}$, $\hat{l}_{1:N_m}^{(m)}$, $\hat{l}_{1:N_e}^{(e)}$ given a mixture x and its associated source separation outputs $\hat{y}^{(s)}$, $\hat{y}^{(m)}$, $\hat{y}^{(e)}$.

In the same vein as some of the remixing approaches presented in Section II, we explore the idea of source remixing towards the three transcription downstream tasks. Our working hypothesis is that, while source separation likely helps the downstream networks to reach a better performance on their

respective tasks, it is imperfect, with the presence of added interferences and artifacts, and adding back a *down-scaled* version of the predicted constituents in each of the target sources (e.g., music and SFX signals in the speech signal for ASR) prior to performing the transcription may help improve performance. More formally, this can be formulated as follows:

$$\tilde{y}^{(i)} = \tau_i^{(s)} \hat{y}^{(s)} + \tau_i^{(m)} \hat{y}^{(m)} + \tau_i^{(e)} \hat{y}^{(e)}, \quad i \in \{s, m, e\}, \quad (2)$$

where $\tilde{y}^{(s)}$, $\tilde{y}^{(m)}$, and $\tilde{y}^{(e)}$ denote the remixed separated sources, and $\tau_i^{(j)}$ denotes the time invariant gain applied to the separated estimated of source j to obtain the remixed source i , e.g., $\tau_m^{(s)}$ is the gain applied to the separated speech estimate $\hat{y}^{(s)}$ when remixing for music transcription. Taking the example of speech transcription, the gain $\tau_s^{(s)}$ applied to the speech estimate $\hat{y}^{(s)}$ is set to always remain at unit level, while the gains $\tau_s^{(m)}$ and $\tau_s^{(e)}$ of the interfering sources $\hat{y}^{(m)}$ and $\hat{y}^{(e)}$, respectively, are adjusted (either individually or jointly) to match a target SNR. In the individual case, the gains are set to

$$\tau_s^{(j)} = \frac{\|\hat{y}^{(s)}\|_2}{\|\hat{y}^{(j)}\|_2} 10^{-\text{snr}_s^{(j)}/20}, \quad j \in \{m, e\}, \quad (3)$$

where $\text{snr}_s^{(m)}$ is the desired SNR of the rescaled speech signal with respect to the rescaled music signal, and $\text{snr}_s^{(e)}$ is defined similarly with respect to the rescaled sound effect signal. Alternatively, we can adjust the gain to reach a desired SNR $\text{snr}_s^{(m+e)}$ of the speech signal with respect to the sum of the two interfering source signals $\hat{y}^{(m+e)} = \hat{y}^{(m)} + \hat{y}^{(e)}$, i.e.,

$$\tau_s^{(m+e)} = \frac{\|\hat{y}^{(s)}\|_2}{\|\hat{y}^{(m+e)}\|_2} 10^{-\text{snr}_s^{(m+e)}/20}. \quad (4)$$

Note that, in all cases, we assume that the power of each separated source signal $\hat{y}^{(j)}$ is roughly equal to the power of the corresponding ground-truth source signal $y^{(j)}$, which holds in practice if the separation quality is reasonable. An overview of our proposed approach is shown in Fig. 2, and we now describe all stages in detail.

B. Source Separation with MRX

In our preliminary work [13], we found that time-frequency (TF) separation models generally worked well, and we observed additional benefit from jointly evaluating multiple TF resolutions to better handle the diverse acoustic characteristics present in mixtures of speech, music, and sound effects. Therefore, in this work, we use the multi-resolution crossnet (MRX) introduced in [13] and shown in Fig. 3 as our main network architecture. MRX takes a time-domain input mixture x and encodes it into I complex spectrograms $X_{W_i} = \text{STFT}_{W_i}(x)$, $i \in \{1, \dots, I\}$ with different STFT resolutions, where W_i denotes the i -th window length in milliseconds. Fig. 3 shows an example with $I = 3$ and $\{W_1, W_2, W_3\} = \{32, 64, 256\}$.

We use the same hop size (e.g., 8 ms in the example of Fig. 3) for all resolutions, so they remain synchronized in time, and N denotes the number of STFT frames for all resolutions. In practice, we set the window size in samples to the nearest power of 2, and the number of unique frequency bins is denoted as F_{W_i} . Each resolution is then passed to a

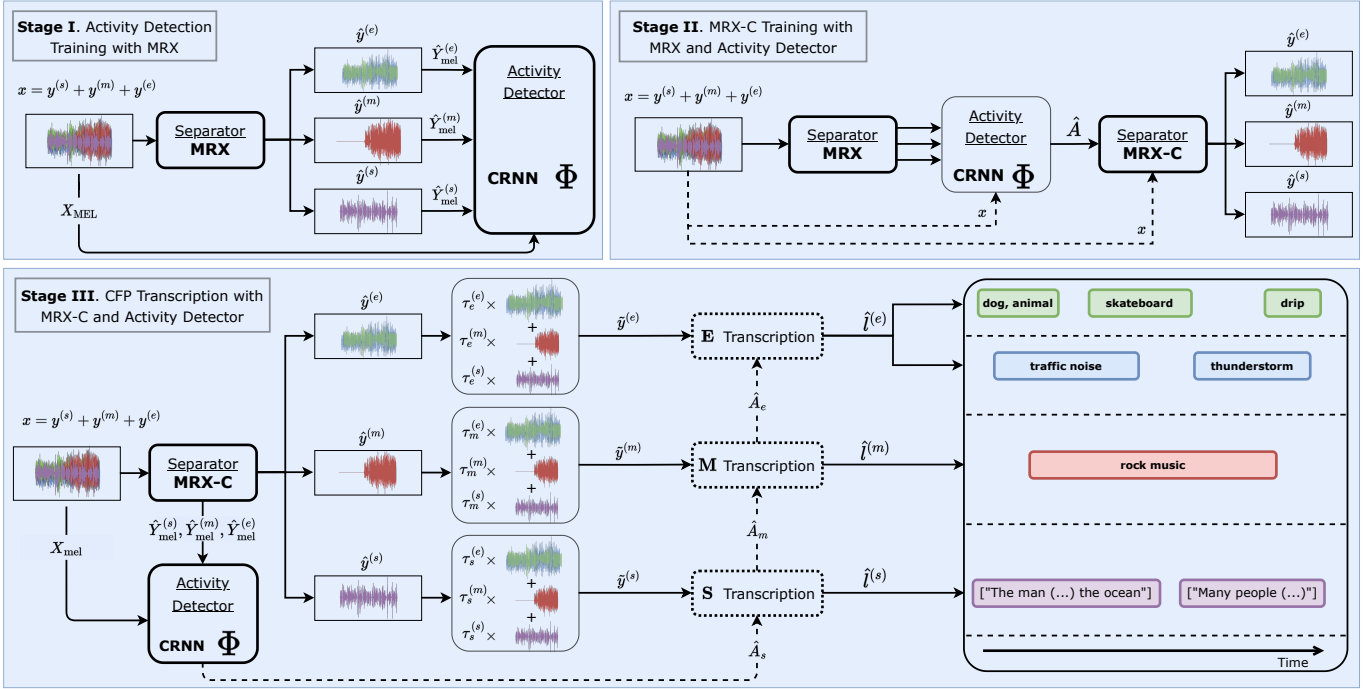


Fig. 2. Overall separation, activity detection, and transcription pipeline. In **Stage I**, the activity detector Φ_{EI} is trained to provide frame-level activity labels for the three CFP classes (music, speech, SFX). In **Stage II**, MRX-C, a label-informed version of MRX, is trained using the activity label provided by Φ_{EI} . Lastly, in **Stage III**, we proceed to total transcription for a given CFP mixture by first estimating its sources and frame-level labels. The transcription for the three sources is then achieved on the source estimates and using the frame labels.

fully connected block to convert the magnitude spectrograms of dimension $N \times F_{W_i}$ into a consistent dimension of 512 across the resolution branches. This allows us to average them together prior to the bidirectional long short-term memory (BLSTM) stacks, whose outputs are averaged once again. MRX was inspired by the Cross-Unmix (XUMX) architecture proposed in [43]. However, in our case, the input averaging is intended to allow the network to efficiently combine inputs with multiple resolutions.

The average inputs and outputs of the BLSTM stacks are concatenated and decoded back into magnitude soft masks $\hat{M}_{W_i}^{(j)}$, one for each of the three sources $j \in \{s, m, e\}$ and each of the I original input resolutions W_i . The decoder consists of two stacks of fully-connected layers, each followed by batch normalization (BN) and rectified linear units (ReLU). For a given source j , each magnitude mask $\hat{M}_{W_i}^{(j)}$ is multiplied element-wise with the original complex mixture spectrogram X_{W_i} for the corresponding resolution, a corresponding time-domain signal $\hat{y}_{W_i}^{(j)}$ is obtained via inverse STFT, and the estimated time-domain signal $\hat{y}^{(j)}$ is obtained by summing the time-domain signals at each resolution:

$$\hat{y}^{(j)} = \sum_{i=1}^I \hat{y}_{W_i}^{(j)} = \sum_{i=1}^I \text{iSTFT}(\hat{M}_{W_i}^{(j)} \odot X_{W_i}). \quad (5)$$

For the cocktail fork problem, the network has to estimate a total of $3I$ masks. Since ReLU is used as the final mask decoder nonlinearity, the network can freely learn weights for each resolution that best reconstruct the time-domain signal. We use the SI-SDR loss function [14], [44] between the estimated signal $\hat{y}^{(j)}$ and the ground-truth signal $y^{(j)}$ for $j \in \{s, m, e\}$.

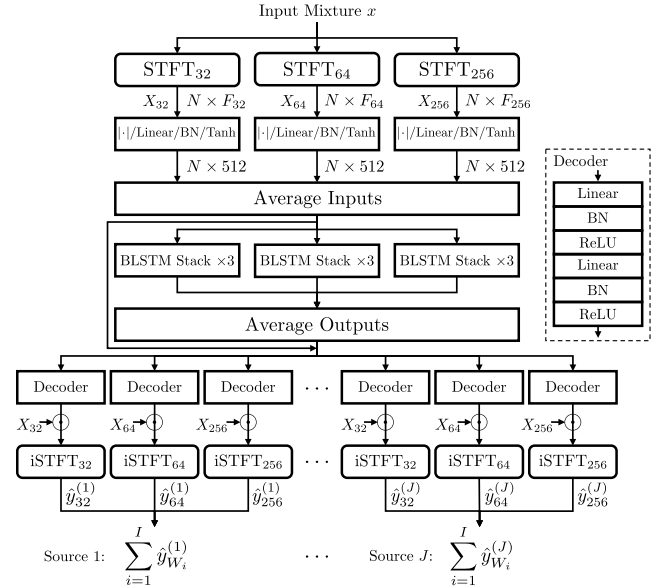


Fig. 3. Multi-resolution CrossNet (MRX) architecture.

We now investigate how source separation and activity detection may benefit from each other.

C. Stage I: Using source separation to improve activity detection

As a first step towards transcribing the separated speech, music, and sound effects stems, we focus on detecting the temporal regions where each source is active. Only these detected temporal regions for each source are then fed to the appropriate downstream classifiers, namely, ASR for speech, genre recognition for music, and audio tagging for sound

effects. End-to-end (E2E) ASR systems are typically sensitive to extended silent portions occurring between utterances as well as extensive utterance length, and may not work optimally when fed unsegmented signals. Similarly, for music and sound effects, we can pool over temporal regions to obtain more accurate tagging performance.

Formally, we define $A = (a_{j,n})_{j,n} \in \{0,1\}^{3 \times N}$ as the ground-truth activity labels, with elements $a_{j,n} = 1$ indicating that source $j \in \{s, m, e\}$ is active at frame n , and $a_{j,n} = 0$ that it is inactive. Note that the ground-truth activity labels are defined by the presence of an active excerpt at frame n as defined in the metadata, rather than its energy, so pauses in a speech utterance or a song may still be considered active. Our goal is to obtain estimated activity labels $\hat{A} = (\hat{a}_{j,n})_{j,n} \in \mathbb{R}^{3 \times N}$, using a neural network Φ applied to an input representation R derived from the mixture of interest:

$$\hat{A} = \Phi(R). \quad (6)$$

The network is trained to minimize the binary cross-entropy between A and \hat{A} over the training set:

$$\begin{aligned} \mathcal{L}(A, \hat{A}) &= \sum_{j,n} \text{BCE}(a_{j,n}, \hat{a}_{j,n}) \\ &= \sum_{j,n} \left(-a_{j,n} \log(\hat{a}_{j,n}) - (1 - a_{j,n}) \log(1 - \hat{a}_{j,n}) \right). \quad (7) \end{aligned}$$

At inference time, we apply a median filter over time and threshold \hat{A} to determine boundaries. We denote as \hat{A}_j for the estimated activities for source j .

For the core architecture of the neural network Φ , we use a convolutional recurrent neural network (CRNN), as such architectures have proven to be highly effective in the context of sound event detection [11], [45], [46], with mel spectrograms as the input representation. The estimated activity labels $\hat{a}_{j,n}$ are in the range $[0, 1]$ as they are the outputs of a sigmoid activation function. Our baseline Φ_B simply takes the mel spectrogram X_{mel} of the mixture as input, and we denote its output as $\hat{A}^{(\text{B,mix})} = \Phi_B(X_{\text{mel}})$.

In [36], multiple source-separation-based algorithms are proposed to improve sound event detection. We explore such approaches for the CFP as shown in Stage I of Fig. 2. In the first approach, **early integration** (EI), the input to a CRNN Φ_{EI} is formed by stacking the mel spectrograms of the mixture and estimated sources along the channel dimension, i.e.,

$$\hat{A}^{(\text{EI})} = \Phi_{\text{EI}}(\text{stack}(X_{\text{mel}}, \hat{Y}_{\text{mel}}^{(s)}, \hat{Y}_{\text{mel}}^{(m)}, \hat{Y}_{\text{mel}}^{(e)})). \quad (8)$$

The input is thus of shape $(4 \times N \times F)$ where F is the number of mel bands, and N the number of time frames. In **middle integration** (MI), X_{mel} and the $\hat{Y}_{\text{mel}}^{(j)}$ are each individually input into the CNN block of a CRNN Φ_{MI} , and their outputs are stacked before being fed to the RNN block. Both CRNNs Φ_{EI} and Φ_{MI} are trained using Eq. (7), with the source separation network kept frozen. Finally, in **late integration** (LI), the baseline CRNN Φ_B is used to process independently the mel spectrograms X_{mel} of the mixture and $\hat{Y}_{\text{mel}}^{(j)}$ of each estimated source to obtain the corresponding output probabilities $\hat{A}^{(\text{B,mix})}$ and $\hat{A}^{(\text{B},j)} = \Phi_B(\hat{Y}_{\text{mel}}^{(j)})$ (for example, $\hat{A}^{(\text{B},s)} \in \mathbb{R}^{3 \times N}$ denotes the estimated activity labels of

all three sources within the separated speech estimate). These probabilities are then combined to obtain the late integration estimates $\hat{A}^{(\text{LI})}$ as:

$$\hat{A}^{(\text{LI})} = \frac{1}{2} \hat{A}^{(\text{B,mix})} + \frac{1}{2} \left(\frac{\hat{A}^{(\text{B},s)} + \hat{A}^{(\text{B},m)} + \hat{A}^{(\text{B},e)}}{3} \right). \quad (9)$$

In this work, all three fusion approaches will be compared in Section V. We found that the early integration approach was the one that worked best for the CFP activity detection. We refer the interested reader to the original work on the topic [36] for further technical details.

D. Stage II: MRX-C – Using activity detection to improve source separation

As shown in Stage II of Fig. 2, we also explore how activity detection can benefit source separation, by conditioning our MRX separation models on the activity detection output. In our conditioning approach, which we call MRX-C, we concatenate class probabilities \hat{A} obtained at the output of one of the configurations of Section III-C with the MRX input mixture STFT at each resolution shown in Fig. 3. Beside concatenation, we also experimented with a FiLM [47] conditioning approach, which consisted in using a FiLM layer placed after each of the MRX STFT operations in Fig. 3. This approach led to very similar results to the concatenation counterpart in preliminary experiments, we thus opted to use the concatenation for all further experiments. Because the frame rate of the CRNN outputs may be lower than the STFT frame rate of MRX due to temporal pooling operations, we use nearest neighbor upsampling to match the frame rates.

We further explore the impact of using the ground-truth class labels during the training and inference stages as the upper bound oracle performance. We also explored running two consecutive iterations, or passes, of Stage I and Stage II of Fig. 2, which we denote MRX-C_{2p}, where the suffix “2p” refers to “2 passes”.

As a comparison to the approaches combining activity detection and separation explored in this section, we also consider a multi-task learning approach that jointly does activity detection and source separation [27], [29], referred to as MRX-MTL. Here, an additional decoder layer with three sigmoid outputs is added in parallel to the separation decoders in Fig. 3 to estimate activity detection. A time average pooling layer with a factor of eight is applied at the input of this decoder to be consistent with the CRNN output resolution, and a weighted binary cross-entropy loss is added to the separation loss for training.

E. Stage III: CFP Transcription

In this section, we focus on downstream transcription tasks using the remixed source separation outputs and source activity boundaries. Specifically, we investigate audio tagging for music (in the form of music genre recognition) and sound effects, and ASR for speech.

1) *Audio Tagging*: To transcribe music and sound effects, we make use of the powerful pre-trained audio tagging model YAMNet¹, which predicts 521 audio event tags, including multiple music genres and sound effects classes. YAMNet is based on the MobileNet convolutional architecture [48], and has been trained on the AudioSet Ontology [49], a human-labeled corpus derived from short Youtube audio segments. YAMNet operates on 960 ms frames with 50% overlap, and outputs a vector of 521 class activity probabilities for each frame. The output layer uses sigmoid activation functions, so multiple classes (i.e., tags) can be active for each frame. In practice, we input the entire soundtrack (either the mixture, separated sources, or remixes) into YAMNet, and then average the class probabilities over segments estimated by the activity detector. Example annotations from a soundtrack are shown in Table I.

For music, we typically do not expect there to be multiple pieces of music playing simultaneously, so we limit the estimated tag to a single one describing the predominant music genre. As we are considering music genre classification, we only take into account the YAMNET outputs corresponding to music genre classes.

In real-world soundtracks, sound effects typically serve two main purposes: *background events* that usually entail longer and lower amplitude sounds that help set the scene (e.g., rain, traffic), and *foreground events* which are shorter and louder to help tell the story (e.g., gun shot, footsteps). For this reason, we further sub-divide the sound effects audio tagging task into foreground (SFX-Fg) and background (SFX-Bg) sub-tasks as illustrated in Table I. Unlike for music genre, we allow sound events to be labeled with multiple tags from the AudioSet ontology (e.g., dog, bark, animal). Furthermore, sound events from both SFX-Bg and SFX-Fg can overlap in time as shown in the example of Table I. While our goal is to separate and annotate a single sound effects stem, because of the widely different characteristics of SFX-Fg and SFX-Bg sound events, we consider them separately when creating synthetic mixtures in Section IV-A and when evaluating audio tagging performance in Section V-B. However, for evaluating source separation and activity detection performance in Section V-A, we consider SFX-Fg and SFX-Bg jointly as a single sound effects class.

2) *Speech Transcription*: We evaluate ASR using a model pre-trained on the LibriSpeech corpus. Specifically, we use a state-of-the-art end-to-end (E2E) model implemented in ESPnet [16] that is based on the Conformer [50] architecture and uses HuBERT [51] input features. In general, a soundtrack contains multiple speech utterances interspersed among non-speech regions, which may cause difficulty when feeding the unsegmented soundtrack (be it the mixture, separated sources, or remixes) into an ASR model. Therefore, we evaluate the E2E ASR model using the approach shown in Fig. 4, where the activity detector Φ_{EI} first retrieves individual utterances which can then be individually input to the E2E ASR model. The individual utterance hypotheses from the E2E ASR model

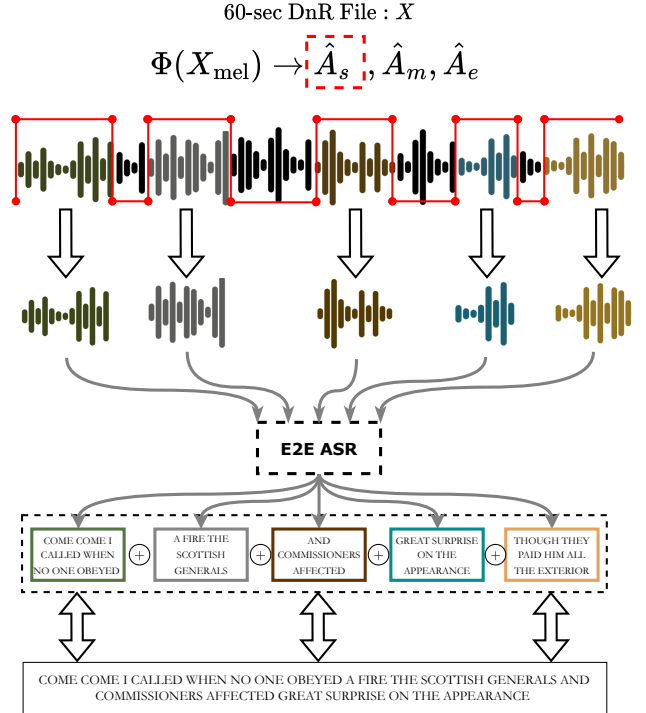


Fig. 4. Illustration of our ASR evaluation pipeline using ESPnet. We first segment the full 60-second utterance into smaller chunks using Φ 's output. Each of the resulting sub-utterances are then individually fed to ESPnet for decoding. Finally their hypotheses are concatenated and compared against the full-utterance reference.

then concatenated and evaluated against the full-file reference.

IV. EXPERIMENTAL SETUP

A. Divide and Remaster (DnR) dataset

The DnR dataset was introduced in our preliminary work [13] with the goal of creating synthetic monophonic soundtracks for training and evaluating source separation algorithms. The dataset is publicly available along with the data creation scripts², and further details on the creation process are available in [13]. DnR is created by mixing speech from LibriSpeech [17], music from the Free Music Archive (FMA) [52], and sound effects from the Freesound Dataset 50k (FSD50K) [53]. All of the DnR building blocks contain audio at sampling rates of 44.1 kHz or greater (by default LibriSpeech is available at 16 kHz, but the original 44.1 kHz mp3 files are available), so we use 44.1 kHz as the default sampling rate for DnR. This enables real-world listening applications, and can easily be downsampled for transcription where high bandwidth is unnecessary. For FSD50k, we manually classify each of the 200 class labels into one of 3 groups: foreground sounds (e.g., dog bark), background sounds (e.g., traffic noise), and speech/musical instruments (e.g., guitar, speech). Speech and musical instrument clips are filtered out to avoid confusion with our speech and music datasets, and we use different mixing rules for foreground and background events.

¹<https://github.com/tensorflow/models/tree/master/research/audioset/yamnet>

²<https://cocktail-fork.github.io>

Many synthetic mixing pipelines for source separation such as wsj0-2mix [3] or speech/music/noise separation [54] create fully-overlapped mixtures, with arbitrarily selected source levels. However, models trained on these datasets may not robustly transfer to real-world situations without full-overlap [55]. Therefore, we took great care in DnR to make mixtures that are sufficiently realistic in terms of source overlap and relative amplitude level between the three classes. In order to ensure that a mixture could contain multiple full speech utterances and feature a sufficient number of onsets and offsets between the different classes, we decided to make each mixture 60 seconds long. We do not allow within-class overlap between speech and music clips, i.e., two music files (or two speech files) will not overlap, but foreground and background sound effects can overlap. This mixing procedure leads to, approximately, 55% of the DnR test set frames having speech, music, and sound effects active, 32% containing two of the three, 10% containing one of the three, and 3% silent frames.

Regarding the relative amplitude levels across the three classes, after analyzing studies such as [56] and informal mixing rules from industries such as motion pictures, video games, and podcasting, we follow an approach where speech is generally found at the forefront of the mix, followed by foreground sound effects, then music, and finally background sound effects. DnR contains 3,406 60-second mixtures for training, 487 for validation, and 973 for testing.

Along with the audio, DnR also contains transcription metadata in a format similar to the example of Table I. The start and end times correspond to where in the 60 s mixture the original clip was inserted. The gain applied to the clip is also available from the metadata, but is not shown in the example of Table I. For speech utterances, we use the unaltered transcription of the sentence from the LibriSpeech metadata. For music, we list all genres from the FMA annotations, but we use only the top-level genre (corresponding to 16 commonly used labels such as “jazz” or “rock”) in our genre recognition experiments. For sound effects, we list all tags from the FSD50K metadata, and also note whether the clip is used as a foreground or background event.

B. Source Separation

1) *XUMX and MRX models*: We consider single-resolution XUMX baselines with various STFT resolutions. We opt to cover a wide range of window lengths $W \in \{32, 64, 128, 256\}$ to assess the impact of resolution on performance. For our proposed MRX model, we use three STFT resolutions of 32, 64, and 256 ms, which we found to work best on the validation set. We use $XUMX_W$ to denote a model with a W ms window. We set the hop size to a quarter of the window size. For the MRX model, we determine hop size based on the shortest window. To parse the contributions of the multi-resolution and multi-decoder features of MRX, we also evaluate an architecture adding MRX’s multi-decoder to the best single-resolution model ($XUMX_{64}$), referred to as $XUMX_{64, \text{multi-dec}}$. This results in an architecture of the same size (i.e., same number of parameters) as our proposed MRX model. In all architectures, each BLSTM layer has 256 hidden units and

input/output dimension of 512, and the hidden layer in the decoder has dimension 512. For all MRX models, we use a sampling rate of 44.1 kHz.

2) *Other benchmarks*: We also evaluate our own implementations of Conv-TasNet [44] and a temporal convolution network (TCN) with mask inference (MaskTCN). The sampling rate is again 44.1 kHz for all models. MaskTCN uses an identical TCN to the one used internally by Conv-TasNet, but the learned encoder and decoder are replaced with STFT and iSTFT operations. For MaskTCN, we use an STFT window/hop of 64/16 ms, and for the learned encoder/decoder of Conv-TasNet, we use 500 filters with a window size of 80 samples and a stride of 40. All TCN parameters in both Conv-TasNet and MaskTCN follow the best configuration of [44]. Additionally, we evaluate Open-Unmix (UMX) [7], the predecessor to XUMX, by training a separate model for each source, but without the parallel branches and averaging operations introduced by XUMX. The Conv-TasNet, UMX, XUMX, and MRX models all use SI-SDR [14], [44] as loss function, while MaskTCN uses the waveform domain L_1 loss.

All models are trained on 9 s chunks, except MaskTCN, trained on 6 s chunks, and Conv-TasNet, trained on 2 s chunks; we found these values to lead to best performance under our GPU memory constraints. All models are trained for 300 epochs using ADAM. The learning rate is initialized to 10^{-3} , and halved if the validation loss is not improved over 3 epochs.

C. Activity Detection

The activity detection CRNN models described in Section III-C follow the DCASE 2020 Task 4 baseline architecture and its extensions as described in [36], based on the publicly available implementation³. We use the ADAM optimizer with a learning rate of 10^{-4} , $\beta_1 = 0.9$, and $\beta_2 = 0.999$. We train the systems for 80 epochs and select the weights returning the lowest loss score on the validation set. Prior to training, the data is preprocessed in the following manner. First, we convert all the 60-second long input mixtures into mel spectrograms with 64 bands, using an FFT size of 2048 and hop-length of 512. During training, we randomly sample 10-second long chunks, which translate to 864 temporal frames, from the resulting features. Note that due to the CRNN pooling operations, the number of frames is reduced by a factor of 8 in contrast to the initial input temporal dimension, reducing the number of temporal frames from 864 at the input down to 108 at the output. The activity predictions are obtained by applying a median filter of size 5 over time and thresholding \hat{A} at 0.5.

D. MRX-C models

The conditioned MRX-C models described in Section III-D follow the same training pipeline as our MRX model described above. For the conditioning portion, the activity labels are obtained for the entirety of the DnR dataset prior to training by using our best performing activity detection model (later described in Section V-A), the early-integration model Φ_{EI} .

³https://github.com/turpaultn/dcase20_task4

During training and inference, the labels are upsampled to the length of the input mixture spectrogram using the nearest neighbor algorithm.

The multi-task learning model MRX-MTL is optimized on a compound objective function consisting of the SI-SDR and binary cross entropy losses, with their respective weight being 1.0 and 10.0.

E. Transcription

1) *Music Genre Recognition*: As YAMNet was originally trained on 16 kHz audio data, all audio involved in our YAMNet experiments is first resampled from 44.1 kHz down to 16 kHz. The audio signal goes through a set of transforms prior to network input. First the magnitude spectrogram is obtained through the STFT transform using a window size of 25ms and hop size of 10ms, the spectrogram is then mapped to 64 mel bins covering a 125-7500Hz range in order to obtain its mel transform. Lastly, a stabilized log operation is applied to it. YAMNet uses 960 ms frames overlapped at 50%. Since the DnR mixtures are 60 s long, we obtain 124 YAMNet frames for each DnR mixture.

One challenge encountered when using YAMNet for the music classification task was the mapping from the 16 top-level genre labels used in the FMA dataset, to one of the 521 sound-event classes from the Audioset ontology used by YAMNet. Ten of the genres had straightforward mapping between FMA and Audioset, namely Pop, Rock, Soul-RnB, Jazz, Country, Electronic, Blues, Hip-hop, Folk, and Classical. The “International” genre used by FMA could be mapped to multiple Audioset genres (e.g., “Music of Asia”, “Music of Africa”), however, this one-to-many mapping provided very poor performance. The YAMNet “Vocal Music” category consistently yielded high confidence when fed music labeled as “International” by FMA, as these audio files typically contained singing with little-to-no instrumental parts. Therefore, we mapped “International” to “Vocal Music”. There are 5 FMA genres that we excluded from the classification task as they were overly broad and had no appropriate Audioset genre counterpart (Experimental, Easy Listening, Instrumental, Spoken, and Old-time/historic). Any clips containing these 5 genres were still used for evaluation of music separation and activity detection, but they were ignored when evaluating music genre recognition performance using YAMNet.

2) *Sound Event Detection*: To filter music and speech clips from FSD50K during the DnR data creation process, we removed a clip if the majority of its tags were related to speech or music. The size of the FSD50K tag vocabulary after this filtering was reduced from 200 to 146. Furthermore, for YAMNet compatibility, we had to exclude 26 additional tags for the followings reasons:

- The class is an actual parent class in AudioSet, therefore not covered in YAMNet output (e.g., “Domestic_sounds_and_home_sounds”).
- The class does not have any examples in AudioSet, therefore it is not covered in YAMNet output (e.g., “Gull_and_seagull”).

- The tag contains a music or speech label (e.g., “Male_speech_and_man_speaking”), even though a majority of its tags are not related to music and speech.

We then split the remaining 120 tags into SFX-Fg (85 tags) and SFX-Bg (35 tags). We treat the classification tasks for SFX-Fg and SFX-Bg separately, as we expect relatively poor performance for SFX-Bg due to its low relative level and long duration, which likely overlaps with multiple foreground events in DnR.

When using YAMNet outputs for the music genre, SFX-Fg, and SFX-Bg tasks, we first filter the 521 class probabilities output by YAMNet at each time frame to only contain those relevant for the given task (i.e., 11 classes for music, 85 for SFX-Fg, and 35 for SFX-Bg). In practice, we would use the boundaries provided by the activity detector to sum the class probabilities across all relevant frames for a given segment, and then output any tags with a class probability above a threshold. However, in our experiments, whose results are described in Section V-B, we aim to evaluate how source separation and remixing can aid soundtrack transcription without having activity detection performance play an out-sized role. We therefore use oracle event boundaries for summing YAMNet probabilities. This also allows us to use threshold independent metrics such as mean average (mAP) precision and area under the ROC curve (AUC), without having to account for missed detections and false alarms. Furthermore, it enables evaluation of the SFX-Fg and SFX-Bg tasks separately, given that our activity detector only outputs overall sound effect boundaries.

3) *Automatic Speech Recognition*: All audio is downsampled to 16 kHz, prior to being input to the HuBERT [51] feature extraction frontend of the pre-trained Conformer-based [50] ESPnet model⁴. As we will demonstrate in Section V-B, inputting an entire unsegmented 60 s DnR file, which contains multiple utterances interspersed with noise/silence regions, leads to highly sub-optimal ASR performance. Therefore, using the activity detector to first segment speech regions, and then passing each segment to the ASR model before concatenating all the outputs (as previously discussed and illustrated in Fig. 4) was essential to obtaining acceptable ASR performance. Since we evaluate the concatenated transcriptions with the ground-truth transcription from the entire soundtrack, we can easily compare performance using estimated activity detection boundaries and oracle boundaries in terms of word error rate (WER) and character error rate (CER) using the noisy mixture, separated speech stem, or remixed speech stem.

V. EXPERIMENTAL RESULTS

A. CFP Separation and activity detection

In this section, we evaluate source separation performance in terms of signal reconstruction metrics for listening applications, as well as the interaction between source separation and sound activity detection.

SI-SDR [14] is perhaps the most widely used objective measure for deep learning-based source separation, and was

⁴https://huggingface.co/espnet/simpleoier_librispeech_asr_train_asr_conformer7_hubert_1l60k_large_raw_en_bpe5000_sp

TABLE II

SI-SDR [DB] RESULTS OF BASELINES AND PROPOSED MODELS ON DnR. MUSHRA SCORE PREDICTIONS FROM THE 2F-MODEL ARE DENOTED AS 2F-M. AN EXTRA COLUMN IS INCLUDED FOR PESQ SCORES ON SPEECH.

Model	Music		SFX		Speech		
	SI-SDR	2f-m.	SI-SDR	2f-m.	SI-SDR	2f-m.	PESQ
No processing	-6.8	10.3	-5.0	10.9	1.0	6.3	2.05
Oracle PSF [57]	11.6	48.3	13.7	47.0	17.8	52.1	4.50
Conv-TasNet [44]	0.3	14.2	2.0	13.5	8.5	19.8	2.35
MaskTCN [44]	1.7	20.5	3.8	21.8	9.7	27.6	2.53
UMX ₆₄ [7]	3.1	22.2	4.4	21.5	11.7	30.9	2.67
XUMX ₃₂ [43]	2.9	21.6	4.7	21.7	11.2	29.9	2.62
XUMX ₆₄ [43]	3.5	22.7	5.1	22.4	11.7	30.6	2.66
XUMX ₁₂₈ [43]	3.7	24.0	5.1	23.2	11.6	30.7	2.66
XUMX ₂₅₆ [43]	2.9	22.6	4.4	22.0	10.5	29.7	2.58
XUMX _{64, multi-dec}	3.5	23.6	5.0	22.5	11.8	31.3	2.72
MRX (proposed)	4.2	25.6	5.7	24.7	12.3	32.9	2.85
MRX-C (proposed)	4.6	26.9	6.1	26.1	12.6	33.4	2.87

used in our preliminary study [13]. However, as shown in [58], SI-SDR is not among the objective metrics most correlated with perceptual quality. Therefore, in this work, we also report results using the 2f-model [19], which combines two mid-level perceptual features from the Perceptual Evaluation of Audio Quality (PEAQ) [59] standard (we used the PEAQ implementation from [60]). The 2f-model was fit to output estimates of MUSHRA scores ranging from 0 to 100, and sound signals were upsampled from 44.1 kHz to 48 kHz for input into the PEAQ model. To avoid any influence due to the scaling of the output, which is particularly an issue for models trained with SI-SDR loss, we normalize all outputs to have the same LUFS as their corresponding ground truth. For the speech source, we also report wideband Perceptual Evaluation of Speech Quality (PESQ) [18], where sound signals are downsampled to 16 kHz.

Table II presents the SI-SDR, PESQ, and 2f values of various models trained and tested on DnR, in addition to the “No Processing” condition (lower bound, using the mixture as estimate) and oracle phase sensitive mask [57] (a form of upper bound). For each model, SI-SDR improvements are fairly consistent across source types, despite the differences in their relative levels in the mix, which can be seen in the “No Processing” SI-SDR. In general, we observe that our proposed multi-resolution (MRX) models outperform all single-resolution baselines on all source types in terms of SI-SDR, PESQ, and the PEAQ-based 2f-model scores. This implies that the network learns to effectively combine information from different STFT resolutions to more accurately reconstruct the target sources. The performance of XUMX_{64, multi-dec} further confirms this hypothesis by performing nearly identically in comparison to XUMX₆₄, showing that the use of multiple decoders alone does not improve performance. We also observe that the single-source models (UMX) tend to perform comparably to the cross-source models (XUMX, MRX) for speech, but perform worse for music and SFX. We speculate that because music and SFX are quieter in the mix, it is harder for the network to isolate them effectively without the support of the other sources, while louder sources (here, speech) do not benefit from joint estimation.

TABLE III

SI-SDR [DB] RESULTS ON THE DNR TEST SET. THE COLUMNS “ORACLE TRAIN” AND “ORACLE TEST” DENOTE WHETHER THE MODEL HAS BEEN TRAINED AND TESTED USING ORACLE OR PREDICTED ACTIVITY LABELS. MRX-C REFERS TO THE MRX CONDITIONED ON ACTIVITY LABELS, MRX-C_{2p} TO MRX CONDITIONED ON ENHANCED ACTIVITY LABELS (2ND ITERATIVE PASS), AND MRX-MTL TO A MULTI-TASK LEARNING VERSION OF MRX WHICH PERFORMS BOTH SS AND ACTIVITY DETECTION.

	Activity labels		SI-SDR			
	Oracle Train	Oracle Test	Music	Speech	SFX	Avg.
No Processing	—	—	-6.8	1.0	-5.0	-3.6
MRX	—	—	4.2	12.3	5.7	7.4
MRX-C	✓	✓	5.1	12.6	6.5	8.1
MRX-C	✓	×	4.4	12.5	5.9	7.6
MRX-C	×	×	4.6	12.6	6.1	7.8
MRX-C _{2p}	×	×	4.6	12.6	6.1	7.8
MRX-MTL	—	—	3.8	11.8	5.3	7.0

TABLE IV

ACTIVITY DETECTION F-MEASURE FOR MUSIC (M), SPEECH (S), AND SOUND EFFECTS (X). THE BASELINE MODEL USES THE MIXTURE AS INPUT, MRX-MTL USES MULTI-TASK LEARNING FOR SOURCE SEPARATION AND EVENT DETECTION, WHILE THE OTHER ROWS INSERT SOURCE SEPARATION OUTPUT AT DIFFERENT LOCATIONS INSIDE THE NETWORK.

	Event-Based F-Measure			Segment-Based F-Measure		
	M	S	X	M	S	X
Baseline	0.77	0.84	0.50	0.97	0.97	0.92
MRX-MTL	0.78	0.87	0.43	0.97	0.97	0.92
Early Integration	0.82	0.88	0.55	0.98	0.98	0.94
Middle Integration	0.81	0.89	0.51	0.97	0.98	0.92
Late Integration	0.73	0.91	0.50	0.97	0.97	0.94

From Table II, we see that concatenating activity detection labels with the input spectrogram in the MRX-C model leads to the best overall separation performance, with some gains observed for all metrics compared to MRX. To further evaluate the impact of including activity detection labels as auxiliary inputs for source separation, Table III compares different settings for activity-conditioned source separation. We note that MRX-C leads to an average SI-SDR improvement of 0.7 dB compared to MRX when using oracle information, meaning the upper bound in expected performance improvement is limited, however larger improvements are observed for the quieter and more difficult to separate sources (i.e., music and SFX). When we switch to the realistic setup that does not use oracle information at test time, we see that training using estimated activity detection probabilities leads to a 0.2 dB improvement compared to using oracle information. We also observe that running multiple passes between activity detection and separation (MRX-C_{2p}) led to no performance improvement. Finally, we also considered a multi-task learning setup, where activity detection was used only as an additional training objective (MRX-MTL), but observed a degradation in performance compared to our plain MRX model.

While we have just seen how activity detection can improve separation, we now turn our focus to how separation can aid activity detection. Table IV displays the sound event detection performance in terms of parent-class F-measure computed using the SED EVAL package [61]. We used a threshold of

0.5, a collar of 750 ms for event-based metrics (with a 20% offset length), and a time resolution of 1 s for segment-based metrics. Compared to a Baseline taking the mixture signal as input, little improvement is observed using multi-task learning of activity detection and source separation (MRX-MTL). When integrating source separation output into our three-class activity detector, we observed the best performance using early or middle integration, which differs from the permutation-invariant analysis in [36] where late integration performed best. We suspect that integrating permutation-invariant source separation outputs at the input or at intermediate layers of the sound event detection network may cause difficulties during training, as the order in which separated outputs are stacked may change between training epochs. However, for problems such as the one studied in this paper, where separated outputs have a fixed ordering, the activity detection network does not have this inconsistency problem and hence early and middle integration perform better than late integration.

B. CFP Transcription and Remixing

1) *Audio Tagging (Music and SFX)*: We report audio tagging performance in terms of mAP and mAUC metrics using the implementation from [62]. Both mAP and mAUC allow us to evaluate audio tagging performance in a threshold-independent way for situations where multiple tags can be active for a single sound file, and are commonly reported for large-scale audio tagging [15], [53], [63]. Table V displays the audio tagging performance for music, SFX-Fg, and SFX-Bg sources using oracle event boundaries. The “Noisy” column is obtained from the original mixture and represents lower-bound performance, and the “GT” column represents upper-bound performance obtained using the ground-truth isolated sources. The scores for music are generally higher than for SFX, because there is only a single genre tag for each music segment, the set of possible music labels is smaller than for SFX, and SFX-Fg and SFX-Bg sounds may overlap as discussed in Section IV-A. The “MRX-C” column from Table V displays audio tagging performance using the separated outputs of MRX-C, our best separation model. It can be seen that for all sources in Table V, source separation improves audio tagging performance compared to the noisy mixture. The improvements for Music and SFX-Fg are larger than those for SFX-Bg, likely because those sources have higher relative levels in the DnR mixes.

The “MRX-C Remix” column in Table V shows the test set results of the systems which obtain the best performance on the validation set for each source in a grid search over remixing gains, as illustrated in Fig. 5. Figure 5 provides a detailed illustration of the impact of different remixing gains on mAP and mAUC for music and SFX. The first row of Figs. 5a, 5b, and 5c shows remixing results on the validation set for the case where interfering source gains are adjusted individually as in (3), while the second row shows the results on the validation and test sets for the case where the gains are adjusted jointly as in (4). For each plot in Fig. 5, the data points denoted by the intersections with the red line indicates the best validation set performance for the given metric over

TABLE V
MAP AND MAUC CLASSIFICATION FOR MUSIC, SFX-FOREGROUND (SFX-FG.), AND SFX-BACKGROUND (SFX-BG.). MRX OUT. RMX DEPICTS THE BEST REMIXING USE-CASE SCENARIO.

Source	Noisy		MRX-C		MRX-C Remix		GT	
	mAP	mAUC	mAP	mAUC	mAP	mAUC	mAP	mAUC
Music	0.233	0.682	0.297	0.723	0.300	0.718	0.336	0.772
SFX-Fg.	0.138	0.739	0.195	0.796	0.192	0.794	0.275	0.836
SFX-Bg.	0.137	0.653	0.143	0.656	0.148	0.666	0.258	0.767

different remixing SNR values. In almost all use-cases, we see that both mAP and mAUC benefit from some source remixing in comparison to using the predicted source as the sole signal input. While music and SFX-Fg show only minor classification improvement with remixing, in the case of SFX-Bg, performance peaks at a much lower SNR, meaning that this source specifically benefits by remixing the other sources for its classification task.

2) *Automatic Speech Recognition*: Table VI presents ASR results using ESPnet for three segmentation methods: using the oracle speech utterance onsets/offsets (“Oracle Boundaries”), using the onsets/offsets detected by Φ_{EI} (“VAD Boundaries”), and no segmentation (“No Boundaries”). We first note that in the “No Boundaries” case, ASR performance degrades dramatically, indicating that our pre-trained ASR model cannot handle 60 s long files containing multiple utterances. We tried training the ASR model from scratch on the DnR dataset speech files, but the performance still significantly lagged that of inputting segmented utterances into the pre-trained ASR model. We also note from Table VI that performance on the original LibriSpeech test set (“Libri. Test-Clean”) and the clean DnR speech submix (“DnR Speech GT”) is comparable in the case of oracle boundary segmentation, which is to be expected since they are identical utterances with slightly different levels added during the DnR data creation process. Performance degrades dramatically when the noisy DnR mixture is used for ASR, which is also to be expected since our model is pre-trained using clean speech. When using the separated speech stem from MRX-C, we obtain a 10.5% absolute (or 60% relative) reduction in WER compared to the noisy mixture when using VAD boundaries.

In the case of ASR, artifacts introduced by the separation process can degrade performance, and only partially separating the speech has been shown to improve performance [33]–[35]. Figure 6 displays the performance in terms of WER and PESQ on the validation set obtained for a grid search over the relative level between the remixed separated speech signal and the interference signals. Note that due to the large computational expense of speech decoding, we perform a grid search only for the combined (SFX+music) interference signal as described in (4). For WER, a minimum is observed at 17.5 dB, which is selected as the gain for the MRX-C Remix system. In the “MRX-C Remix” row of Table VI, we observe that by remixing the separated interference signal back with the estimated speech, WER is reduced from 7.0% down to 6.3% using VAD boundaries. In the case of PESQ, the non-remixed speech MRX-C output performs best. That is, MRX-C

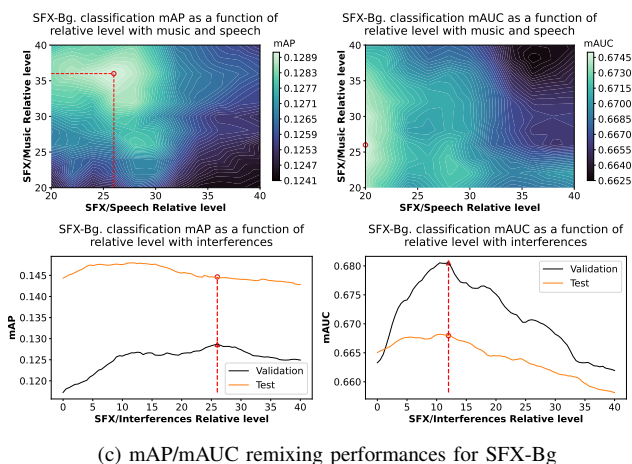
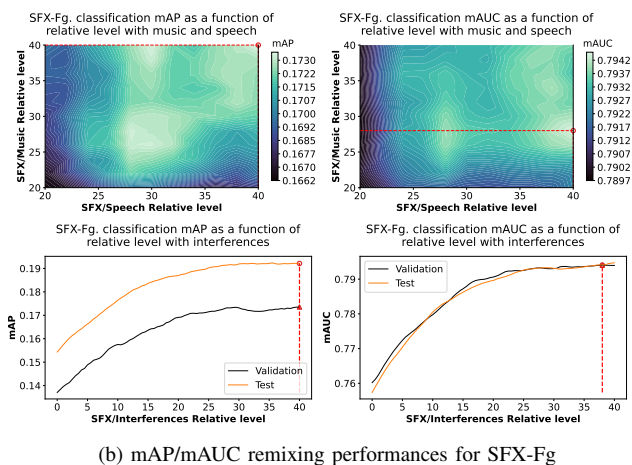
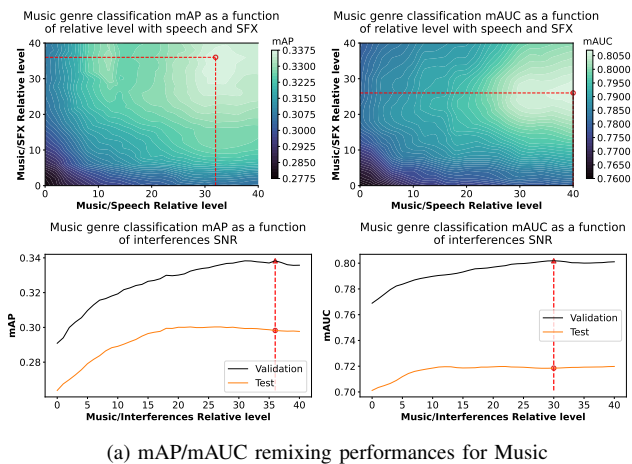


Fig. 5. mAUC and mAP music genre classification mappings as a function of Music SNR against interfering sources. The contour plots (first row) show the classification performance on the validation set depending on the SNR against each interfering source individually, while the second row shows the performance on the validation and test sets for the scenario where the two interfering sources are combined in order to compute the SNR. The red symbols along the red dashed lines denote the performance on the test set (circle) where the best validation set performance is observed (triangle).

TABLE VI
WER AND CER (%) RESULTS USING DIFFERENT BOUNDARY TYPES AND INPUT WAVEFORMS. THE PERFORMANCE ON LIBRISPEECH TEST-CLEAN IS ALSO INCLUDED FOR REFERENCE.

Data	Oracle Boundaries		VAD Boundaries		No Boundaries	
	WER	CER	WER	CER	WER	CER
Libri. Test-Clean	1.8	0.5	—	—	—	—
DnR Speech GT	1.8	1.3	2.0	1.4	27.5	19.5
Noisy Mixture	16.4	12.3	17.5	13.0	81.7	63.7
MRX-C	6.5	4.1	7.0	4.5	36.6	25.5
MRX-C Remix	5.7	3.7	6.3	4.1	45.5	32.2

reports a PESQ score of 2.87, while the best scenario denoted by our remixing experiments is found at SNR set to 40 dB with a PESQ score of 2.85. This underscores the importance of remixing for transcription, while the benefit for listening end-goals is unclear (at least in terms of PESQ), but this is an important topic of future work.

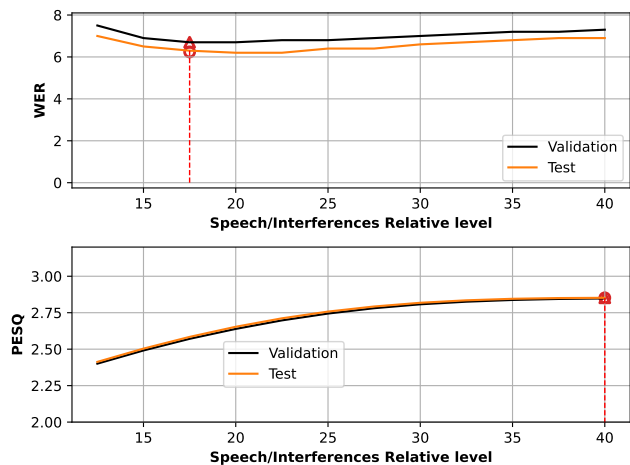


Fig. 6. WER (%) and PESQ performances on speech as a function of speech SNR against interfering sources (Music and SFX). Once again the symbols along the red dashed lines denote the points where the best validation set performance is observed (triangle) as well as the test-set performance associated with it (circle).

VI. CONCLUSIONS AND DISCUSSIONS

In this work, we extended our previous work on the cocktail fork problem by tackling transcription for each of the three sources involved: audio tagging for music and sound effect, and ASR for speech. We proposed an activity detection system for the three parent classes and showcased its benefits towards both the separation and transcriptions tasks; we first demonstrated how the system could help improve separation by using the activity labels as conditioning information. Secondly, we described how the integration of an activity detection mechanism was essential in order to tackle real-world soundtrack transcription tasks. We led an extensive investigation to show how source remixing could help towards transcription and demonstrated that mixing back the interfering signals with the isolated source estimates could help improve performance on their associated transcription downstream task.

In the present work, we explored how transcription benefited from source-remixing strategies, due to the imperfect nature of our separator output. Beside negatively impacting the transcription downstream tasks, the presence of separation artifacts undoubtedly deteriorates the listening experience as well. Moving forward, we aim to explore source-remixing strategies that could minimize perceptual artifacts for the separator output and enhance the listening experience. While we approached this work from a fully supervised angle, for both separation and transcription, taking advantage of the large amount of “real-world” unlabeled data available is an important topic for our future work.

REFERENCES

- [1] D. Wang and J. Chen, “Supervised speech separation based on deep learning: An overview,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 26, no. 10, pp. 1702–1726, 2018.
- [2] C. K. Reddy, V. Gopal, R. Cutler, E. Beyrami, R. Cheng, H. Dubey, S. Matushevych, R. Aichner, A. Aazami, S. Braun *et al.*, “The Interspeech 2020 Deep Noise Suppression challenge: Datasets, subjective testing framework, and challenge results,” in *Proc. Interspeech*, Oct. 2020.
- [3] J. R. Hershey, Z. Chen, and J. Le Roux, “Deep clustering: Discriminative embeddings for segmentation and separation,” in *Proc. ICASSP*, Mar. 2016, pp. 31–35.
- [4] L. Drude, J. Heitkaemper, C. Boeddeker, and R. Haeb-Umbach, “SMS-WSJ: Database, performance measures, and baseline recipe for multi-channel source separation and recognition,” in *arXiv preprint arXiv:1910.13934*, 2019.
- [5] G. Wichern, J. Antognini, M. Flynn, L. R. Zhu, E. McQuinn, D. Crow, E. Manilow, and J. Le Roux, “WHAM!: Extending speech separation to noisy environments,” in *Proc. Interspeech*, Sep. 2019.
- [6] Z. Rafii, A. Liutkus, F.-R. Stöter, S. I. Mimitakis, and R. Bittner, “The MUSDB18 corpus for music separation,” Dec. 2017. [Online]. Available: <https://doi.org/10.5281/zenodo.1117372>
- [7] F.-R. Stöter, S. Uhlich, A. Liutkus, and Y. Mitsufuji, “Open-Unmix - a reference implementation for music source separation,” *Journal of Open Source Software*, 2019.
- [8] E. Manilow, G. Wichern, P. Seetharaman, and J. Le Roux, “Cutting music source separation some Slakh: A dataset to study the impact of training data quality and quantity,” in *Proc. WASPAA*, Oct. 2019.
- [9] I. Kavalerov, S. Wisdom, H. Erdogan, B. Patton, K. Wilson, J. Le Roux, and J. R. Hershey, “Universal sound separation,” in *Proc. WASPAA*, Oct. 2019.
- [10] E. Tzinis, S. Wisdom, J. R. Hershey, A. Jansen, and D. P. W. Ellis, “Improving universal sound separation using sound classification,” in *Proc. ICASSP*, May 2020.
- [11] F. Pishdadian, G. Wichern, and J. Le Roux, “Finding strength in weakness: Learning to separate sounds with weak supervision,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 28, pp. 2386–2399, 2020.
- [12] 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. Interspeech*, Oct. 2020.
- [13] D. Petermann, G. Wichern, Z.-Q. Wang, and J. Le Roux, “The cocktail fork problem: Three-stem audio separation for real-world soundtracks,” in *Proc. ICASSP*, May 2021.
- [14] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, “SDR—half-baked or well done?” in *Proc. ICASSP*, May 2019.
- [15] 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. ICASSP*, Mar. 2017.
- [16] C. Li, J. Shi, W. Zhang, A. S. Subramanian, X. Chang, N. Kamo, M. Hira, T. Hayashi, C. Boeddeker, Z. Chen, and S. Watanabe, “ESPnet-SE: End-to-end speech enhancement and separation toolkit designed for ASR integration,” in *Proc. SLT*, Jan. 2021, pp. 785–792.
- [17] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “LibriSpeech: An ASR corpus based on public domain audio books,” in *Proc. ICASSP*, Apr. 2015, pp. 5206–5210.
- [18] A. Rix, J. Beerends, M. Hollier, and A. Hekstra, “Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs,” in *Proc. ICASSP*, vol. 2, May 2001, pp. 749–752.
- [19] T. Kastner and J. Herre, “An efficient model for estimating subjective quality of separated audio source signals,” in *Proc. WASPAA*, Oct. 2019, pp. 95–99.
- [20] E. Manilow, P. Seetharaman, and B. Pardo, “Simultaneous separation and transcription of mixtures with multiple polyphonic and percussive instruments,” in *Proc. ICASSP*, May 2020, pp. 771–775.
- [21] J. F. Woodruff, B. Pardo, and R. B. Dannenberg, “Remixing stereo music with score-informed source separation,” in *Proc. ISMIR*, Oct. 2006.
- [22] H. Yang, K. Zhen, S. Beack, and M. Kim, “Source-aware neural speech coding for noisy speech compression,” in *Proc. ICASSP*, Jun. 2021, pp. 706–710.
- [23] A. S. Subramanian, X. Wang, M. K. Baskar, S. Watanabe, T. Taniguchi, D. Tran, and Y. Fujita, “Speech enhancement using end-to-end speech recognition objectives,” in *Proc. WASPAA*, Oct. 2019, pp. 229–233.
- [24] X. Chang, T. Maekaku, Y. Fujita, and S. Watanabe, “End-to-end integration of speech recognition, speech enhancement, and self-supervised learning representation,” in *arXiv preprint arXiv:2204.00540*, 2022.
- [25] N. Moritz, G. Wichern, T. Hori, and J. Le Roux, “All-in-one transformer: Unifying speech recognition, audio tagging, and event detection,” in *Proc. Interspeech*, Oct. 2020.
- [26] O. Slizovskaia, L. Kim, G. Haro, and E. Gomez, “End-to-end sound source separation conditioned on instrument labels,” in *Proc. ICASSP*, May 2019, pp. 306–310.
- [27] Y.-N. Hung and A. Lerch, “Multitask learning for instrument activation aware music source separation,” in *Proc. ISMIR*, Oct. 2020.
- [28] E. Verteletskaya and K. Sakniov, “Voice activity detection for speech enhancement applications,” *Acta Polytechnica*, vol. 50, Jan. 2010.
- [29] X. Tan and X.-L. Zhang, “Speech enhancement aided end-to-end multitask learning for voice activity detection,” in *Proc. ICASSP*, Jun. 2021, pp. 6823–6827.
- [30] N. L. Westhausen, R. Huber, H. Baumgartner, R. Sinha, J. Rennie, and B. T. Meyer, “Reduction of subjective listening effort for TV broadcast signals with recurrent neural networks,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 3541–3550, 2021.
- [31] M. Torcoli, J. Paulus, T. Kastner, and C. Uhle, “Controlling the remixing of separated dialogue with a non-intrusive quality estimate,” in *Proc. WASPAA*, Oct. 2021.
- [32] O. Kalinli, M. L. Seltzer, J. Droppo, and A. Acero, “Noise adaptive training for robust automatic speech recognition,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, pp. 1889–1901, 2010.
- [33] Y. Koizumi, S. Karita, A. Narayanan, S. Panchapagesan, and M. Bacchi-ani, “SNRi target training for joint speech enhancement and recognition,” *arXiv preprint arXiv:2111.00764*, 2021.
- [34] H. Sato, T. Ochiai, M. Delcroix, K. Kinoshita, N. Kamo, and T. Moriya, “Learning to enhance or not: Neural network-based switching of enhanced and observed signals for overlapping speech recognition,” in *Proc. ICASSP*, May 2022, pp. 6287–6291.
- [35] K. Iwamoto, T. Ochiai, M. Delcroix, R. Ikeshita, H. Sato, S. Araki, and S. Katagiri, “How bad are artifacts?: Analyzing the impact of speech enhancement errors on ASR,” in *Proc. Interspeech*, Sep. 2022.
- [36] N. Turpault, S. Wisdom, H. Erdogan, J. R. Hershey, R. Serizel, E. Fonseca, P. Seetharaman, and J. Salamon, “Improving sound event detection in domestic environments using sound separation,” in *Proc. DCASE*, Nov. 2020, pp. 205–209.
- [37] H. Yang, S. Firodiya, N. J. Bryan, and M. Kim, “Don’t separate, learn to remix: End-to-end neural remixing with joint optimization,” in *Proc. ICASSP*, May 2022, pp. 116–120.
- [38] C. Gupta, E. Yilmaz, and H. Li, “Automatic lyrics alignment and transcription in polyphonic music: Does background music help?” in *Proc. ICASSP*, May 2020, pp. 496–500.
- [39] C.-Y. Chiu, J. Ching, W.-Y. Hsiao, Y.-H. Chen, A. W.-Y. Su, and Y.-H. Yang, “Source separation-based data augmentation for improved joint beat and downbeat tracking,” in *Proc. EUSIPCO*, Aug. 2021, pp. 391–395.
- [40] L. Lin, Q. Kong, J. Jiang, and G. Xia, “A unified model for zero-shot music source separation, transcription and synthesis,” in *Proc. ISMIR*, Nov. 2021.
- [41] P. Chandna, H. Cuesta, D. Petermann, and E. Gómez, “A deep-learning based framework for source separation, analysis, and synthesis of choral ensembles,” *Frontiers in Signal Processing*, vol. 2, 2022.
- [42] H. Wierstorf, D. Ward, R. Mason, E. M. Graiss, C. Hummersone, and M. Plumbley, “Perceptual evaluation of source separation for remixing music,” in *Proc. AES Convention*, Oct. 2017.
- [43] R. Sawata, S. Uhlich, S. Takahashi, and Y. Mitsufuji, “All for one and one for all: Improving music separation by bridging networks,” in *Proc. ICASSP*, Jun. 2021, pp. 51–55.

- [44] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time-frequency magnitude masking for speech separation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [45] J. Baumann, P. Meyer, T. Lohrenz, A. Roy, M. Papendieck, and T. Fingscheidt, "A new DCASE 2017 rare sound event detection benchmark under equal training data: CRNN with multi-width kernels," in *Proc. ICASSP*, Jun. 2021, pp. 865–869.
- [46] D. De Benito-Gorrón, D. Ramos, and D. T. Toledano, "A multi-resolution CRNN-based approach for semi-supervised sound event detection in DCASE 2020 challenge," *IEEE Access*, vol. 9, pp. 89 029–89 042, 2021.
- [47] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "FiLM: Visual reasoning with a general conditioning layer," *Proc. AAAI*, Feb. 2018.
- [48] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [49] J. F. Gemmeke, D. P. W. 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. ICASSP*, Mar. 2017.
- [50] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, "Conformer: Convolution-augmented transformer for speech recognition," *arXiv preprint arXiv:2005.08100*, 2020.
- [51] W. Hsu, B. Bolte, Y. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "HuBERT: Self-supervised speech representation learning by masked prediction of hidden units," *arXiv preprint arXiv:2106.07447*, 2021.
- [52] M. Defferrard, K. Benzi, P. Vandergheynst, and X. Bresson, "FMA: A dataset for music analysis," in *Proc. ISMIR*, Oct. 2017.
- [53] E. Fonseca, X. Favory, J. Pons, F. Font, and X. Serra, "FSD50K: An open dataset of human-labeled sound events," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 30, pp. 829–852, 2022.
- [54] L. Zhang, C. Li, F. Deng, and X. Wang, "Multi-task audio source separation," *arXiv preprint arXiv:2107.06467*, 2021.
- [55] T. Menne, I. Sklyar, R. Schlüter, and H. Ney, "Analysis of deep clustering as preprocessing for automatic speech recognition of sparsely overlapping speech," in *Proc. Interspeech*, Sep. 2019, pp. 2638–2642.
- [56] S. Chaudhuri, J. Roth, D. P. Ellis, A. Gallagher, L. Kaver, R. Marvin, C. Pantofaru, N. Reale, L. G. Reid, K. Wilson *et al.*, "AVA-speech: A densely labeled dataset of speech activity in movies," in *Proc. Interspeech*, Sep. 2018.
- [57] H. Erdogan, J. Hershey, S. Watanabe, and J. Le Roux, "Phase-sensitive and recognition-boosted speech separation using deep recurrent neural networks," in *Proc. ICASSP*, Apr. 2015.
- [58] M. Torcoli, T. Kastner, and J. Herre, "Objective measures of perceptual audio quality reviewed: An evaluation of their application domain dependence," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 29, pp. 1530–1541, 2021.
- [59] C. Colomes, C. Schmidmer, T. Thiede, and W. C. Treurniet, "Perceptual quality assessment for digital audio: PEAQ-the new ITU standard for objective measurement of the perceived audio quality," *Journal of the Audio Engineering Society*, Sep. 1999.
- [60] P. Kabal, "An examination and interpretation of ITU-R BS. 1387: Perceptual evaluation of audio quality," Dept. Electrical & Computer Engineering, McGill University, Tech. Rep., 2002.
- [61] A. Mesaros, T. Heittola, and T. Virtanen, "Metrics for polyphonic sound event detection," *Applied Sciences*, vol. 6, no. 6, p. 162, 2016.
- [62] H. Chen, W. Xie, A. Vedaldi, and A. Zisserman, "VGGSound: A large-scale audio-visual dataset," in *Proc. ICASSP*, May 2020, pp. 721–725.
- [63] S. Hershey, S. Chaudhuri, D. P. W. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold, M. Slaney, R. J. Weiss, and K. W. Wilson, "CNN architectures for large-scale audio classification," in *Proc. ICASSP*, Mar. 2017, pp. 131–135.