# CONVOIFILTER: A CASE STUDY OF DOING COCKTAIL PARTY SPEECH RECOGNITION

*Thai-Binh Nguyen[1], Alexander Waibel[1,2]*

[1]Karlsruhe Institute of Technology
[2]Carnegie Mellon University
`thai-binh.nguyen@kit.edu`

## ABSTRACT

This paper presents an end-to-end model designed to improve automatic speech recognition (ASR) for a particular speaker in a crowded, noisy environment. The model utilizes a single-channel speech enhancement module that isolates the speaker's voice from background noise (ConVoiFilter) and an ASR module. The model can decrease ASR's word error rate (WER) from 80% to 26.4% through this approach. Typically, these two components are adjusted independently due to variations in data requirements. However, speech enhancement can create anomalies that decrease ASR efficiency. By implementing a joint fine-tuning strategy, the model can reduce the WER from 26.4% in separate tuning to 14.5% in joint tuning. We openly share our pre-trained model to foster further research `hf.co/nguyenvulebinh/voice-filter`.

***Index Terms***— ASR, Speech Enhancement, Voice Filter

## 1. INTRODUCTION

In the ideal environment with a close speaking microphone, the speech recognition performance of current ASR systems can surpass humans, with a common word error rate usually below 5% [1]. However, in a realistic environment, the ASR system has to deal with complex acoustic conditions like noise, reverberation, cross-talk (known as the cocktail party setting) [2, 3]. As a result, it makes the model performance drops dramatically. For example, in CHiME-5 competition, WER of barely below 80% achieved by the baseline system; using a robust back-end, approximately 60% WER is achieved [4].

Unlike machines, humans do an outstanding job of ignoring interfering signals and focusing on what we want to hear [5, 6]. As deep learning gains popularity in speech signal processing, much impressive progress has been proposed to help enhance speech signals [7] like denoising, dereverberation, source separation, and neural beamforming. Among speech enhancement techniques, masked-based is among the most popular and effective. In masking approaches, rather than estimating the enhanced signal directly, we estimate a mask, then multiply it with the noisy signal to get the enhanced signal. Depending on the type of input/output we can have waveform masking and spectral masking. Our study based on spectral masking since it's much faster than waveform masking.

Speech enhancement techniques generally use for blind signal enhancement. In our case study, we want to develop a robot to communicate and take orders from its master. So, we know precisely to whom the robot needs to listen, which can provide critical information that helps our speech recognition system work better, especially in complex acoustic situations like cocktail parties.

There are some related studies for this circumstance, known as the *speaker extraction* problem, including DENet [8], SpeakerBeam [9], and VoiceFilter [10, 11]. However, our proposal has a few significant distinctions from them: (1) we utilize an x-vector pre-trained model [12] instead of i-vector or d-vector in [8, 9, 10, 11] since the x-vector brings better results in our experiment. (2) We use scale-invariant source-to-noise ratio (SI-SNR) [13] as a loss function because it is a speech enhancement evaluation metric and a training target that makes optimizing and choosing the best model more precise. (3) Different from [8, 9], we focus on improving WER like [10, 11] but joint tuning ASR and speaker extraction model rather than optimizing the loss function. (4) We make a pre-trained self-supervised model based on wav2vec2 [14] architecture that works better for the noise acoustic condition. (5) For the mask estimation model, we use Conformer block [15] rather than LSTM and CNN in [10, 11]. Furthermore, we introduce a cross-extraction mechanism between the reference signal and noisy signal for speaker embedding, which enhances the performance of our model, as demonstrated in our experiments.

## 2. MODEL DESCRIPTION

Our system consists of two primary modules: speaker extraction, which enhances the target speaker's voice, and the ASR module. In the following subsections, we will describe each module and our approach to jointly tuning them.

### 2.1. Target speaker's voice enhancement

Figure 1 shows an overview of our target speaker's extraction module, ConVoiFilter. This module aims to remove all noise and interfering speech from the noisy audio input, producing a clean utterance for the target speaker.

Firstly, an embedding vector identifies the target speaker ($e_{ref} \in \mathbb{R}^{d\_emb}$) extracted from their audio recordings (reference utterance) using a speaker encoder module. Moreover, we perform cross-extraction of speaker embedding from the noisy audio ($e_{noisy} \in \mathbb{R}^{d\_emb}$). Afterwards, we concatenate the two embeddings (from the reference utterance and the noisy audio) and feed them through a feed-forward network to produce the final presentation for the target speaker ($e = FFN(e_{ref}, e_{noisy}) \in \mathbb{R}^{d\_emb}$). Without cross-extraction,
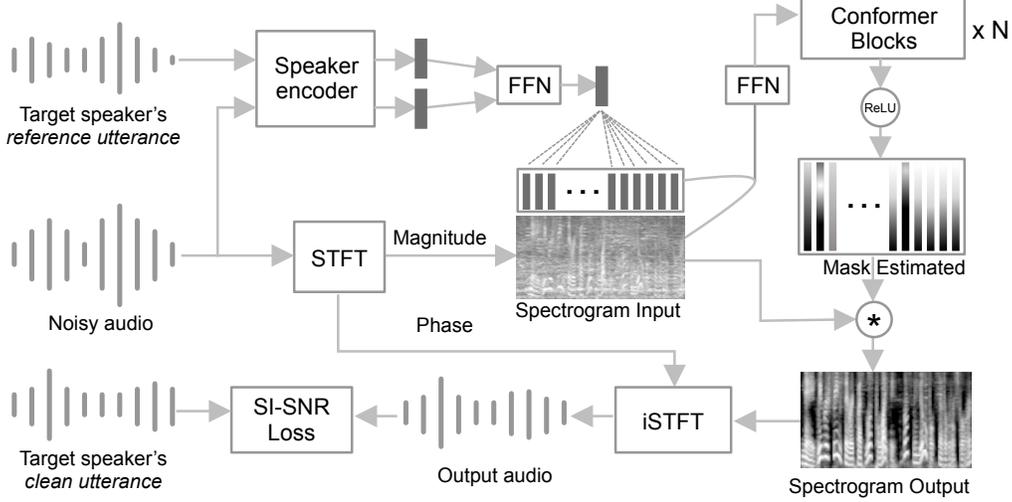
**Fig. 1**. Overview of the ConVoiFilter model.

$e = e_{ref}$. We experimented with two different speaker encoder models: x-vector [12] and i-vector [16].

Secondly, we utilize the short-time Fourier transform (STFT) to process the noisy audio and generate a magnitude spectrogram ($F \in \mathbb{R}^{S*(\frac{n\_fft}{2}+1)}$) that will aid in the mask estimation process. In this step, we preserve the phase ($P$) for the inverse STFT (iSTFT), which is used to reconstruct the output audio. The variable $S$ denotes the number of time frames, and $n\_fft$ represents the size of the Fourier transform. To create the features required to guide the mask estimation model, we horizontally concatenate the target speaker's embedding $E$ with the magnitude spectrogram $F$, resulting in $\hat{G} = [F \quad e] \in \mathbb{R}^{S*(\frac{n\_fft}{2}+1+d\_emb)}$.

Next, $N$ conformer blocks used to transforms $\hat{G}$ into a mask $M \in \mathbb{R}^{S*(\frac{n\_fft}{2}+1)}$ of the same shape as the magnitude spectrogram $F$. Since $\hat{G}$ and $M$ is in different shape and we want to maintain the original conformer blocks, we need to transform $\hat{G}$ into $G = FFN(\hat{G}) \in \mathbb{R}^{S*(\frac{n\_fft}{2}+1)}$ which has the same shape with $M$. A conformer block (formula 1) consists of four modules stacked together (described in [15]) where $x_i$ is the input to conformer block $i$ ($x_0 = G$). We selected conformer because it incorporates convolution and multi-head self-attention, both of which are effective in utilizing contextual information, which is crucial for detecting interfering signals. The mask's purpose is to amplify or attenuate the amplitude of certain frequencies in the spectrogram, and therefore, it must be greater than or equal to zero. The ReLU function is applied to the output of the conformer block ($M = ReLU(y_N)$).

$$
\begin{aligned}
\tilde{x}_i &= x_i + \frac{1}{2}\text{FFN}(x_i) \\
x'_i &= \tilde{x}_i + \text{MHSA}(\tilde{x}_i) \\
x''_i &= x'_i + \text{Conv}(x'_i) \\
y_i &= \text{Layernorm}(x''_i + \frac{1}{2}\text{FFN}(x''_i))
\end{aligned}
\tag{1}
$$

Finally, the estimated mask $M$ is multiplied element-wise with the magnitude spectrogram $F$ to obtain an enhanced magnitude spectrogram. Phase information $P$ is combined with this enhanced magnitude spectrogram to reconstruct the output audio $iSTFT(M \odot F, P)$. We evaluate audio quality using the SI-SNR [13] loss function, which compares it with the clean utterance of the target speaker.

### 2.2. Automatic speech recognition

Self-supervised learning of speech representations [17] has recently shown its effectiveness in utilizing unlabeled speech data, resulting in outperforming the state-of-the-art (SoTA) in many automatic speech recognition (ASR) datasets. For our study, we utilized the pre-trained wav2vec2 model [14] to construct our ASR model. The wav2vec2 model acts as a speech encoder, and for the decoder, we used an RNN transducer [18]. Despite having a speech enhancement module to eliminate noise from the audio, the output may still contain noise. To address this issue, we utilized the self-supervised learning capabilities of wav2vec2 and created a pre-trained model by incorporating noise and room reverb into the unlabeled data (see section 3.1 for dataset details). Our subsequent experiment demonstrated that this approach significantly enhances the system's accuracy.

### 2.3. Joint fine-tuning strategy

ConVoiFilter is expected to produce only the target speaker's voice. However, in practice, speech enhancement always comes up with unknown artifacts. In a naive way, we can directly connect the enhancing module to the ASR module and optimize their total loss. However, there are two reasons why it does not work. Firstly, the enhancement works in a high resolution of the input (e.g., an audio 15s, rate of 16kHz, STFT has hop size of 128 will output around 2000 time-frames). It makes the mask estimation module hard to work. Whereas with the wav2vec model, the same audio will output about 700 time-frames. Secondly, ConVoiFilter quickly fails the ASR module at the beginning of optimization because its output is too noisy.

We handle these issues with a chunk-merging strategy. Below is the pseudo-code of the forward function. First, long

audio is split into smaller chunks (line 4) to optimize the enhancing module, then the output (line 6) to optimize the ASR module. Audio input is padded into integer times the chunk's size to ensure all chunks are the same after splitting.

```
1  def forward(noisy_audio, spk_embed,
2               clean_audio, label):
3    chunk_size = 5 #5s each chunk
4    noisy_chunks = split(noisy_audio, chunk_size)
5    output_chunks = enhance(noisy_chunks, spk_embed)
6    output_audio = merge(output_chunks)
7    enh_loss = si_snr(output_audio, clean_audio)
8    if snr_loss < threshold: #enhancing is well
9      output_text = asr(output_audio)
10   else:
11     output_text = asr(clean_audio)
12   asr_loss = transducer_loss(output_text, label)
13   loss = enh_loss + asr_loss
14   return loss
```

In the beginning, the enhancement module can output randomly. So, we use a threshold to decide if the enhancement module works well; then, the ASR module will learn from its output; else, the ASR model will learn from the clean audio.

## 3. EXPERIMENTAL SETUP

### 3.1. Data preparation

In our setup, we need a clean utterance of a specific speaker, noisy audio containing that utterance, and that speaker's embedding. Although CHiME-5 [19] is like a cocktail party dataset; however, the clean utterance is no warranty. It shows through the WERs for the development set using the binaural microphones (clean utterance of a speaker) reported around 47.9%. Instead, we train and evaluate the system using our generated data. We use audio from the LibriSpeech dataset [20] (2338 speakers for training, 73 speakers for testing). The ambient noise dataset includes MUSAN and WHAM [21, 22] (a total of 189 hours including music, speech, and environmental noise, 169 hours for training, 20 hours for testing). The reverb dataset is from Room RIR and BUT Speech@FIT [23, 24] (2650 room impulse response signals, 2350 signals for training, 300 signals for testing).
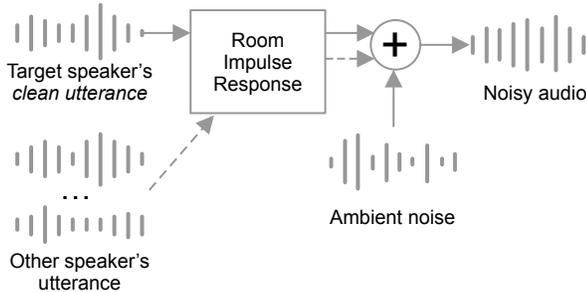


**Fig. 2**. Data preparation pipeline

Figure 2 shows our data generation pipeline. Noisy audio is cumulative of a few utterances and ambient noise. An utterance can be convolved with a room impulse response. We use random variables in the generating process to make the dataset more diverse. The number of other speakers' utterances is random, from 0 to 3. Ambient noise is added in 80% number of times. The room impulse response is applied 30% number of times. All other speaker's utterances and ambient sound is normalized (based on the target speaker's clean utterance) with SNR value around 1 to 20dB.

Both training processes (self-supervised wav2vec2 and ConVoiFilter model) use this data pipeline. One difference is that when we train self-supervised wav2vec2, we do not use other speakers' utterances (the dashed arrow in figure 2 - means no cross-talk) because it makes the data too noisy and can fail the wav2vec2 model.

### 3.2. Model setup

In our experiment, the ConVoiFilter model has 4 layer conformers with a hidden size is 1024, STFT has $n\_fft = 512$, hope size is 128, speaker embedding is x-vector[12] model has $d\_emb = 512$. A speaker's embedding is calculated from multiple utterances of that speaker. For the ASR model, we use the wav2vec2-base architecture with 12 hidden layers; the hidden size is 768.

### 3.3. Evaluation

Our system was evaluated using five different types of data. The first is the "Clean Audio" which consists of the original audio. The second type is the "Noisy Audio" which is the output of the data processing pipeline (depicted in Figure 2) applied to the clean audio. The third type is the "No cross-talk", a subset of noisy audio, which contains only single speaker samples with both ambient noise and speech reverberation. The fourth type is the "Ambient noise", a subset of noisy audio, which contains samples with cross-talk and ambient noise only. Finally, the fifth type is the "Reverberation", a subset of noisy audio, which contains samples with cross-talk and reverberation only.

We evaluated our system using four different model settings to assess the WER on various types of data (table 1). The first two settings consisted of ASR models only, which aimed to measure the ASR model's ability to handle noisy data. The first ASR model, named ASR_based, was initialized from the pre-trained wav2vec2 base model [14], which was trained with 960 hours of Librispeech data. The second ASR model, ASR_noisy, was initialized from our pre-trained wav2vec2 base model, which was trained with the same 960 hours of data, but augmented with noise and reverb data. The remaining two settings incorporated a speech enhancement module. The third was a cascade model, in which ConVoiFilter and ASR were trained independently. The final model was end-to-end, where ConVoiFilter and ASR were jointly trained. The first two ASR models were trained with noisy audio (without cross-talk). In contrast, the cascade and end-to-end models were trained with noisy audio that may have cross-talk.

## 4. RESULTS

The WER for different model settings on various data types is presented in Table 1. Generally, the end-to-end models perform better than other models across most input sets, except for clean audio. In cases with cross-talk, the ConVoiFilter model demonstrates its effectiveness by using speaker embedding to extract the target speaker's voice, resulting in cleaner audio and a significant improvement in WER (the result in the cross-talk column, rows 2 and 4). The "Noisy

| Model setting | Clean Audio | No cross-talk | Cross-talk | | Noisy Audio |
| --- | --- | --- | --- | --- | --- |
| | | | Ambient noise | Reverberation | |
| ASR_based | 2.22 | 16.21 | 50.72 | 90.87 | 80.04 |
| ASR_noisy | **2.03** | 12.12 | 45.12 | 84.14 | 75.19 |
| Cascade ConVoiFilter + ASR_noisy | 3.51 | 11.59 | 20.24 | 30.30 | 26.40 |
| End-to-end ConVoiFilter-ASR_noisy | 3.36 | **9.41** | **13.23** | **25.14** | **14.51** |

**Table 1**. %WERs for different types of model settings. The ConVoiFilter model uses x-vector as speaker encoder

| System | Overlap ratio in % | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | 0S | 0L | 10 | 20 | 30 | 40 |
| Baseline [25] | 8.4 | 8.3 | 11.6 | 15.8 | 18.7 | 21.7 |
| Whisper-large | **3.64** | **3.4** | 8.86 | 15.64 | 23.55 | 32.73 |
| ConVoiFilter + Whisper-large | 5.35 | 5.59 | **7.32** | **13.28** | **14.46** | **16.83** |

**Table 2**. %WERs for LibriCSS utterance-wise evaluation

| Method | SI-SNR | SDR |
| --- | --- | --- |
| No Enhancement | 1.04 | 1.14 |
| ConVoiFilter | **13.97** | **15.14** |
| x-vector $\rightarrow$ i-vector | 10.02 | 11.14 |
| Conformer $\rightarrow$ bi-LSTM | 8.12 | 9.58 |
| SI-SNR loss $\rightarrow$ MSE loss | 9.11 | 10.81 |

**Table 3**. Ablation study over each change in the ConVoiFilter model. Source to distortion ratio (SDR) and Scale invariant signal to noise ratio (SI-SNR) in dB.

Audio" column displays the overall model performance, with only the ASR model achieving the best WER at 75.19% (ASR_noisy). When combined with ConVoiFilter model, the WER can be significantly reduced to 26.40%. The end-to-end model can further improve the performance, achieving a WER of 14.51%. Additionally, our pre-trained wav2vec2 model (with noisy audio) proves its worth, with ASR_noisy outperforming ASR_based in all input sets.

Table 2 presents the WER for the LibriCSS[25] dataset, a 10-hour real-recorded dataset derived from the LibriSpeech corpus featuring speaker conversations. The baseline comprises BLSTM ASR, speaker separation, and MVDR based on 7-channels. We compare this baseline with a robust ASR model (Whisper large [26]) and also explore a combination of our ConVoiFilter with Whisper. The results clearly indicate that ConVoiFilter significantly improves Whisper, particularly in reducing WER in overlapping audio.

Table 3 presents an ablation study comparing our proposal to other studies on the speech enhancement model for a target speaker (speaker extraction problem). The effectiveness is measured using two standard metrics, SI-SNR and SDR, expressed in dB, where a higher value indicates better performance. Our ConVoiFilter differs from the recent VoiceFilter [10, 11] in three key aspects, namely the speaker encoder (x-vector), mask estimation model (conformer), and loss function (SI-SNR loss). Table 3 demonstrates the improvement resulting from each change we made. Firstly, replacing bi-LSTM with Conformer resulted in the most significant gain (5.56 points in SDR). Secondly, the SI-SNR loss function outperformed the MSE loss used in [10]. We speculate that the direct computation of the audio signal by the SI-SNR loss provides better optimization signals than the spectrogram-based MSE loss. Thirdly, although the x-vector and i-vector were trained with the same dataset, our experiment showed that the x-vector provided better results than the i-vector.

Figure 3 illustrates the effectiveness of cross-extraction of speaker embedding. In this example, the input is a mixture of two people (labeled as per_1 and per_2) and ambi-

ent noise. When we use the speaker embedding from per_1, both the model with and without cross-extraction can extract the target speaker as per_1. However, if the speaker embedding is obtained from a random person, only the model with cross-extraction can output a blank speech signal. Without the cross-extraction mechanism, the model extracts the wrong speech signal (the speech signal of per_2).
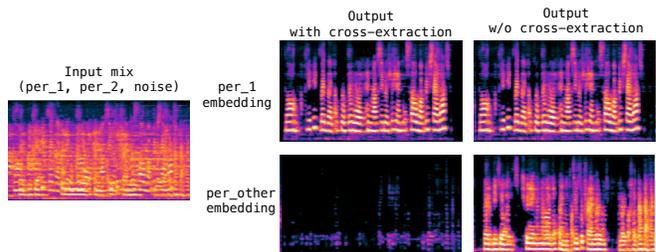


**Fig. 3**. An illustration spectrogram demonstrating the extraction of the target speaker's voice.

## 5. CONCLUSION

This paper details a case study on cocktail party speech recognition. Instead of recognizing all speakers, our system focuses on enhancing the target speaker's voice before conducting speech recognition. Through rigorous experiments, we showcased the effectiveness of our improved end-to-end model. Noteworthy enhancements include a cross-extraction speaker encoder, an improved mask estimation model, and an optimized loss function. We also publicly share our pre-trained ConVoiFilter to support ongoing research.

## 6. ACKNOWLEDGMENT

# 7. REFERENCES

[1] Thai-Son Nguyen, Sebastian Stüker, and Alex Waibel, "Super-human performance in online low-latency recognition of conversational speech," in *Interspeech 2021*. 2021, ISCA.

[2] Jie Yang, Alex Waibel, and et al., "Multimodal people id for a multimedia meeting browser," in *Proceedings of the seventh ACM international conference on Multimedia (Part 1)*, 1999, pp. 159–168.

[3] Michael Bett, Alex Waibel, and et al., "Multimodal meeting tracker.," in *RIAO*. Citeseer, 2000, pp. 32–45.

[4] Cătălin Zorilă, Christoph Boeddeker, Rama Doddipatla, and Reinhold Haeb-Umbach, "An investigation into the effectiveness of enhancement in asr training and test for chime-5 dinner party transcription," in *2019 IEEE (ASRU)*. IEEE, 2019.

[5] Alex Waibel, Hartwig Steusloff, Rainer Stiefelhagen, and Kym Watson, *Computers in the Human Interaction Loop*, pp. 3–6, Springer London, London, 2009.

[6] Alex Waibel, *Beyond CHIL*, pp. 367–371, Springer London, London, 2009.

[7] Reinhold Haeb-Umbach, Jahn Heymann, Lukas Drude, Shinji Watanabe, Marc Delcroix, and Tomohiro Nakatani, "Far-field automatic speech recognition," *Proceedings of the IEEE*, vol. 109, no. 2, pp. 124–148, 2021.

[8] J. Wang, Jie Chen, Dan Su, Lianwu Chen, Meng Yu, Yanmin Qian, and Dong Yu, "Deep extractor network for target speaker recovery from single channel speech mixtures," in *INTERSPEECH*, 2018.

[9] Marc Delcroix and et al., "Single channel target speaker extraction and recognition with speaker beam," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 5554–5558.

[10] Quan Wang and et al., "VoiceFilter: Targeted Voice Separation by Speaker-Conditioned Spectrogram Masking," in *Proc. Interspeech 2019*, 2019, pp. 2728–2732.

[11] Quan Wang and et al., "VoiceFilter-Lite: Streaming Targeted Voice Separation for On-Device Speech Recognition," in *Proc. Interspeech 2020*, 2020, pp. 2677–2681.

[12] Juan M. Coria and et al., "A comparison of metric learning loss functions for end-to-end speaker verification," in *Statistical Language and Speech Processing*, Luis Espinosa-Anke, Carlos Martín-Vide, and Irena Spasić, Eds., Cham, 2020, pp. 137–148, Springer International Publishing.

[13] Yi Luo and Nima Mesgarani, "Tasnet: Time-domain audio separation network for real-time, single-channel speech separation," *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 696–700, 2018.

[14] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," in *NeurIPS 2020*, 2020.

[15] Anmol Gulati and et al., "Conformer: Convolution-augmented Transformer for Speech Recognition," in *Proc. Interspeech 2020*, 2020, pp. 5036–5040.

[16] Li Wan, Quan Wang, Alan Papir, and Ignacio Lopez Moreno, "Generalized end-to-end loss for speaker verification," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 4879–4883.

[17] Abdelrahman Mohamed and et al., "Self-supervised speech representation learning: A review," *IEEE Journal of Selected Topics in Signal Processing*, pp. 1–34, 2022.

[18] Alex Graves, "Sequence transduction with recurrent neural networks," 2012.

[19] Jon Barker, Shinji Watanabe, Emmanuel Vincent, and Jan Trmal, "The fifth 'chime' speech separation and recognition challenge: Dataset, task and baselines," *CoRR*, vol. abs/1803.10609, 2018.

[20] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, "Librispeech: An asr corpus based on public domain audio books," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2015, pp. 5206–5210.

[21] David Snyder, Guoguo Chen, and Daniel Povey, "MUSAN: A Music, Speech, and Noise Corpus," 2015, arXiv:1510.08484v1.

[22] Gordon Wichern and et al., "Wham!: Extending speech separation to noisy environments," *CoRR*, vol. abs/1907.01160, 2019.

[23] Tom Ko and et al., "A study on data augmentation of reverberant speech for robust speech recognition," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 5220–5224.

[24] Igor Szöke, Miroslav Skácel, Ladislav Mošner, Jakub Paliesek, and Jan Černocký, "Building and evaluation of a real room impulse response dataset," *IEEE Journal of Selected Topics in Signal Processing*, vol. 13, no. 4, pp. 863–876, 2019.

[25] Zhuo Chen and et al., "Continuous speech separation: Dataset and analysis," in *ICASSP 2020*, 2020, pp. 7284–7288.

[26] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever, "Robust speech recognition via large-scale weak supervision," 2022.