Fully open-sourced music source separation and speech quality enhancement systems

Abstract

This poster introduces a music source separation (MSS) system and a speech quality enhancement (SQE) system. Both systems are fully open-sourced. (1) Music source separation aims to separate different sound sources (e.g., drums, bass) from a mixture audio file. MSS has several applications such as karaoke and music remixing. Our study proposed a new model to predict complex-valued ideal ratio masks with deep ResUNet architecture and channel-wise subband features. In ISMIR 2021 music demixing (MDX) challenge, our system ByteMSS ranked 2nd in the vocal track and 5th on average score.

(2)Speech quality enhancement is a crucial topic in improving online communication or recorded speech quality. Our study proposed a new model to predict complex-valued ideal ratio masks with deep ResUNet architecture and channel-wise subband features. VoiceFixer achieves state-of-the-art results on High-Fidelity speech restoration, dereverberation, de-clipping, enhancement, and equalization in one pass. Subjective evaluations show that our system has clear advantages over baselines and achieves good restoration quality on real-world test cases.

Introduction

Music source separation (MSS)

Music source separation (MSS) is a task to separate audio mixtures into individual sources such as vocals, drums, accompaniment, etc. MSS is an important topic for music information retrieval (MIR) since it can be used for several downstream MIR tasks including melody extraction, pitch estimation, , music transcription, music remixing, and so on. MSS also has several direct applications such as Karaoke and music remixing.

In this study, we propose a new MSS model, channel-wise subband phase-aware ResUNet (CWS-PResUNet), to decompose signals into subbands and estimate an unbound complex ideal ratio mask (cIRM) for each source. By combining CWS-PResUNet and Demucs, our ByteMSS system ranks the 2nd on vocals score and 5th on average score in the 2021 ISMIR Music Demixing (MDX) Challenge limited training data track (leaderboard A).

Speech Quality Enhancement (SQE)

Human speech often suffers from distortions such as background noise, room reverberations, or clipping from low-quality devices. Those distortions degrade the perceptual quality of human listeners. Speech restoration is a task to restore degraded speech to high-quality speech, which is useful in a wide range of applications such as online meeting and hearing aids.

In this work we introduce VoiceFixer, a unified framework for highfidelity speech quality enhancement. VoiceFixer restores speech from multiple distortions~(e.g., noise, reverberation, and clipping) and can expand degraded speech~(e.g., noisy speech) with a low bandwidth to 44.1 kHz full-bandwidth high-fidelity speech. We design VoiceFixer based on (1) an analysis stage that predicts intermediate-level features from the degraded speech, and (2) a synthesis stage that generates waveform using a neural vocoder. Both objective and subjective evaluations show that VoiceFixer is effective on severely degraded speech, such as real-world historical speech recordings.

Methodology of MSS

ConvBlock



As is shown in Figure 7, the model is a symmetric architecture containing a down-sampling and an upsampling path with skipconnections between the same level. It accepts magnitude spectrogram as input and estimates mask estimation, phase variations, and direct magnitude prediction. The complex spectrogram can be reconstructed with the equation on the left.

Figure 7: Proposed MSS model (CWS-PResUNet)

 $\hat{S}' = \mathsf{relu}(|X'| \odot \mathsf{sigmoid}(\hat{M}) + \hat{Q}) \exp^{j(\angle X' + \angle \hat{ heta})}$

Haohe Liu¹

Doctoral Research Student

¹Centre for Vision, Speech and Signal Processing (CVSSP), University of Surrey Email: <u>hl01486@surrey.ac.uk;</u> ORCID: 0000-0003-1036-7888



ctices				
PyPI package	es we built for the py	/thon users. [14, 15, 1	6]	Our MSS sys
n projects	Q		🕑 haoheliu 👻	margin and ISMIR MDX
0.0.17			✓ <u>Latest version</u>	Table 1 : Eva benchmark
cefixer 🗳			Released: Nov 6, 2021	Models
no videos. [13]	Figure 6: We and datasets	e open sourced our s. [5, 6, 7]	r pre-trained mod	X-UMX D3Net Demucs CWS-PResUNet ByteMSS
recordings.	 VoiceFixer Model Protection Protection<	Al CheckPoint	New version 12 12 12 13 14 15 16 17 17 18 18 19 19 10 10 10 10 11 12 15 16 17 18 19 10 10 10 11 12 13 14 14 15 16 17 18 19 10 10 10 10 11 12 13 14 15 16 17 18 19 10 10 10 10 10 11 12 14 15 16 17 18 19 19 10	VoiceFixer a enhanceme attentions f has been de Figure 10: S on high-fide benchmark Target - Oracle - VoiceFixer - Baseline-UNet - 1
does he imply? dges. hut is a kind of food. hut is delicious. Bilateral Anterior STG Left Posterior STG eft Inferior Frontal Gyrus Speech estimation Waveform Uwaveform	In speech restoration breaks the convent two-stage system: The first equation of VoiceFixer where a into a representation denotes the synthe synthesize z to the two-stage processin human perception Compared to the convert waveform, VoiceFixes spectrogram as z, we restoring multiple of	on, our proposed Va ional one-stage system $f: x \mapsto z,$ $g: z \mapsto \hat{s}.$ denotes the analysist distorted speech x on z . The second equivalent of speech of VoiceFix restored speech. The ng, VoiceFixer mine of speech describes onventional speech ate on spectrogram for uses the low diminant which alleviates the distortions simultar	biceFixer stem into a stage of is mapped uation ixer, which hrough the ics the ed in Figure 8.	[1] <u>https://gith</u> [2] <u>https://gith</u> [3] <u>https://gith</u> [4] <u>https://gith</u> [5] <u>https://zen</u> [6] <u>https://zen</u> [7] <u>https://zen</u> [7] <u>https://zen</u> [8] <u>https://zen</u> [8] <u>https://zen</u> [8] <u>https://col</u> <u>2BaeAcGie</u> [9] <u>https://col</u> <u>uWPIMCMf</u> [10] <u>https://col</u> <u>uWPIMCMf</u> [10] <u>https://py</u> [13] <u>https://py</u> [15] <u>https://py</u> [16] <u>https://py</u>





Results

rstem (ByteMSS) outperforms SOTA methods by a large ranks the 2nd in vocal separation performance in 2021 challenge.

• 02 kuielab

luation result on MSS K MUSDB18HQ

Drums	Bass	Other	Average	
6.47	5.43	4.64	5.79	
7.01	5.25	4.53	6.01	
6.57	6.53	5.14	6.28	
6.38	5.93	5.84	6.77	
6.57	6.53	5.84	6.97	
	Drums 6.47 7.01 6.57 6.38 6.57	DrumsBass6.475.437.015.25 6.576.53 6.385.936.576.53	DrumsBassOther6.475.434.647.015.254.53 6.576.53 5.146.385.93 5.84 6.576.535.84	DrumsBassOtherAverage6.475.434.645.797.015.254.536.01 6.576.53 5.146.286.385.93 5.84 6.776.576.535.84 6.97

ass	Other	Average							
.43	4.64	5.79	• 03	Music_AI	6.882	7.273	7.371	5.091	7.79
.25	4.53	0.01		👪 Kazane_Ryo_n					
.53	5.14	6.28	• 04		6.649	6.993	7.018	4.901	7.68
.93	5.84	6.77							
.53	5.84	6.97	▲ 05	ByteMSS	6.514	6.602	6.545	4.830	8.07
A performance on multiple speech quality									

achieves SOTA performance on multiple speech quality ent benchmarks. Our open-sourced tools receive a lot of from the speech researcher. The voicefixer package on PyPI ownloaded for more than 13k times in the last six months.

Subjective evaluation result elity speech restoration





DR (Song) SDR (Bass) SDR (Drums) SDR (Other) SDR (Vocals

Relevant links

nub.com/haoheliu/voicefixer_main

- hub.com/haoheliu/voicefixer
- hub.com/haoheliu/2021-ISMIR-MSS-Challenge-CWS-PResUNet <u>hub.com/haoheliu/torchsubband</u>
- nodo.org/record/5546723#.YkRLiprMKAk
- nodo.org/record/5600188#.YkRLpJrMKAk
- nodo.org/record/5175846
- ab.research.google.com/drive/1E2yJLWN8MH6GJUw15cj490E

ab.research.google.com/drive/1HYYUepIsI2aXsdET6P_AmNVX

aoheliu.github.io/nvsr

- oheliu.github.io/voicefixer
- uggingface.co/spaces/akhalig/VoiceFixer/tree/main
- ww.youtube.com/watch?v=d_j8UKTZ7J8
- pi.org/project/voicefixer/0.0.17/
- /pi.org/project/torchsubband/
- <u>/pi.org/project/ssr-eval/</u>