

SUB-BAND AND FULL-BAND INTERACTIVE U-NET WITH DPRNN FOR DEMIXING CROSS-TALK STEREO MUSIC

Han Yin¹, Mou Wang², Jisheng Bai^{1,3}, Dongyuan Shi³, Woon-Seng Gan³, Jianfeng Chen¹

¹ School of Marine Science and Technology, Northwestern Polytechnical University, Xi'an, China

² Institute of Acoustics, Chinese Academy of Sciences, Beijing, China

³ School of Electrical & Electronic Engineering, Nanyang Technological University, Singapore

ABSTRACT

This paper presents a detailed description of our proposed methods for the ICASSP 2024 Cadenza Challenge. Experimental results show that the proposed system can achieve better performance than official baselines.

Index Terms— Stereo music demixing, sub-band, full-band, U-Net, DPRNN

1. INTRODUCTION

How to provide hearing-impaired people with a better listening experience is a challenging problem. One approach is to develop a signal processing system that allows a personalized rebalancing of the music. The ICASSP 2024 Cadenza Challenge [1] encourages participants to separate music into different tracks (vocals, bass, drums and other), and then intelligently remix the tracks in a personalized manner to enhance the listening experience for people with hearing loss.

In this paper, we propose a sub-band and full-band interactive U-Net model to demix the cross-talk stereo music. The model extracts dynamic audio information from the sub-bands and the full-band respectively, and then uses convolutional layers to adaptively fuse the embeddings from different branches. Furthermore, a neural beamformer composed of multilayer perceptrons (MLPs) is applied to attenuate the effect of cross-talk.

2. METHODS

2.1. Proposed Demixing Model

Fig. 1 shows the overall pipeline of the proposed demixing model. Short-time Fourier transform (STFT) is applied on the original music to get the full-band spectrogram, denoted as $\mathbf{X} \in \mathbb{C}^{C \times T \times F}$, where C , T and F are the number of channels, frames and frequency points respectively. The full band is divided into two predefined sub-bands, denoted as $\mathbf{X}_1 \in \mathbb{C}^{C \times T \times F_1}$ and $\mathbf{X}_2 \in \mathbb{C}^{C \times T \times F_2}$, where F_1 and F_2 are predefined sub-band ranges.

\mathbf{X} , \mathbf{X}_1 and \mathbf{X}_2 are fed into three U-Nets with dual-path recurrent neural network (DPRNN) [2] to extract dynamic audio information, an interactive block is then applied to fuse the masks generated by different branches. Furthermore, MLPs in the neural beamformer are used to adaptively adjust the phase of the estimated complex-valued mask to mitigate the effect of cross-talk on the listening experience.

2.1.1. U-Net with DPRNN

As shown in Fig. 2, this module consists of encoder blocks, a DPRNN block and decoder blocks. We use two-dimensional convolutional layers in the encoder block to extract features from the original spectrogram. The DPRNN block, which consists of two bi-directional LSTMs with a residual connection, is used to alternately model the time and frequency information. In decoder blocks, we use two-dimensional deconvolution layers to restore the time-frequency domain mask.

2.1.2. Interactive Block

As shown in Fig. 1, the sub-bands and the full band are separately passed through a UNet with DPRNN to generate the corresponding masks, denoted as $\mathbf{H}_1 \in \mathbb{R}^{4C \times T \times F}$, $\mathbf{H}_2 \in \mathbb{R}^{4C \times T \times F_1}$ and $\mathbf{H} \in \mathbb{R}^{4C \times T \times F_2}$.

In the proposed interactive block, the full band mask is firstly divided into two parts, as formulated in:

$$\begin{aligned} \mathbf{H}^a &= \mathbf{H}[:, :, : F_1 + F_2] \\ \mathbf{H}^b &= \mathbf{H}[:, :, F_1 + F_2 :] \end{aligned} \quad (1)$$

Then, \mathbf{H}_1 , \mathbf{H}_2 and \mathbf{H}^a are concatenated to produce a new sub-band mask $\mathbf{S} \in \mathbb{R}^{8C \times T \times (F_1 + F_2)}$, which is fed into a convolutional layer, resulting in $\mathbf{S}_1 \in \mathbb{R}^{4C \times T \times (F_1 + F_2)}$. \mathbf{H}^b and \mathbf{S}_1 are concatenated to generate a new full-band mask, which is passed through two convolutional layers to estimate the real and imaginary parts of the complex-valued mask. Finally, we pass the real part and the imaginary part of the estimated mask through an MLP in the frequency domain respectively for adaptive adjustment in the neural beamformer. The kernel size in all convolutional layers is set to 1.

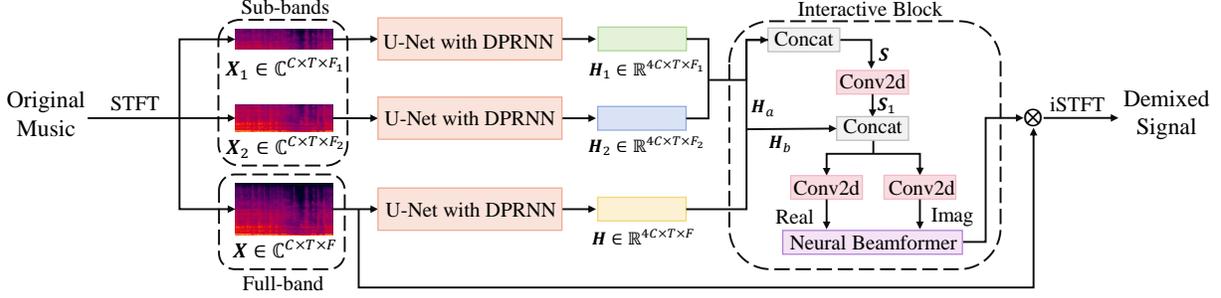


Fig. 1. The architecture of the proposed sub-band and full-band interactive U-Net with DPRNN for demixing cross-talk stereo music

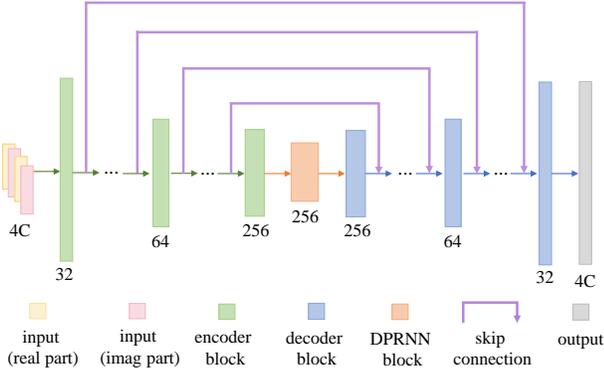


Fig. 2. The architecture of the U-Net with DPRNN

2.2. Loss Function

In the proposed system, the loss function is formulated as:

$$Loss = |\mathbf{y}, \hat{\mathbf{y}}| + |\mathbf{Y}, \hat{\mathbf{Y}}| + |\text{real}(\mathbf{Y}), \text{real}(\hat{\mathbf{Y}})| + |\text{imag}(\mathbf{Y}), \text{imag}(\hat{\mathbf{Y}})| \quad (2)$$

where \mathbf{y} and $\hat{\mathbf{y}}$ are the estimated signal and reference signal respectively, \mathbf{Y} and $\hat{\mathbf{Y}}$ are corresponding spectrograms.

3. EXPERIMENTS AND DISCUSSION

Adam optimizer is used and the initial learning rate is set to 0.001. Batch size is set to 2 and the number of gradient accumulations is 6. The parameter size of the proposed model is 152.4 MB. Computational resources and configurations of the sub-bands are presented in Table 1.

Table 2 shows the HAAQI results of different systems on the evaluation set and the validation set respectively. For each track, we train the proposed model multiple times and obtain the top 3 models. ‘‘Ensemble’’ is the weighted average result of the three models, while no ensemble represents the result generated by the best model. Results show that the proposed model can achieve HAAQI scores comparable to HDmucs [3], and the ensemble can further improve the performance.

Table 1. Configuration of sub-bands and corresponding GPU memories for training (Input audio: 4s, FFT size: 2048, Hop size: 600)

Track	Sub-band 1		Sub-band 2		Training GPU memory
	Start	End	Start	End	
Vocals	0 kHz	4 kHz	4 kHz	10 kHz	9.52 GB
Drums	0 kHz	6 kHz	6 kHz	10 kHz	9.49 GB
Bass	0 kHz	1 kHz	1 kHz	6 kHz	9.23 GB
Other	0 kHz	7 kHz	7 kHz	11 kHz	9.41 GB

Table 2. HAAQI Results of different models on the evaluation set and validation set

System	Evaluation Set			Validation Set		
	Left	Right	Overall	Left	Right	Overall
OpenUnmix [4]	0.507	0.515	0.511	0.599	0.594	0.596
HDmucs [3]	0.566	0.574	0.570	0.669	0.667	0.668
Proposed	0.564	0.572	0.568	0.667	0.664	0.665
Proposed(ensemble)	0.581	0.589	0.585	0.686	0.682	0.684

4. CONCLUSIONS

This paper presents our proposed demixing system based on U-Net and DPRNN in the ICASSP 2024 Cadenza Challenge. Experimental results show that the proposed system outperforms OpenUnmix [4] and HDmucs [3], achieving an overall HAAQI of 0.585 on the evaluation set.

5. REFERENCES

- [1] G. Roa Dabike et al., ‘‘Overview and results of the icassp cadenza challenge: music demixing/remixing for hearing aids,’’ in *IEEE Int. Conf. Acoust. Speech and Signal Process. (ICASSP)*, 2024.
- [2] Y. Luo et al., ‘‘Dual-path rnn: efficient long sequence modeling for time-domain single-channel speech separation,’’ in *IEEE Int. Conf. Acoust. Speech and Signal Process. (ICASSP)*, 2020, pp. 46–50.
- [3] A. D efossez, ‘‘Hybrid spectrogram and waveform source separation,’’ *arXiv preprint arXiv:2111.03600*, 2021.
- [4] F. R. St oter et al., ‘‘Open-unmix-a reference implementation for music source separation,’’ *J. Open Source Softw.*, vol. 4, no. 41, pp. 1667, 2019.