



The NWPU-ByteAudio System for CHiME-7 Task 2 UDASE Challenge

Zihan Zhang^{1,2}, Runduo Han¹, Ziqian Wang¹, Xianjun Xia², Yijian Xiao², Lei Xie^{1*}

¹Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science,
Northwestern Polytechnical University, Xi'an, China

²ByteDance, China

zhzhang@mail.nwpu.edu.cn, lxie@nwpu.edu.cn

Abstract

This paper describes the NWPU-ByteAudio system for CHiME-7 Task 2 - unsupervised domain adaptation for conversational speech enhancement (UDASE). To better make use of the in-domain mixture data, we improve the self-supervised learning (SSL) approach RemixIT with MetricGAN discriminator, resulting in an updated version called *RemixIT-G*. Under the RemixIT-G framework, we take Uformer+ as the speech enhancement model, which is an improved version of Uformer updated with the MetricGAN discriminator as well. We also apply an unsupervised noise adaptation model to generate noisy speech in the target domain. A perceptual contrast stretching (PCS) method is used to further improve the auditory perception quality of the enhanced speech. Our approach has achieved an SI-SDR of 12.95 and an OVRL-MOS of 3.07 in the CHiME-7 task 2 evaluation set and ranked the 1st place in the challenge.

Index Terms: Speech enhancement, unsupervised domain adaptation, RemixIT, MetricGAN

1. Introduction

In recent years, deep neural network (DNN) based speech enhancement has achieved superior performance over the traditional signal processing based methods [1]. However, the widely adopted supervised learning of neural models requires a large number of paired data for training. Since it is not feasible to capture paired noisy and clean speech in real-world scenarios, noisy speech is usually simulated by mixing clean speech and noise. This leads to a mismatch between the simulated training data and the real-world data in real applications [2]. *Unsupervised domain adaptation* has been recently proposed to solve this problem [2, 3, 4, 5].

Wisdom et al. proposed a completely unsupervised method called MixIT [2], which separate the mixtures into a variable number of sources and remix it to approximate the original mixtures. Tzinis et al. proposed a continuous self-training scheme in which a pre-trained teacher model on out-of-domain data infers estimated pseudo-target signals for in-domain mixtures [4]. Saijo and Ogawa further improved the framework by proposing Self-Remixing, which does not change the number of sources in inputs, thus alleviates the MixIT's mismatch problem [5]. However, these methods are difficult to achieve the results of supervised training.

The CHiME-7 Challenge unsupervised domain adaptation for conversational speech enhancement (UDASE) task aims to use unlabeled data to overcome the performance drop caused by domain mismatch for speech enhancement models trained on simulated data. In other words, it focuses on improving neural

speech enhancement models with the help of in-domain unlabeled data and out-domain labeled data.

In this challenge, we submitted a *self-supervised learning* (SSL) approach based on RemixIT [4]. RemixIT follows a continuous teacher-student learning scheme, in which a teacher model, trained with out-of-domain data, infers estimated pseudo-target signals for in-domain mixtures. By permuting the estimated clean and noise signals and remixing them together, a new set of bootstrapped mixtures and corresponding pseudo-targets are generated to train the student network. The teacher periodically refines its estimates using the latest student models. Moreover, we explore the efficacy of MetricGAN [6] discriminator in our approach. Specifically, we use Uformer [7] as the backbone of the teacher network but update it with MetricGAN+ [8] for better generalization capabilities. Similarly, we improve RemixIT by incorporating MetricGAN-U [9] during the training of the student network. To get better speech and noise estimates from the pre-trained teacher using unlabeled data, we adopt unsupervised noise adaptation by data simulation [10]. Moreover, to make use of the in-domain noise for model training, a voice activity detection (VAD) module [11] is used to automatically extract the noise segments from the CHiME-5 training set. Lastly, to further enhance the auditory quality of the enhanced speech, we apply perceptual contrast stretching (PCS) [12] during training and decoding.

2. Proposed approach

We first introduce our neural speech enhancement model and then describe the improved self-supervised learning scheme. Next, unsupervised noise adaptation and perceptual contrast stretching are introduced. Finally, we introduce the loss function used.

2.1. Uformer+: improved Uformer with MetricGAN

We choose *Uformer* as our enhancement model. Uformer [7] is a Unet-based dilated dual-path conformer network working in both complex and magnitude domains for simultaneous speech enhancement and dereverberation. As shown in Fig. 1, Uformer has two distinct branches – the magnitude branch and the complex branch. The primary focus of the magnitude branch is the suppression of noise, while the complex branch serves as an auxiliary module to compensate for the possible loss of spectral details and the phase mismatch. Time attention (TA) and dilated convolution (DC) are used to extract local and global contextual information. The frequency attention (FA) is used to model dimensional information. These three sub-modules and the two branches modeling effectively improve the speech enhancement performance.

In the challenge, we apply MetricGAN [6] in the training of the Uformer model to improve its generalization ability. The

*: Corresponding author.

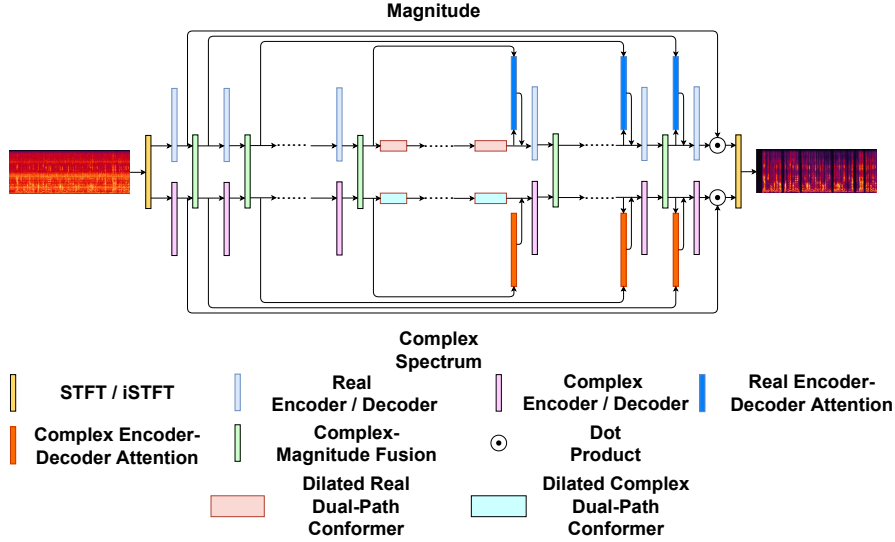


Figure 1: The overall architecture of Uformer [7].

main idea of MetricGAN is to mimic the behavior of a target evaluation function such as PESQ [13] with a neural network and thus such a non-differentiable metric can be used during model training. Specifically, we use MetricGAN+ [8], an improved version, to predict PESQ or DNS-MOS [14]. We call our updated model *Uformer+*.

2.2. RemixIT-G: improved RemixIT with MetricGAN

RemixIT [4] is a self-supervised learning (SSL) strategy based on pseudo-labeling and continual training. It remixes the clean speech and noise estimated by a pre-trained teacher model to obtain the pseudo-labels which are then used to train a student model. The teacher model is constantly updated with the weight of the student model during training.

We use *Uformer+* introduced in Section 2.1 as the teacher model, which is pre-trained on the out-of-domain data. To further improve performance, we add the MetricGAN-U [9] discriminator in *RemixIT* to learn the remixed pseudo-labels. This updated version is named as *RemixIT-G*.

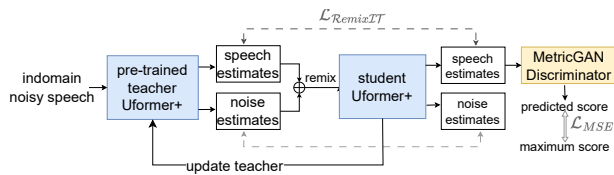


Figure 2: *RemixIT-G*: improved *RemixIT* with *MetricGAN*.

The structure of *RemixIT-G* is shown in Fig. 2. To predict PESQ without clean labels, we pre-train a *MetricGAN* discriminator on out-of-domain data. In fact, the quality net is trained simultaneously with the teacher model. The enhanced speech of the student model is fed into the discriminator which predicts PESQ. We calculate the loss with the maximum score, which can make the student model to focus on the improvement of perceived speech quality to obtain better performance.

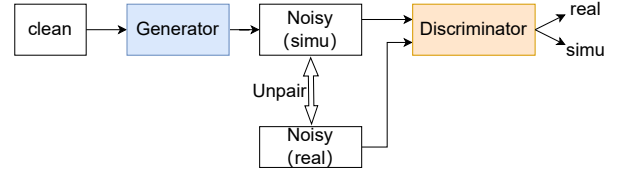


Figure 3: The structure of *UNA-GAN*.

2.3. Unsupervised noise adaptation

We use a data-simulation-based method called *UNA-GAN* [10] to generate noisy speech in the target domain. The structure of *UNA-GAN* is shown in Fig. 3. The generator is used to map clean speech to noisy speech. The goal of the generator is to integrate only the target noise features into the clean magnitude without changing the original clean speech. The discriminator determines whether the input is real or simulated.

In this design, the generator learns to incorporate noise, which matches the target domain, into clean speech. To prevent the generator from generating too much noise and overwriting the clean speech, contrast learning is used to maximize the mutual information between the paired clean and noisy magnitude spectra. Specifically, we sample 256 blocks in the simulated noisy spectrum and the same 256 blocks in the same position of the clean spectrum are also sampled. Blocks in corresponding positions are viewed as positive examples, and other mismatch pairs are treated as negative examples. Such selected blocks are reshaped via two linear layers with 256 units followed by the ReLU activation [15]. Besides the input magnitude, the same sampling operation is performed in the feature layer. Finally, the positive and negative training examples are used to calculate the cross entropy loss.

2.4. Perceptual contrast stretching

Perceptual contrast stretching (PCS) [12] is derived based on the critical band importance function. It stretches the contrast of target features in the data based on a set of auditory weights. The weights are designed according to the critical band impor-

tance [16]. As a pre-processing step, PCS is applied to the input noisy speech during training. Likewise, PCS is also applied to the enhanced speech after inference as a post-processing step.

2.5. Loss function

Our loss function consists of two parts during training the Uformer+ and the RemixIT-G student model, the speech loss and the GAN loss. To get a better speech enhancement performance, we use a speech loss that combines scale-invariant source-to-noise ratio (SI-SNR), time domain mean absolute error (MAE) and frequency domain mean squared error (MSE). The SI-SNR loss is defined as Eq. 1.

$$\mathcal{L}_{SI-SNR} = -20 \log_{10} \frac{\|\mathbf{s}\|^2}{\|\mathbf{s} - \hat{\mathbf{s}}\|^2} \quad (1)$$

The MAE loss is defined as Eq. 2.

$$\mathcal{L}_{MAE} = \sum_t |\mathbf{s}(t) - \hat{\mathbf{s}}(t)| \quad (2)$$

The MSE loss contains the magnitude MSE and the complex MSE, which is shown as Eq. 3 and Eq. 4.

$$\mathcal{L}_{mag} = \sum_{t,f} \left| |\mathbf{S}(t, f)| - |\hat{\mathbf{S}}(t, f)| \right|^2 \quad (3)$$

$$\mathcal{L}_{cplx} = \sum_{t,f} \left| \mathbf{S}(t, f) - \hat{\mathbf{S}}(t, f) \right|^2 \quad (4)$$

The final speech loss is shown as Eq. 5.

$$\mathcal{L}_{speech} = 5 * \mathcal{L}_{SI-SNR} + \mathcal{L}_{MAE} + \mathcal{L}_{mag} + \mathcal{L}_{cplx} \quad (5)$$

For the GAN loss, we use the same loss as MetricGAN-U [9], the predicted metric are normalized between 0 and 1 and try to maximum score. The total loss is defined as Eq. 6.

$$\mathcal{L} = \mathcal{L}_{speech} + \mathcal{L}_{G(MetricGAN)} \quad (6)$$

3. Experiments

3.1. Dataset

For supervised training, we use speech signals from LibriSpeech [17] to generate clean labels. To cope with the possible 1-3 speakers, we choose to generate clean labels for 1, 2, and 3 speakers in the ratio of 70%, 20%, and 10%, respectively. Noise signals are from WHAM! [18] noise dataset and CHiME-5 train set. We adopt a voice activity detection (VAD) model [11]¹ to automatically extract the in-domain noise from the unlabeled CHiME-5 train set.

We use an online data generation strategy, randomly selecting a signal-to-noise ratio (SNR) to combine clean label and noise signal, while SNR randomly ranges from 0dB to 25 dB. A total of 50,000 room impulse responses (RIRs) are generated using the HYB method [19], with random room size 5×3×3 to 8×5×5 and RT60 0.2-1s. Totally 30% of the clean speech is convoluted with RIRs to simulate reverbed signals.

The other part of the paired noisy-clean data is generated by UNA-GAN, which is used to finetune the teacher model. We use CHiME-5 train set as unpaired noisy speech to train the UNA-GAN, thus making it capable of generating in-domain noise, even if the noise overlaps with human speaking and cannot be simply separated using VAD. And for the RemixIT-G training step, we only use the unlabeled CHiME-5 train set to adapt the student.

¹<https://github.com/marianne-m/brouhaha-vad>

Table 1: Results on reverberant LibriCHiME-5 dev set

Model	Predicted metric	SI-SDR (dB)
Unprocessed	-	6.57
Sudo rm-rf	-	8.23
Uformer+	PESQ	8.83
Uformer+	DNS-MOS	8.66

Table 2: Results on reverberant LibriCHiME-5 eval set

Model	Training strategy	SI-SDR (dB)
Unprocessed	-	6.59
Sudo rm-rf	fully-supervised	7.80
Sudo rm-rf	RemixIT	9.44
Sudo rm-rf	RemixIT + VAD	10.05
Uformer+	fully-supervised	8.79
Uformer+	UNA-GAN finetune	9.37
Uformer+	RemixIT	12.04
Uformer+	RemixIT-G	12.95

3.2. Experimental setup

For the Uformer model, we use a 512-point short time Fourier transform (STFT) with a window size of 32 ms and a window shift of 16 ms. The channel numbers for the encoder and decoder layers are [8, 16, 32, 64, 128, 128]. For each 2D convolutional layers (Conv2D), the kernel size and stride are set to (2, 5) and (1, 2). Same as [20], we use a power compression for magnitude, and the compression variable is 0.5. For the MetricGAN discriminator, we use the same structures, a convolutional neural network (CNN) consisting of 4 Conv2D layers with 15 filters and a kernel size of (5, 5). The discriminator is used as Quality-Net to predict PESQ. The parameters of the discriminator are only updated when training on the labeled out-of-domain data.

For the UNA-GAN, the generator contains symmetric Conv2D layers with kernel size (3, 3) for downsampling and upsampling, respectively. The ResNet block [21] is repeated 9 times to learn the depth representation, each ResNet block consists of two Conv2D layers with kernel size (3, 3) and a dropout layer. After ResNet block are 3 self-attention layers [22], which can capture the global feature of the speech signals. The discriminator of UNA-GAN aim to determine whether the input noisy signal is real or simulated. It consists of 5 Conv2D layers with kernel size (4, 4) and each Conv2D layer was followed by a LeakyReLU activation function. The stride for the first three Conv2D layers is (3, 3) and for the last two Conv2D layers is (1,1).

3.3. Experimental results

To prove the performance of Uformer+ and select the optimal training metric for MetricGAN, we first trained Uformer+ using the LibriMix dataset, following the same data setup as the baseline fully-supervised Sudo rm-rf model [23]. We test the SI-SDR [24] on the reverberant LibriCHiME-5 dev set, the results are shown in Table 1. Compared to the Sudo rm-rf, Uformer+ obtains an absolute improvement of 0.6dB in SI-SDR. For MetricGAN, it appears that the optimal training metric is PESQ.

Table 2 shows the SI-SDR metric on the reverberant LibriCHiME-5 eval set and Table 3 shows the DNSMOS met-

Table 3: Results on CHiME-5 eval subset

Model	Training strategy	OVRL	BAK	SIG
Unprocessed	-	2.84	2.92	3.48
Sudo rm-rf	fully-supervised	2.88	3.59	3.33
Sudo rm-rf	RemixIT	2.82	3.64	3.26
Sudo rm-rf	RemixIT + VAD	2.84	3.62	3.28
Uformer+	fully-supervised	3.03	3.88	3.35
Uformer+	UNA-GAN finetune	3.05	3.91	3.36
Uformer+	RemixIT	3.04	3.94	3.37
Uformer+	RemixIT-G	3.07	3.93	3.39

ric on the CHiME-5 eval subset. It is worth mentioning that Uformer is a causal model with 9.46M parameters. We have the following conclusions. First, Uformer+ has a better performance than Sudo rm-rf. The increase of the parameters and the use of MetricGAN have brought benefits. Second, UNA-GAN is capable of learning a clean-to-noisy transformation and adapting the speech enhancement model to the target noise by simulated data. Third, RemixIT brings a large improvement on SI-SDR. Through our improvements, RemixIT-G can also improve the score of DNSMOS with the help of MetricGAN-U.

4. Conclusions

In our approach designed for the UDASE task, we improved the Uformer speech enhancement model and the RemixIT framework with MetricGAN. We have experimentally shown that the improvements are effective. We also explored domain adaptation methods using generative models to simulate in-domain data, and clear gains were achieved. The results show that by combining various domain adaptation methods, significantly better results than supervised methods can be obtained on unlabeled out-of-domain datasets.

5. References

- [1] P. Ochieng, “Deep neural network techniques for monaural speech enhancement: State of the art analysis,” *arXiv preprint arXiv:2212.00369*, 2022.
- [2] S. Wisdom, E. Tzinis, H. Erdogan, R. Weiss, K. Wilson, and J. Hershey, “Unsupervised sound separation using mixture invariant training,” *Advances in Neural Information Processing Systems*, 2020.
- [3] K. Saijo and T. Ogawa, “Remix-cycle-consistent learning on adversarially learned separator for accurate and stable unsupervised speech separation,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022.
- [4] E. Tzinis, Y. Adi, V. K. Ithapu, B. Xu, P. Smaragdis, and A. Kumar, “Remixit: Continual self-training of speech enhancement models via bootstrapped remixing,” *IEEE Journal of Selected Topics in Signal Processing*, 2022.
- [5] K. Saijo and T. Ogawa, “Self-remixing: Unsupervised speech separation via separation and remixing,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023.
- [6] S.-W. Fu, C.-F. Liao, Y. Tsao, and S.-D. Lin, “Metricgan: Generative adversarial networks based black-box metric scores optimization for speech enhancement,” in *International Conference on Machine Learning*. PMLR, 2019.
- [7] Y. Fu, Y. Liu, J. Li, D. Luo, S. Lv, Y. Jv, and L. Xie, “Uformer: A unet based dilated complex & real dual-path conformer network for simultaneous speech enhancement and dereverberation,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022.
- [8] S.-W. Fu, C. Yu, T.-A. Hsieh, P. Plantinga, M. Ravanelli, X. Lu, and Y. Tsao, “Metricgan+: An improved version of metricgan for speech enhancement,” *arXiv preprint arXiv:2104.03538*, 2021.
- [9] S.-W. Fu, C. Yu, K.-H. Hung, M. Ravanelli, and Y. Tsao, “Metricgan-u: Unsupervised speech enhancement/dereverberation based only on noisy/reverberated speech,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022.
- [10] C. Chen, Y. Hu, H. Zou, L. Sun, and E. S. Chng, “Unsupervised noise adaptation using data simulation,” in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023.
- [11] M. Lavechin, M. Métais, H. Titeux, A. Boissonnet, J. Copet, M. Rivière, E. Bergelson, A. Cristia, E. Dupoux, and H. Bredin, “Brouhaha: multi-task training for voice activity detection, speech-to-noise ratio, and c50 room acoustics estimation,” *arXiv preprint arXiv:2210.13248*, 2022.
- [12] R. Chao, C. Yu, S.-W. Fu, X. Lu, and Y. Tsao, “Perceptual contrast stretching on target feature for speech enhancement,” *arXiv preprint arXiv:2203.17152*, 2022.
- [13] I.-T. Recommendation, “Perceptual evaluation of speech quality (pesq): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs,” *Rec. ITU-T P. 862*, 2001.
- [14] C. K. Reddy, V. Gopal, and R. Cutler, “Dnsmos p. 835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022.
- [15] C. Chen, N. Hou, Y. Hu, S. Shirol, and E. S. Chng, “Noise-robust speech recognition with 10 minutes unparalleled in-domain data,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022.
- [16] C. Pavlovic, “Sii—speech intelligibility index standard: Ansi s3. 5 1997,” *the Journal of the Acoustical Society of America*, 2018.
- [17] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: an asr corpus based on public domain audio books,” in *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2015.
- [18] G. Wichern, J. Antognini, M. Flynn, L. R. Zhu, E. McQuinn, D. Crow, E. Manilow, and J. L. Roux, “Wham!: Extending speech separation to noisy environments,” *arXiv preprint arXiv:1907.01160*, 2019.
- [19] E. Bezzam, R. Scheibler, C. Cadoux, and T. Gisselbrecht, “A study on more realistic room simulation for far-field keyword spotting,” in *APSIPA ASC*. IEEE, 2020.
- [20] A. Li, C. Zheng, R. Peng, and X. Li, “On the importance of power compression and phase estimation in monaural speech dereverberation,” *JASA express letters*, 2021.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- [22] C. Chen, N. Hou, D. Ma, and E. S. Chng, “Time domain speech enhancement with attentive multi-scale approach,” in *2021 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*. IEEE, 2021.
- [23] E. Tzinis, Z. Wang, and P. Smaragdis, “Sudo rm-rf: Efficient networks for universal audio source separation,” in *2020 IEEE 30th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2020.
- [24] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, “Sdr—half-baked or well done?” in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019.