

Boosting Objective Scores of Speech Enhancement Model through MetricGAN Post-Processing

Szu-Wei Fu^{1,2*}, Chien-Feng Liao^{2*}, Tsun-An Hsieh^{2*}, Kuo-Hsuan Hung^{2*}, Syu-Siang Wang²,
Cheng Yu², Heng-Cheng Kuo², Ryandhimas E. Zezario^{1,2}, You-Jin Li², Shang-Yi Chuang²,
Yen-Ju Lu², Yu Tsao²

¹ Department of Computer Science and Information Engineering, National Taiwan University, Taipei, Taiwan

² Research Center for Information Technology Innovation, Academia Sinica, Taipei, Taiwan

{jasonfu, yu.tsao}@citi.sinica.edu.tw

Abstract

The Transformer architecture has shown its superior ability than recurrent neural networks on many different natural language processing applications. Therefore, this study applies a modified Transformer on the speech enhancement task. Specifically, the positional encoding may not be necessary and hence is replaced by convolutional layers. To further improve PESQ scores of enhanced speech, the L_1 pre-trained Transformer is fine-tuned by MetricGAN framework. The proposed MetricGAN can be treated as a general post-processing module to further boost interested objective scores. The experiments are conducted using the data sets provided by the organizer of the Deep Noise Suppression (DNS) challenge. Experimental results demonstrate that the proposed system outperforms the challenge baseline in both subjective and objective evaluation with a large margin.

Index Terms: speech enhancement, PESQ, MetricGAN, DNS challenge

1. Introduction

For commercial speech related applications, such as automatic speech recognition (ASR), hearing aids systems, and VoIP services heavily rely on clear sound provided by robust speech enhancement (SE) systems [1-4]. An SE system aims to reduce the background noises from noisy speech signals and further improve the quality and the intelligibility of enhanced ones. Traditional approaches utilize the statistical attributes of speech signals to perform enhancement under many circumstances. However, these approaches require certain premises. For instance, a widely used denoise approach, Wiener filtering [5], performs well in many conditions, but the input must be guaranteed as a stationary process, which may not be fulfilled in real world situation.

Deep learning algorithms are known for their powerful capability of learning transformations, thus the learned features are usually more representative than the handcrafted ones. In recent years, deep learning algorithms are incorporated into the SE task. Some approaches [6-10] conduct SE on time-frequency acoustic features provided by short time Fourier transform (STFT). Lu et al. [3] used a deep denoising auto-encoder (DDAE) to directly estimate the

enhanced speech. A main drawback of DDAE is that global time information is not considered because it merely depends on the frames nearby to predict an enhanced frame, not regarding the entire sequence. To avoid from this problem, some approaches operate SE in sequential modeling manner. Weninger et al. [11] and Maas et al. [12] utilize recurrent neural network (RNN) for SE system and further improve the robustness of ASR systems. RNNs such as long short-term memory (LSTM) and gated recurrent unit (GRU) handles multiple gates and uses hidden states to capture correlation within a sequence. However, because of sequential processing, it is difficult for RNNs to learn long-range dependencies between symbols. In addition, the computation of gates is inefficient due to time dependency between each other.

To solve these problems, Transformer model is proposed and has been shown to achieve state-of-the-art results in various natural language processing tasks [13]. Specifically, to model long-range dependencies, the sequence relation between all time-steps is learned by the attention mechanism [14] which can be parallel computed. Unlike RNNs, the Transformer processes an input sequence in parallel, which can significantly increase training and inference efficiency. However, Kim et al. [15] find that original Transformer does not show improvements in the task of speech enhancement. They hence propose Gaussian-weighted self-attention and surpass LSTM-based model. In this study, we find that the positional encoding in Transformer may not be necessary for SE and hence is replaced by convolutional layers.

To further boost the objective scores of speech enhanced by the modified Transformer model, we apply it as the generator of the previously proposed MetricGAN [16]. Through some training techniques proposed by [17], the MetricGAN framework is used as a post-processing module. Specifically, the SE model (generator) is first pre-trained by conventional loss function (e.g., L_1 or L_2 loss) until converge, then the surrogate loss from the discriminator further guides the generator training to achieve a better solution. Because previous researches [16, 17] have already successfully applied BLSTM and convolutional-BLSTM as the generator of MetricGAN, respectively; this framework can be treated as a general module to improve the performance of a trained deep SE model.

**Equal Contribution*

2. Transformer Model for Causal Speech Enhancement

In this paper, Transformer is developed as the backbone architecture of our SE system, some modifications have been made to fit the denoising task. Also, causal setting is adopted in order to achieve the real-time processing requirement. First, the original Transformer consists of encoder and decoder networks for sequence-to-sequence learning. For the SE task, since the input sequence length is identical to the output sequence, the decoder part is omitted in our system. Second modification is that, in order to inject some kind of relative location information to the frames in the sequence, causal convolutional layers are utilized instead of the original positional encoding. There are many choices of positional encodings, learned and fixed [18]. For the learned ones, it requires the input sequence to be fixed length, thus unable to adapt to sequence lengths that are longer than the ones encountered during training. For the fixed ones, although it may allow the model to adapt to variable sequence length, the hand-crafted fixed feature may not be rich enough for the embedding. Hence, for both flexibility and model capacity, convolutional layers are chosen for capturing location information. Finally, a future masking mechanism is applied to the scaled dot-product attention in the Multi-Head Self-Attention (MHSA) layer for causality, where the attention weights are set to zero for all future frames. More formally, three linear layers transform the input argument of MHSA into queries $Q_h \in \mathbb{R}^{T \times d_k}$, keys $K_h \in \mathbb{R}^{T \times d_k}$, and values $V_h \in \mathbb{R}^{T \times d_k}$, where T , h and d_k are the sequence length, head index, and the feature dimension, respectively. Then the masked scaled dot-product is computed as

$$\text{Attention}(Q_h, K_h, V_h) = \text{softmax}\left(M + \frac{Q_h K_h^T}{\sqrt{d_k}}\right) V_h \quad (1)$$

M is the future masking, which is an upper triangular matrix where the upper entries are set to negative infinity (excluding main diagonal), i.e.

$$M = \begin{bmatrix} 0 & \dots & -\infty \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix} \quad (2)$$

In this way, future frames will not be considered since upper entries of the attention weight become zero after the softmax function.

The rest of the architecture is implemented as a standard Transformer shown in Figure 1, which is composed of N attention blocks. In each attention block, the first sub-layer is the masked MHSA, next is a feed forward network with two fully-connected layers. Both sub-layers are followed by a residual connection to the input and a layer normalization. Herein, moments for layer normalization are computed across the channel dimension only, thus obeying the causal setting. Finally, the Transformer output is projected back to frequency dimension using a fully-connected layer with ReLU activation, and the L_1 loss is computed with the clean speech.

2.1. Implementation

The speech waveforms were recorded at 16 kHz sampling rate. The short time Fourier transform (STFT) with a hamming window size of 32 ms and a hop size of 16 ms were applied to transform the speech waveforms into 257-points spectral features. In the preliminary experiments we found that

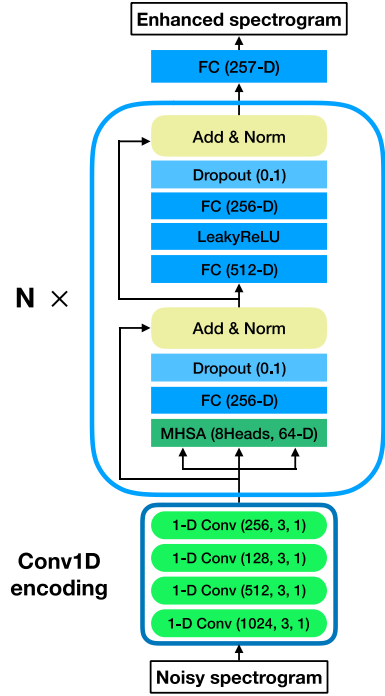


Figure 1: The proposed Transformer architecture with 1-D convolutional encoding. 1-D Conv is in the format (output channels, kernel size, stride) and FC (output channels) denotes the fully-connected layer. Add & Norm is the residual connection followed by layer normalization. Finally, each MHSA layer consists of 8 heads and 64 dimension per head.

compressing each coefficient to a tighter range produces better results, hence the \log_{1p} function ($\log_{1p}(x) = \log(1+x)$) was adopted on the magnitude spectrogram. During testing, the enhanced spectral features were synthesized back to the waveform signals via the inverse STFT and an overlap-add procedure. The phases of the noisy signals were used for the waveform generation. The Adam optimizer was used with learning rate of $5e-5$, and an early stopping was performed based on the validation set to prevent overfitting.

3. Fine-tuning the Enhancement Model by MetricGAN

Because the evaluation of this challenge is based on the ITU-T P.808 subjective evaluation of speech quality, a quality-related loss function may be a good choice for training the speech enhancement model. However, most of the quality metrics (e.g. PESQ) are too complicated to be directly applied as an objective function. To solve this problem, [19] employed a deep model (called Quality-Net [20]) to learn the behavior of PESQ function. Then, Quality-Net is served as a surrogate of PESQ function to guide the learning of the enhancement model. Although Quality-Net loss can further improve the PESQ score of enhanced speech, the gradient provided by Quality-Net is only accurate for the first few learning iteration. In other words, Quality-Net is easily fooled (estimated quality scores increase but true scores decrease) [19].

To solve this problem, [16] proposed a learning framework such that Quality-Net and enhancement model are alternatively updated. This method is called MetricGAN, because its goal is for black-box metric scores optimization

and the architecture is similar to generative adversarial networks (GAN). In the following, we briefly introduce the training of MetricGAN:

Here, we first introduce a function $Q'(I)$ to represent the normalized evaluation metric (between 0 and 1) to be optimized, where I is the input of the metric. For example, for PESQ and STOI, I is the pair of enhanced speech $G(x)$ that we want to evaluate and the corresponding clean speech y . Therefore, to ensure that discriminator network (D) behaves similar to Q' , the objective function of D is:

$$L_{D(\text{MetricGAN})} = \mathbb{E}_{x,y}[(D(y,y) - 1)^2 + (D(G(x),y) - Q'(G(x),y))^2] \quad (3)$$

where $0 \leq Q'(G(x),y) \leq 1$.

The training of generator network (G) can completely rely on the adversarial loss:

$$L_{G(\text{MetricGAN})} = \mathbb{E}_x[(D(G(x),y) - s)^2] \quad (4)$$

where s is the desired assigned score. For example, to generate clean speech, we can simply assign s to be 1.

Although the original MetricGAN is trained from scratch [16], Kawanaka *et al.* [17] proposed some training techniques to make MetricGAN a post-processing method of a trained speech enhancement model. After fine-tuning by the surrogate loss in MetricGAN, the interested metric scores can be stably further improved. In this study, we apply some of the tricks to fine-tune the Transformer model (pre-trained by L_1 loss) for further boosting the PESQ scores. The overall flow chart is shown in Figure 2.

4. Experiments

4.1. Dataset

The dataset used in this experiment is provided by Deep Noise Suppression Challenge [21]. The default configuration is used to generate noisy-clean paired speech data. To reduce the training time, we only randomly choose 10000 training utterances to train our model. The synthetic test set without reverberation is selected as the validation set to evaluate the performance of different models. The subjective speech quality evaluation is based on the blind test set.

4.2. Model Structure

The pre-trained Transformer described in Section 2 is used as the SE model in this experiment. The parameters are first pre-trained with L_1 -based signal approximation (SA) [11] loss. The discriminator (Quality-Net) herein is a CNN with four two-dimensional (2-D) convolutional layers with the number of filters and kernel size as follows: [15, (5, 5)], [25, (7, 7)], [40, (9, 9)], and [50, (11, 11)]. To handle the variable-length input, a 2-D global average pooling layer was added, so that the features were fixed with 50 dimensions (50 is the number of feature maps in the previous layer). Three fully connected layers were added subsequently, each with 50 and 10 LeakyReLU nodes, and 1 linear node, respectively. To make Quality-Net a smooth function, we constrained it to be 1-Lipschitz continuous by spectral normalization [22].

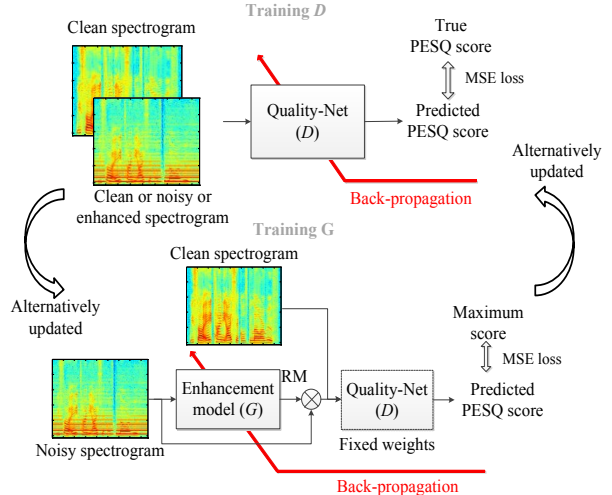


Figure 2: Flow chart of MetricGAN training.

Table 1: Performance comparisons of different models on the synthetic test set without reverberation.

	PESQ	STOI
Noisy	2.454	0.915
NSNet [23]	2.692	0.906
Transformer (PE)	2.429	0.894
Proposed Transformer	2.966	0.932
Transformer+MetricGAN	3.104	0.946

4.3. Experimental Results of Objective Evaluation

To verify the effectiveness of the proposed framework, the standard PESQ function was used to measure the speech quality and the score ranges from -0.5 to 4.5. We also presented STOI [24] for speech intelligibility evaluation and the score ranges from 0 to 1. Both the two metrics are the higher the better. Table 1 presents the results of the average PESQ and STOI scores on the validation set for the Noise Suppression Net (NSNet) [8] baseline and proposed method, which fine-tunes the parameters of Transformer model by the MetricGAN framework. As shown in Table 1, the significant performance drop in the Transformer with additive fixed positional encoding, denoted as Transformer (PE), echoes our hypothesis that the Transformer requires a better mechanism to inject location information. And the performance of the proposed Transformer is much better than NSNet. When we pre-trained the Transformer model with the L_1 loss and subsequently post-processing by MetricGAN, we could further improve both PESQ and STOI scores with a large margin. Note that, because we only fine-tune the parameters, the computation load and model size keeps unchanged.

Figure 3 shows the fine-tuning process of proposed MetricGAN. The PESQ score roughly converge to 2.97 by using the L_1 loss to train the Transformer model. When the surrogate loss in MetricGAN is applied, the score can further be improved by 0.13.

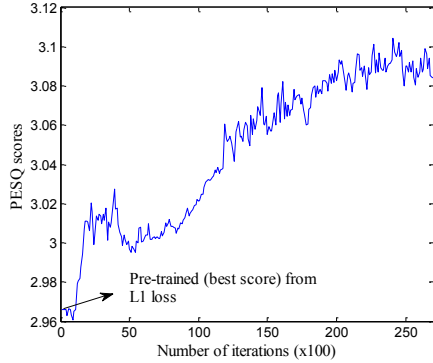


Figure 3: PESQ scores of fine-tuning process by MetricGAN framework.

Table 2: Computational complexity of proposed model.

	Number of parameters	Inference time (ms/frame)
Proposed Transformer	5,953,920	0.256

4.4. Computational Complexity

In this section, we report the computational complexity in terms of number of parameters and the time it takes to infer a frame. As shown in the table 2, the number of weights in the proposed Transformer model is roughly 5.9M, and it takes 0.256 millisecond (ms) to process a frame of 32ms long (this is based on the average processing time of the whole blind test set) by an Intel Core i5 CPU quad core machine clocked at 2.4 GHz.

4.5. Spectrogram Comparison

Next, we present the spectrograms of a clean utterance in the synthetic test set, the same utterance corrupted by traffic noise, enhanced speeches by Transformer with L_1 loss and fine-tuned by the proposed MetricGAN in Fig. 4. From Fig. 4(c), we observe that although L_1 loss can guide Transformer effectively remove the background noise, some noise still exists (as shown in the blue dashed rectangle). Comparing Fig. 4(c) and (d), we can find that the remaining noise can be further removed by the MetricGAN post-processing.

4.6. Subjective Evaluation

The organizer of DNS challenge conducted P.808 subjective evaluation of the submitted enhanced speech. 10 qualified judges rated each clip, which resulted in 95% confidence interval (CI) of about 0.02 on the overall Mean Opinion Scores (MOS). The blind test set can be further split into noisy speech without reverberation (noreverb), noisy real recordings (realrec) and noisy reverberant speech (reverb) categories. Table 3 shows the MOS scores of noisy, NSNet baseline and the proposed Transformer model fine-tuned by MetricGAN. From this table we can observe that our proposed model can significantly outperform the baseline.

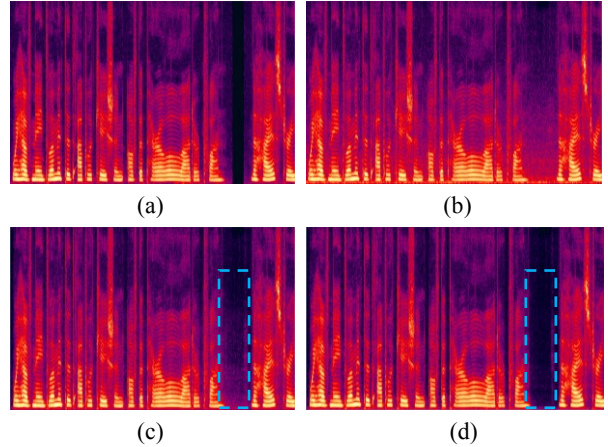


Figure 4: Spectrograms of an utterance in the synthetic test set: (a) clean speech, (b) noisy speech (traffic noise), (c) enhanced speech by Transformer with L_1 loss (d) enhanced speech by Transformer+MetricGAN.

Table 3: MOS scores of the blind test set.

	noreverb	realrec	reverb	Overall
Noisy	3.32	2.97	2.78	3.01
NSNet [23]	3.49	3.00	2.64	3.03
Transformer+ MetricGAN	3.63	3.18	2.83	3.21

5. Discussion

The proposed MetricGAN fine-tuning framework can be treated as a universal post-processing for speech enhancement model. For example, although we apply a Transformer as the generator (enhancement model) in this study, it also works for other models such as BLSTM [16] and convolutional BLSTM [17]. In addition to T-F mask estimation; this method can also improve the objective scores of mapping-based enhancement model. However, to avoid generating additional artifact (we observe that it may generate some high frequency noise when we use PESQ or STOI function as Q' . This may be because these two functions [24] ignore the signal difference in high frequency range.), we suggest this post-processing is better applied on mask estimation based deep model.

6. Conclusion

In this paper, we apply a modified Transformer model as the generator of MetricGAN. To further boost the interested objective scores, the Transformer model is first trained with conventional L_1 loss, and then fine-tuned by the surrogate loss provided by the discriminator. Experimental results demonstrate that the proposed framework outperformed the challenge baseline in terms of both objective scores and subjective evaluation. Through MetricGAN, we anticipate that the mismatch between the human auditory perception and the loss used in training a speech enhancement model can be effectively reduced.

7. References

- [1] P. C. Loizou, *Speech enhancement: theory and practice*: CRC press, 2013.
- [2] M. Kolbæk, Z.-H. Tan, and J. Jensen, "Speech intelligibility potential of general and specialized deep neural network based speech enhancement systems," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, pp. 153-167, 2017.
- [3] X. Lu, Y. Tsao, S. Matsuda, and C. Hori, "Speech enhancement based on deep denoising autoencoder," in *Proc. Interspeech*, 2013, pp. 436-440.
- [4] M. H. Soni, N. Shah, and H. A. Patil, "Time-frequency masking-based speech enhancement using generative adversarial network," in *Proc. ICASSP*, 2018, pp. 5039-5043.
- [5] P. Scalart, "Speech enhancement based on a priori signal to noise estimation," in *Proc. ICASSP*, 1996, pp. 629-632.
- [6] L. Sun, J. Du, L.-R. Dai, and C.-H. Lee, "Multiple-target deep learning for LSTM-RNN based speech enhancement," in *Hands-free Speech Communications and Microphone Arrays (HSCMA)*, 2017, pp. 136-140.
- [7] D. S. Williamson, Y. Wang, and D. Wang, "Complex ratio masking for joint enhancement of magnitude and phase," in *Proc. ICASSP*, 2016, pp. 5220-5224.
- [8] Y. Xu, J. Du, L.-R. Dai, and C.-H. Lee, "A regression approach to speech enhancement based on deep neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing* vol. 23, pp. 7-19, 2015.
- [9] D. Wang and J. Chen, "Supervised speech separation based on deep learning: An overview," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2018.
- [10] Z. Zhao, S. Elshamy, H. Liu, and T. Fingscheidt, "A CNN postprocessor to enhance coded speech," in *Proc. IWAENC*, 2018.
- [11] F. Weninger, H. Erdogan, S. Watanabe, E. Vincent, J. Le Roux, J. R. Hershey, *et al.*, "Speech enhancement with LSTM recurrent neural networks and its application to noise-robust ASR," in *Proc. LVA/ICA*, 2015, pp. 91-99.
- [12] A. Maas, Q. V. Le, T. M. O'neil, O. Vinyals, P. Nguyen, and A. Y. Ng, "Recurrent neural networks for noise reduction in robust ASR," in *Proc. Interspeech*, 2012.
- [13] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, *et al.*, "Attention is all you need," in *Proc. NIPS*, 2017, pp. 5998-6008.
- [14] Y. Koizumi, K. Yaiabe, M. Delcroix, Y. Maxuxama, and D. Takeuchi, "Speech enhancement using self-adaptation and multi-head self-attention," in *Proc. ICASSP*, 2020, pp. 181-185.
- [15] J. Kim, M. El-Khamy, and J. Lee, "Transformer with Gaussian weighted self-attention for speech enhancement," in *Proc. ICASSP*, 2020.
- [16] 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 *Proc. ICML*, 2019.
- [17] M. Kawanaka, Y. Koizumi, R. Miyazaki, and K. Yatabe, "Stable training of DNN for speech enhancement based on perceptually-motivated black-box cost function," in *Proc. ICASSP*, 2020.
- [18] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," in *Proc. ICML*, 2017, pp. 1243-1252.
- [19] S.-W. Fu, C.-F. Liao, and Y. Tsao, "Learning with learned loss function: Speech enhancement with quality-net to improve perceptual evaluation of speech quality," *IEEE Signal Processing Letters*, 2019.
- [20] S.-W. Fu, Y. Tsao, H.-T. Hwang, and H.-M. Wang, "Quality-Net: An end-to-end non-intrusive speech quality assessment model based on blstm," in *Proc. Interspeech*, 2018.
- [21] C. K. Reddy, E. Beyrami, H. Dubey, V. Gopal, R. Cheng, R. Cutler, *et al.*, "The INTERSPEECH 2020 Deep Noise Suppression Challenge: Datasets, Subjective Speech Quality and Testing Framework," *arXiv preprint arXiv:2001.08662*, 2020.
- [22] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," *arXiv preprint arXiv:1802.05957*, 2018.
- [23] Y. Xia, S. Braun, C. K. Reddy, H. Dubey, R. Cutler, and I. Tashev, "Weighted speech distortion losses for neural-network-based real-time speech enhancement," in *Proc. ICASSP*, 2020.
- [24] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time-frequency weighted noisy speech," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, pp. 2125-2136, 2011.