

Explicit Estimation of Magnitude and Phase Spectra in Parallel for High-Quality Speech Enhancement

Ye-Xin Lu, Yang Ai, *Member, IEEE*, Zhen-Hua Ling, *Senior Member, IEEE*

Abstract—Phase information has a significant impact on speech perceptual quality and intelligibility. However, existing speech enhancement methods encounter limitations in explicit phase estimation due to the non-structural nature and wrapping characteristics of the phase, leading to a bottleneck in enhanced speech quality. To overcome the above issue, in this paper, we proposed MP-SENet, a novel Speech Enhancement Network that explicitly enhances Magnitude and Phase spectra in parallel. The proposed MP-SENet comprises a Transformer-embedded encoder-decoder architecture. The encoder aims to encode the input distorted magnitude and phase spectra into time-frequency representations, which are further fed into time-frequency Transformers for alternatively capturing time and frequency dependencies. The decoder comprises a magnitude mask decoder and a phase decoder, directly enhancing magnitude and wrapped phase spectra by incorporating a magnitude masking architecture and a phase parallel estimation architecture, respectively. Multi-level loss functions explicitly defined on the magnitude spectra, wrapped phase spectra, and short-time complex spectra are adopted to jointly train the MP-SENet model. A metric discriminator is further employed to compensate for the incomplete correlation between these losses and human auditory perception. Experimental results demonstrate that our proposed MP-SENet achieves state-of-the-art performance across multiple speech enhancement tasks, including speech denoising, dereverberation, and bandwidth extension. Compared to existing phase-aware speech enhancement methods, it further mitigates the compensation effect between the magnitude and phase by explicit phase estimation, elevating the perceptual quality of enhanced speech. Remarkably, for the speech denoising task, the proposed MP-SENet yields a PESQ of 3.60 on the VoiceBank+DEMAND dataset and 3.62 on the DNS challenge dataset.

Index Terms—Speech enhancement, encoder-decoder, magnitude prediction, phase prediction, explicit estimation.

I. INTRODUCTION

IN real-life scenarios, speech signals captured by devices are inevitably distorted by intrusive noises and their self-reflections from the surfaces of surrounding objects. For the bandwidth-limited scenario, when these captured signals are transmitted between narrowband communication devices, only their low-frequency components are typically preserved. These issues immensely degrade the performances of automatic speech recognition, hearing aids, and telecommunication systems [2], [3]. In this regard, speech enhancement (SE) aims to improve the intelligibility and overall perceptual quality of distorted speech signals. With the development of deep

learning and deep neural network (DNN), numerous DNN-based SE methods have been proposed to specifically address one or more of these issues, with the aim of recovering the original speech signals from the distorted ones.

For speech signal processing, the phase information determines the temporal and spectral characteristics of the speech signal, exerting a substantial impact on the naturalness and intelligibility of the reconstructed speech. In convolutional SE methods, restoring phase information has long been underestimated, with distorted phase being used as an optimal estimator of the original phase [4]. However, recent studies demonstrated the noteworthy contribution of precise knowledge of the phase spectrum to speech quality [5]. The phase information ensures coherence across time frames and between adjacent frequency bins in the time-frequency (TF) domain. Accurately preserving it during enhancement helps maintain the temporal characteristics and harmonic structures of the original speech.

Despite the significance of phase reconstruction in SE tasks, existing SE methods still face challenges in precisely modeling and optimizing the phase due to its non-structural nature and wrapping characteristics. Time-domain SE methods [6]–[10] utilized neural networks to learn the mapping from distorted waveforms to original ones, implicitly restored both the magnitude and phase information without explicit modeling. For the TF-domain phase-aware SE methods, some focused on enhancing the short-time complex spectrum, which implicitly recovered both magnitude and phase within the complex domain [11]–[18]. Recently, several advancements in the form of multi-stage decoupling-style methods have been proposed, which first explicitly enhanced the magnitude spectrum and subsequently conducted complex spectrum refinement to implicitly enhance the phase [19]–[22]. While both time-domain and TF-domain methods can partially recover phase information through modeling and optimizing waveforms or complex spectra, Wang *et al.* [23] pointed out that optimizing loss functions defined solely in the time or complex domain would result in implicit compensation issue between magnitude and phase. Although introducing additional magnitude-domain loss can effectively alleviate this issue [23], the lack of explicit phase optimization would still result in magnitude compensating for inaccurate and discontinuous phase prediction, undermining the harmonic structures of the magnitude and impacting the enhanced speech quality. Therefore, explicit phase modeling and optimization are particularly important for SE models to produce precise magnitude and phase with slighter compensation effect, thereby further elevating the perceptual quality and intelligibility of enhanced speech.

To this end, we propose MP-SENet, a TF-domain monaural

The authors are with the National Engineering Research Center of Speech and Language Information Processing, University of Science and Technology of China, Hefei 230027, China (e-mail: yxlu0102@mail.ustc.edu.cn; yangai@ustc.edu.cn; zhling@ustc.edu.cn).

This work is the extended version of our conference paper [1] published at Interspeech 2023.

SE model that performs parallel estimations of magnitude and phase spectra. The MP-SENet adopts an encoder-decoder architecture, connecting the encoder and decoder with Transformers [24] across both time and frequency dimensions, enabling the capture of alternating time and frequency dependencies. The encoder encodes the input distorted magnitude and wrapped phase spectra derived from short-time Fourier transform (STFT) into TF-domain representations for subsequent decoding. The parallel magnitude mask decoder and phase decoder estimate the original magnitude and phase spectra, respectively, and the enhanced speech waveforms are finally reconstructed using inverse STFT (iSTFT). To realize explicit phase enhancement, we follow our previous work [25] to utilize the parallel estimation architecture and anti-wrapping loss to explicitly model and optimize the phase spectra, respectively. Additionally, we utilize magnitude loss, complex spectral loss, and STFT consistency loss to train the MP-SENet model effectively. Furthermore, to address the incomplete correlation between these objective losses and human auditory perception, we incorporate a metric discriminator [26], [27] to help enhance the speech perceptual quality. Experimental results demonstrate that MP-SENet achieves state-of-the-art (SOTA) performance across speech denoising, dereverberation, and bandwidth extension (BWE) tasks. It is noteworthy that compared to previous methods with implicit phase optimization, our MP-SENet further alleviates the compensation effect between magnitude and phase by explicitly modeling and optimizing the phase, resulting in more accurate magnitude-phase prediction and harmonic restoration.

The main contribution of this work lies in proposing a comprehensive solution for phase-aware speech enhancement, which involves explicit modeling and optimization of both magnitude and phase. This enables the model to predict more accurate magnitude and phase with slighter compensation, thereby elevating the quality of enhanced speech to a new level. To the best of our knowledge, we are the first to achieve direct enhancement of the phase. Additionally, our proposed MP-SENet exhibits the ability to handle three distinct SE tasks, i.e., speech denoising, dereverberation, and BWE with a unified framework, showcasing its generalization capability across different SE tasks.

The rest of this paper is organized as follows. Section II briefly reviews several phase-aware DNN-based SE methods and the relevant backgrounds of three aforementioned SE tasks. In Section III, we give details of our proposed MP-SENet framework. The experimental setup is presented in Section IV, while Section V gives the results and analysis. Finally, we give conclusions in Section VI.

II. RELATED WORK

A. Phase-aware SE Methods

In traditional SE methods, only the magnitude spectrum was enhanced and then combined with the distorted phase spectrum to reconstruct the enhanced speech [28], [29]. Upon recognizing the importance of phase reconstruction in improving speech perceptual quality, plenty of phase-aware SE methods have been proposed to recover phase information

from distorted speech or spectral features. In this paper, we primarily focus on the TF-domain phase-aware SE methods and broadly classify them into three categories for discussion.

- The first category of TF-domain phase-aware SE methods enhances the phase within the complex spectrum and can be further divided into complex spectral masking-based and complex spectral mapping-based methods. Complex spectral masking [11]–[16] is a commonly used method that predicts the complex-valued mask applied to the real and imaginary parts of the distorted short-time complex spectrum to mask the distorted components. While in recent studies [17], [18], it seems that complex spectral mapping, which directly maps the distorted complex spectrum to the clean one, tends to have better performance than spectral masking, especially in low signal-to-noise ratio (SNR) conditions. For both complex spectral masking and mapping, the phase information was implicitly restored when optimizing the complex spectrum.
- The second category is the decouple-style SE methods [19], [20], [22], which decomposes the complex spectrum estimation problem into two sub-stages, including magnitude estimation and complex spectral refinement, respectively. Specifically, in the first stage, the distorted magnitude spectrum undergoes initial enhancement, followed by combining with the distorted phase spectrum to obtain a coarse complex spectrum estimation. In the second stage, an auxiliary network predicts the residual complex spectrum which is then added to the preliminary enhanced complex spectrum to obtain a more precise complex spectrum estimation. This category of methods aims to compensate for the distorted phase by the complex spectral refinement, thereby restoring the intact phase information.
- Both of the above types implicitly model the phase within the complex spectrum, while the last category of phase-aware SE methods incorporates a dual-stream architecture to explicitly model magnitude and phase [30]. Similar to the decouple-style methods, the distorted magnitude is firstly enhanced through magnitude masking in the magnitude stream, and the phase stream outputs a complex-valued feature map with two channels, corresponding to the real and the imaginary parts. Differently, the magnitude of the complex feature map is normalized to 1, and consequently, the complex feature map only contains the phase information. In this way, the phase spectrum can be explicitly modeled by the complex spectrum features.

In summary, existing TF-domain phase-aware SE methods implicitly or explicitly model the phase by predicting the complex spectrum. However, they still optimize the phase implicitly in the complex spectral domain, leading to a compensation effect [23] between magnitude and phase. Even with the introduction of magnitude-domain loss, phase optimization remains challenging, resulting in magnitude compensation for inaccurate and discontinuous phase prediction, thereby affecting its harmonic structure and further impacting speech perceptual quality. Thus, there is scope for improvement in terms of speech perceptual quality and intelligibility by undertaking both explicit phase modeling and optimization.

B. Speech Enhancement Tasks

In this paper, we mainly focus on three SE tasks, including speech denoising, dereverberation, and BWE. In this subsection, we will sequentially introduce the backgrounds and relevant methods for these three tasks.

1) *Speech denoising*: Speech denoising aims to remove the additive noise signal from the noisy speech signal. Classic signal-processing-based speech denoising methods include spectral subtraction [31], Wiener filtering [32], statistical model-based methods [33], subspace algorithms [34], [35], etc. These traditional methods often required prior information and faced challenges in dealing with non-stationary noise. With the renaissance of deep learning, DNN-based denoising methods demonstrate strong non-stationary noise suppression capabilities and achieve better SE performance. Existing DNN-based denoising methods can be classified into two categories, i.e., time-domain methods and TF-domain methods.

Time-domain denoising methods aim to predict clean speech waveforms directly from noisy ones without any frequency transformation. One of the most notable methods in this category is SEGAN [6], which employs generative adversarial networks (GANs) to learn the mapping from noisy speech waveforms to clean ones at the waveform level. While this category of methods successfully avoids the phase estimation problem encountered in TF-domain methods, they sacrifice the ability to capture speech phonetics in the frequency domain. Consequently, this trade-off can lead to artifacts in enhanced speech waveforms. As a result, the performance of time-domain denoising methods is generally considered inferior to that of TF-domain methods.

TF-domain denoising methods are predominantly based on spectral mapping or masking. Mapping-based methods map the spectral representations of degraded speech to those of clean speech, including magnitude mapping with noisy phase retained [28], [29], complex spectral mapping [17], [18], and other phase-contained spectral mappings [36]. Masking-based methods initially enhance the magnitude spectrum using the ideal binary mask (IBM) [37], ideal ratio mask (IRM) [38], [39], and spectral magnitude mask (SMM) [40] to selectively attenuate the noise components based on the mask values. Recognizing the significance of phase information, phase-sensitive mask (PSM) [41] and complex ideal ratio mask (cIRM) [42] were proposed based on SMM and IRM, respectively. PSM remains a real-valued mask with additional phase estimation, allowing the estimated magnitude to compensate for the noisy phase, while cIRM is a complex-valued mask used to directly enhance the complex spectrum. In addition, the decoupling-style methods and dual-stream magnitude and phase estimation methods we mentioned in section II-A also belong to this category of methods.

2) *Speech dereverberation*: Speech dereverberation aims to remove the convolutive effects or reflections caused by the surfaces of surrounding objects, which result in smearing effects across the time and frequency domains of the original speech waveforms. The quantity and intensity of these reflections are influenced by factors such as room size, surface properties, and microphone distance. These reflections can be effectively modeled using the room impulse response (RIR) filter.

In the early research, signal processing techniques were employed to eliminate the convolutive effects caused by reverberation via statistical modeling of the RIRs. For instance, the weighted prediction error (WPE) algorithm [43] employs long-term linear prediction to estimate the impact of RIRs on reverberant speech waveforms, which assumes that the early reflections of RIRs, along with the direct path, are crucial for improved recognition, while late reflections are considered diffused and detrimental to speech intelligibility.

DNN-based speech dereverberation methods are generally similar to the approaches used in denoising tasks, and there are models capable of simultaneously handling both denoising and dereverberation tasks. Since reverberation manifests as smearing effects in the TF-domain spectrum, IBM [44] and IRM [45] based methods predict masks applied to the reverberant magnitude spectra, combined with the unprocessed reverberant phase spectra, to reconstruct the anechoic speech waveforms. Following the speech denoising task, cIRM is further employed to implicitly reconstruct the phase [13]. Furthermore, some mapping-based methods have also achieved promising dereverberation performance [46], [47].

3) *Speech bandwidth extension*: Speech BWE aims to predict the high-frequency components from the narrowband speech signal to obtain a wideband one, thereby extending the frequency bandwidth and enhancing the quality and intelligibility of the speech signal. Therefore, it can also be regarded as an SE task. Traditional speech BWE methods utilized signal processing techniques to predict high-frequency residual signals and spectral envelopes. However, these traditional methods suffered from bottlenecks in terms of model capabilities, leading to the production of overly smoothed spectral parameters [48], whereas recent DNN-based methods have demonstrated superior performance.

DNN-based speech BWE methods can also be classified into time-domain methods, TF-domain methods, and hybrid-domain methods. Time-domain BWE methods [7], [10] directly predict the high-resolution speech waveforms from low-resolution ones. For TF-domain BWE methods, the recovery of the high-frequency phase is the major challenge. Early studies directly replicate [49] or mirror inverse [50] the narrowband phase spectrum to obtain the upper-band phase spectrum. Vocoder-based methods [51] adopt a neural vocoder to restore the wideband speech waveform from the extended mel-spectrogram, and further replace the low-frequency components of the vocoder's output with the original low-frequency speech signal. Similar to speech denoising and dereverberation, decouple-style methods [22] employ an additional complex spectral refinement to compensate for the narrowband phase, implicitly recovering the high-frequency phase. Additionally, the hybrid-domain methods [52] combine both time-domain and TF-domain information with a dual-branch architecture. One branch predicts the wideband magnitude spectrum from the narrowband one, and the other predicts the wideband waveform from the input narrowband one. Finally, the prediction of the wideband phase spectrum is derived from the wideband waveform output of the time-domain branch.

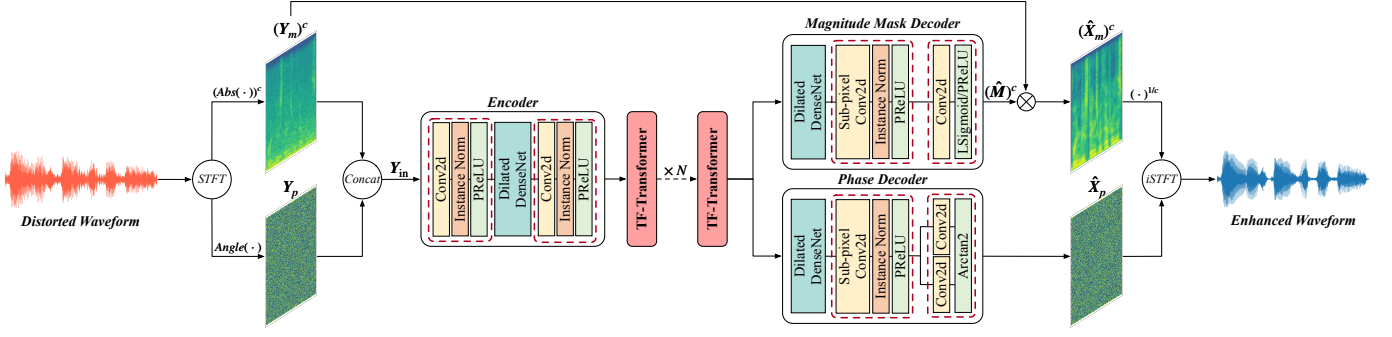


Fig. 1. Overall structure of the proposed MP-SENet for three SE tasks, with the denoising task used as an illustration. The “ $Abs(\cdot)$ ” and “ $Angle(\cdot)$ ” represent the magnitude and phase calculation functions, respectively. The “ $(\cdot)^c$ ” and “ $(\cdot)^{1/c}$ ” denote the magnitude compression and decompression operations, respectively. The “ $Concat$ ” denotes the spectral concatenation operation and the “ $Arctan2$ ” denotes the two-argument arc-tangent function.

III. METHODOLOGY

A. Model Structure

1) *Overview*: An overview of the model structure of the proposed MP-SENet is illustrated in Fig. 1. For all the speech denoising, dereverberation, and BWE tasks, they share the same architecture. Generally, the MP-SENet employs an encoder-decoder architecture to enhance the distorted speech waveform $y \in \mathbb{R}^L$ and recover the original speech waveform $x \in \mathbb{R}^L$ in the TF domain, where L represents the length of the waveform. To begin, we apply STFT on y to extract the short-time complex spectrum $Y \in \mathbb{C}^{T \times F}$, and subsequently calculate the magnitude spectrum $Y_m \in \mathbb{R}^{T \times F}$ and the wrapped phase spectrum $Y_p \in \mathbb{R}^{T \times F}$, where $Y = Y_m \cdot e^{jY_p}$, and T and F denote the total number of frames and frequency bins, respectively. Before feeding them to the encoder, we apply the power-law compression [53] on Y_m and get the compressed magnitude spectrum $(Y_m)^c \in \mathbb{R}^{T \times F}$, where c is the compression factor. The input feature $Y_{in} \in \mathbb{R}^{T \times F \times 2}$ of MP-SENet is the concatenation of the $(Y_m)^c$ and Y_p . The encoder encodes Y_{in} into a compressed TF-domain representation, and subsequently, the TF-domain representation is processed by N time-frequency transformer (TF-Transformer) blocks to alternately capture time and frequency dependencies. In the final stages, the parallel magnitude mask decoder and phase decoder predict the compressed magnitude spectrum $(\hat{X}_m)^c \in \mathbb{R}^{T \times F}$ and wrapped phase spectrum $\hat{X}_p \in \mathbb{R}^{T \times F}$ from the TF-domain representation, respectively. Subsequently, the enhanced waveform \hat{x} is reconstructed through iSTFT. Further insights into the encoder, TF-Transformer block, magnitude mask decoder, and phase decoder are elaborated below.

2) *Encoder*: As illustrated in Fig. 1, the encoder transforms the input feature Y_{in} into a TF-domain representation with C channels, T temporal dimensions, and $F' = F/2$ frequency dimensions. Specifically, the encoder comprises a cascade of a convolutional block, a dilated DenseNet [54], and another convolutional block. Each convolutional block comprises a 2D convolutional layer, an instance normalization [55], and a parametric rectified linear unit (PReLU) activation [56]. The first convolutional block increases the channel numbers of the TF-domain representation to C , and the second convolutional block halves the feature dimension F to F' to reduce the subsequent computational complexity in the TF-Transformer

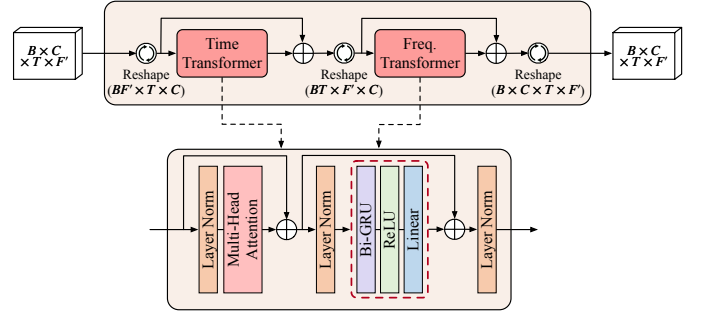


Fig. 2. The diagram of the TF-Transformer block used in Fig 1, where B , C , T , and F' represent the batch size, the number of channels, the number of frames, and the number of frequency bins, respectively. The TF-Transformer block is a cascade of a Time-Transformer and a Freq.-Transformer, both of which share the same architecture with inputs of different shapes.

blocks. The dilated DenseNet utilizes four convolutional layers with dilation sizes of 1, 2, 4, and 8 to extend the receptive field along the time axis, facilitating context aggregation at different resolutions. Additionally, it incorporates dense connections at each convolutional layer, connecting to all previous layers. This not only helps prevent the vanishing gradient problem but also enhances the parameter efficiency of the network.

3) *TF-Transformer*: The dual-path attention-based structure [22], [57], [58] has demonstrated remarkable performance in SE tasks, primarily due to its ability to capture both time and frequency dependencies. Draw inspiration from the two-stage Conformers in CMGAN [22], we proposed the TF-Transformer block, which can sequentially capture long-range correlations exhibited in the temporal and spectral domains. As illustrated in Fig. 2, the TF-Transformer block takes an intermediate TF-domain representation as input, with a shape of $B \times C \times T \times F'$, where B denotes the batch size. Initially, this representation is reshaped to $BF' \times T \times C$ and passed through the Time-Transformer layer to capture temporal dependencies. Subsequently, it is reshaped to $BT \times F' \times C$ and fed into the Frequency-Transformer layer to capture frequency dependencies. Finally, the processed TF-domain representation is reshaped back to $B \times C \times T \times F'$ as the output of the TF-Transformer block.

For each TF-Transformer block, both the Time-Transformer layer and the Frequency-Transformer layer share the same

architecture, which is a gated recurrent unit (GRU) based Transformer [57] without positional encoding. The GRU-based transformer consists of a multi-head self-attention (MHSA) [24] and a GRU-based position-wise feed-forward network (FFN), augmented with layer normalizations [59] and residual connections. The MHSA with M heads enables the model to collectively focus on the global information from different representation subspaces at various positions. The GRU-based position-wise FFN comprises a bidirectional GRU (Bi-GRU) layer, a rectified linear unit (ReLU) activation [60], and a linear layer to learn the local information [61].

4) *Magnitude mask decoder*: As illustrated in Fig 1, the magnitude mask decoder predicts a compressed magnitude mask $(\hat{M})^c \in \mathbb{R}^{T \times F}$ from the TF-domain representation processed by the TF-Transformer blocks and multiplies it with the compressed distorted magnitude spectrum $(Y_m)^c$ to obtain the enhanced compressed magnitude spectrum $(\hat{X}_m)^c$.

The magnitude mask decoder comprises a dilated DenseNet, a deconvolutional block, and a mask estimation architecture. The deconvolutional block is composed of a 2D sub-pixel convolutional layer [62], an instance normalization, and a PReLU activation. The sub-pixel convolutional layer is used to upsample the frequency dimensions of the TF-domain representation back to F while avoiding checkerboard artifacts compared to transposed convolution. The mask estimation architecture comprises a 2D convolutional layer and an activation function. The convolutional layer is used to reduce the output channel numbers of the deconvolutional block to 1, and the activation function is for mask estimation. For speech denoising or dereverberation, we use the learnable sigmoid (LSigmoid) function [27] to estimate a boundary magnitude mask:

$$\text{LSigmoid}(t) = \frac{\beta}{1 + e^{1-\alpha t}}, \quad (1)$$

where β is set to 2.0, and $\alpha \in \mathbb{R}^F$ is a trainable frequency-related parameter, which enables the model to adaptively adjust the shape of the activation function across frequency bands. For speech BWE, we use the PReLU activation to predict an unbounded high-frequency magnitude mask.

5) *Phase decoder*: Considering the phase wrapping problem, we design the phase decoder to directly predict the enhanced wrapped phase spectrum \hat{X}_p from the TF-domain representation. The phase decoder and the magnitude mask decoder share a similar architecture. To address the challenges posed by the non-structural nature and wrapping characteristics of the phase, we follow our previous work [25] and incorporate the parallel estimation architecture after a dilated DenseNet and a deconvolutional block in the phase decoder. The parallel estimation architecture first adopts two parallel 2D convolutional layers to output the pseudo-real part component and pseudo-imaginary part component and then activates these two components to obtain the enhanced wrapped phase spectrum \hat{X}_p using the two-argument arc-tangent (Arctan2) function.

B. Training Criteria

1) *Spectrum-based losses*: We first define loss functions on the magnitude spectra, phase spectra, and short-time complex

spectra to jointly train the proposed MP-SENet in the spectral domain. These loss functions are described as follows.

- **Magnitude Spectrum Loss**: The magnitude spectrum loss is defined as the mean square error (MSE) loss between the original compressed magnitude spectrum $(X_m)^c \in \mathbb{R}^{T \times F}$ and the enhanced compressed magnitude spectrum $(\hat{X}_m)^c$ to explicitly optimize the magnitude spectrum:

$$\mathcal{L}_{\text{Mag.}} = \mathbb{E}_{(X_m, \hat{X}_m)} \left[\left\| (X_m)^c - (\hat{X}_m)^c \right\|_F^2 \right], \quad (2)$$

where $\|\cdot\|_F$ represents the Frobenius norm.

- **Phase Spectrum Loss**: Considering the phase wrapping property, the absolute distance between the original phase spectra $X_p \in \mathbb{R}^{T \times F}$ and the enhanced phase spectra \hat{X}_p may not accurately reflect their actual distance, revealing the inappropriateness of conventional phase losses (e.g., absolute L^p distance) for phase optimization. Thus, to explicitly optimize the phase spectrum, we utilize the anti-wrapping function we proposed in [25] to prevent error expansion caused by phase wrapping, which is defined as:

$$f_{\text{AW}}(t) = \left| t - 2\pi \cdot \text{round}\left(\frac{t}{2\pi}\right) \right|, t \in \mathbb{R}. \quad (3)$$

On the basis of the anti-wrapping function, we first define the instantaneous phase (IP) loss between the original phase spectrum X_p and the enhanced phase spectrum \hat{X}_p :

$$\mathcal{L}_{\text{IP}} = \mathbb{E}_{(X_p, \hat{X}_p)} \left[\left\| f_{\text{AW}}(X_p - \hat{X}_p) \right\|_1 \right]. \quad (4)$$

Considering the phase continuity along both the time and frequency axis, we calculate the group delay (GD) spectrum and the instantaneous angular frequency (IAF) spectrum of the phase, which correspond to the derivatives of the phase spectrum along the frequency axis and the time axis, respectively. We further define the GD loss between the anti-wrapped original and enhanced GD spectra as well as the IAF loss between the anti-wrapped original and enhanced IAF spectra:

$$\mathcal{L}_{\text{GD}} = \mathbb{E}_{\Delta_{\text{DF}}(X_p, \hat{X}_p)} \left[\left\| f_{\text{AW}}(\Delta_{\text{DF}}(X_p - \hat{X}_p)) \right\|_1 \right], \quad (5)$$

$$\mathcal{L}_{\text{IAF}} = \mathbb{E}_{\Delta_{\text{DT}}(X_p, \hat{X}_p)} \left[\left\| f_{\text{AW}}(\Delta_{\text{DT}}(X_p - \hat{X}_p)) \right\|_1 \right], \quad (6)$$

where Δ_{DF} and Δ_{DT} represent the differential operator along the frequency axis and temporal axis, respectively. The final phase spectrum loss is the sum of IP loss, GD loss, and IAF loss:

$$\mathcal{L}_{\text{Pha.}} = \mathcal{L}_{\text{IP}} + \mathcal{L}_{\text{GD}} + \mathcal{L}_{\text{IAF}}. \quad (7)$$

- **Complex Spectrum Loss**: To further optimize the magnitude and phase spectra within the complex domain, we define the complex spectrum loss between the real and imaginary parts of the original complex spectrum $X = X_m \cdot e^{jX_p} \in \mathbb{C}^{T \times F}$ and the enhanced complex spectrum $\hat{X} = \hat{X}_m \cdot e^{j\hat{X}_p} \in \mathbb{C}^{T \times F}$ as follows:

$$\mathcal{L}_{\text{Com.}} = \mathbb{E}_{(X, \hat{X})} \left[\left\| \text{Real}(X - \hat{X}) \right\|_F^2 + \left\| \text{Imag}(X - \hat{X}) \right\|_F^2 \right], \quad (8)$$

where $\text{Real}(\cdot)$ and $\text{Imag}(\cdot)$ represent operations for extracting the real and imaginary parts from the complex spectrum, respectively.

- **STFT Consistency Loss:** The previous losses are directly defined on the magnitude spectrum, phase spectrum, and their complex combination output by the model, without considering the inconsistency introduced by the iSTFT and STFT [63], [64]. Therefore, to ensure consistency between $\hat{\mathbf{X}}$ and its consistent complex spectrum $\text{STFT}(\text{iSTFT}(\hat{\mathbf{X}}))$, we define the STFT consistency loss between the real and imaginary parts of them as follows:

$$\mathcal{L}_{\text{Con.}} = \mathbb{E}_{\hat{\mathbf{X}}} \left[\left\| \text{Real}(\hat{\mathbf{X}} - \text{STFT}(\text{iSTFT}(\hat{\mathbf{X}}))) \right\|_{\text{F}}^2 + \left\| \text{Imag}(\hat{\mathbf{X}} - \text{STFT}(\text{iSTFT}(\hat{\mathbf{X}}))) \right\|_{\text{F}}^2 \right] \quad (9)$$

2) *MetricGAN-based losses:* The primary goal of speech enhancement is to improve the quality and intelligibility of distorted speech considering human perception [65]. However, specific objective metrics associated with human perception, such as the perceptual evaluation of speech quality (PESQ) [66] and short-time objective intelligibility (STOI) [67], are non-differentiable. Thus they cannot be directly used as loss functions. To overcome this issue, MetricGAN [26] introduces a metric discriminator as a learned surrogate of the evaluation metrics. The metric discriminator iteratively estimates a surrogate loss that approaches the sophisticated metric surface, enabling the SE model to leverage this surrogate to determine the optimization direction toward achieving improved speech perceptual quality.

We here adopt the metric discriminator from [22] for adversarial training, which uses the PESQ as the target objective metric. The metric discriminator uses pairs of reference and enhanced speech waveforms as inputs and outputs the corresponding PESQ score of the enhanced speech from their magnitude. Since the original PESQ score ranges from -0.5 to 4.5, we scale them to the range of 0 to 1 with linear normalization. With the metric discriminator, we subsequently define the metric loss to guide the MP-SENet model training in generating speech waveforms corresponding to as high PESQ scores as possible. The discriminator loss and the corresponding generator metric loss are described as follows:

$$\mathcal{L}_{\text{D}} = \mathbb{E}_{\mathbf{x}} \left[\left\| D(\mathbf{x}, \mathbf{x}) - 1 \right\|_{\text{F}}^2 \right] + \mathbb{E}_{(\mathbf{x}, \hat{\mathbf{x}})} \left[\left\| D(\mathbf{x}, \hat{\mathbf{x}}) - Q_{\text{PESQ}} \right\|_{\text{F}}^2 \right], \quad (10)$$

$$\mathcal{L}_{\text{Metric}} = \mathbb{E}_{(\mathbf{x}, \hat{\mathbf{x}})} \left[\left\| D(\mathbf{x}, \hat{\mathbf{x}}) - 1 \right\|_{\text{F}}^2 \right], \quad (11)$$

where D denotes the metric discriminator and $Q_{\text{PESQ}} \in [0, 1]$ denotes the real scaled PESQ score between \mathbf{x} and $\hat{\mathbf{x}}$.

3) *The final loss for model training:* The final generator loss is the linear combination of the aforementioned spectrum-based losses and the generator metric loss:

$$\mathcal{L}_{\text{G}} = \lambda_1 \mathcal{L}_{\text{Mag.}} + \lambda_2 \mathcal{L}_{\text{Pha.}} + \lambda_3 \mathcal{L}_{\text{Com.}} + \lambda_4 \mathcal{L}_{\text{Con.}} + \lambda_5 \mathcal{L}_{\text{Metric}}, \quad (12)$$

where $\lambda_1, \lambda_2, \dots, \lambda_5$ are hyperparameters. The training criteria of the MP-SENet is to jointly minimize \mathcal{L}_{G} and \mathcal{L}_{D} .

IV. EXPERIMENTAL SETUP

A. Datasets

1) *Speech denoising:* We used the publicly available dataset [68] for the speech denoising experiments, including the VoiceBank+DEMAND dataset and the Deep Noise Suppression (DNS) Challenge 2020 dataset [69].

The VoiceBank+DEMAND dataset includes pairs of clean and noisy audio clips with a sampling rate of 48 kHz. All the audio clips were resampled to 16 kHz in the experiments. The clean audio set is selected from the Voice Bank corpus [70], which consists of 11,572 audio clips from 28 speakers for training and 824 audio clips from 2 unseen speakers for testing. The clean audio clips are mixed with 10 types of noise (8 types from the DEMAND database [71] and 2 artificial types) at SNRs of 0dB, 5dB, 10dB, and 15 dB for the training set and 5 types of unseen noise (i.e., living room, office space, bus, open area cafeteria, and public square) from the DEMAND database at SNRs of 2.5 dB, 7.5dB, 12.5 dB, and 17.5 dB for the test set.

The DNS Challenge 2020 dataset includes 500 hours of clean audio clips from 2150 speakers and over 180 hours of noise audio clips. All the clean audio clips were split into segments of 10 seconds. Following the official script¹, we generated 3,000 hours of noisy audio clips with SNRs from -5dB to 15dB for training. For model evaluation, the DNS Challenge 2020 dataset provided non-blind test sets that contained data with or without reverberation. Since we used a dedicated dataset for the dereverberation task, we here only adopted the test set without reverberation. The test set contained 4,270 audio clips spoken by 20 speakers and corresponding noisy audio clips with SNR levels uniformly sampled between 0 dB and 25 dB.

2) *Speech dereverberation:* We used the Reverb Challenge dataset [72] for the speech dereverberation experiments, exclusively employing its single-channel data for the monaural speech enhancement conducted in this study. The training set consisted of the clean WSJCAM0 [73] training set and a multi-condition training set, which was generated from the clean WSJCAM0 training data by convolving these clean utterances with 24 measured RIRs with the reverberation time (T_{60}) ranging from 0.2s to 0.8s and adding recorded background noise at a fixed SNR of 20 dB. The validation set and the test set each consisted of two different parts, namely the ‘SimData’ and the ‘RealData’. The ‘SimData’ consists of artificially distorted versions of utterances taken from the WSJCAM0 corpus. These data were created by convolving the clean WSJCAM0 signals with RIRs measured in three different rooms (small, medium, and large) with T_{60} of approximately 0.25s, 0.5s, and 0.7s, and microphone distance of a near condition (0.5m) and a far condition (2m). Consistent with the training set, recorded background noise was added to the reverberated data at a fixed SNR of 20 dB. The ‘RealData’ consists of real reverberant speech recordings from the MC-WSJ-AV dataset [74], which were captured in a noisy and reverberant meeting room with T_{60} of approximately 0.7s and microphone distance of 1m and 2.5m.

¹<https://github.com/microsoft/DNS-Challenge>.

TABLE I

EXPERIMENTAL RESULTS FOR SPEECH DENOISING METHODS EVALUATED ON THE VOICEBANK+DEMAND DATASET. “-” DENOTES THE RESULT THAT IS NOT PROVIDED IN THE ORIGINAL PAPER.

Method	Year	Structure	#Param.	WB-PESQ	CSIG	CBAK	COVL	STOI
Noisy	-	-	-	1.97	3.35	2.44	2.63	0.91
SEGAN [6]	2017	Waveform	43.2M	2.16	3.48	2.94	2.80	0.92
DEMUCS [8]	2021	Waveform	33.5M	3.07	4.31	3.40	3.63	0.95
SE-Conformer [9]	2021	Waveform	-	3.13	4.45	3.55	3.82	0.95
MetricGAN [26]	2019	Magnitude	-	2.86	3.99	3.18	3.42	-
MetricGAN+ [27]	2021	Magnitude	-	3.15	4.14	3.16	3.64	-
DPT-FSNet [15]	2021	Complex	0.88M	3.33	4.58	3.72	4.00	0.96
TridentSE [16]	2023	Complex	3.03M	3.47	4.70	3.81	4.10	0.96
DB-AIAT [21]	2021	Magnitude+Complex	2.81M	3.31	4.61	3.75	3.96	-
CMGAN [22]	2022	Magnitude+Complex	1.83M	3.41	4.63	3.94	4.12	0.96
PHASEN [30]	2020	Magnitude+Phase	20.9M	2.99	4.21	3.55	3.62	-
MP-SENet	2024	Magnitude+Phase	2.26M	3.60	4.81	3.99	4.34	0.96

3) *Speech bandwidth extension*: We used the VCTK dataset [75] which contains 44 hours of speech recordings from 108 speakers with a sampling rate of 48 kHz for the speech BWE experiments. We followed the previous works [7], [22], [52] to use the first 100 speakers for training and the remaining 8 speakers for testing. To generate the wideband-narrowband speech pairs, We first resampled all the audio clips to the sampling rate of 16 kHz as the wideband speech. We further followed the implementation of AudioUNet² to subsampled the 16 kHz wideband speech into narrowband ones with a subsampling rate, and upsampled them back to 16 kHz using the spline interpolation. We set the subsampling rate to 2 and 4 in our implementation, indicating the extensions from 8 kHz and 4 kHz to 16 kHz, respectively.

B. Model Configuration

In the training stage, all the audio clips were sliced into 2-second segments. To extract input features from raw waveforms using STFT, the FFT point number, Hanning window size, and hop size were set to 400, 400 (25 ms), and 100 (6.25 ms), and consequently, the number of frequency bins $F = 201$. For compressing the magnitude spectrum, the compression factor c was set to 0.3. For training our proposed MP-SENet, we set the batch size B , the number of channels C , the number of TF-Transformer blocks N , and the number of heads in the MHSA M to 4, 64, 4, and 4, respectively. The hyperparameters of the final generator loss $\lambda_1, \lambda_2, \dots, \lambda_5$ were set to 0.9, 0.3, 0.1, 0.1, and 0.05 with the empirical trial. All the models were trained until 500k steps using the AdamW optimizer [76], with $\beta_1 = 0.8$, $\beta_2 = 0.99$, and weight decay of 0.01. The learning rate was set initially to 0.0005 and scheduled to decay with a factor of 0.99 every epoch.³

V. RESULTS AND ANALYSIS

A. Speech Denoising

1) *Evaluation metrics*: For the evaluation on the VoiceBank+DEMAND dataset, five commonly used objective eval-

uation metrics were chosen to evaluate the enhanced speech quality, including wideband PESQ (WB-PESQ), STOI, and three composite measures [77] (i.e., CSIG, CBAK, and COVL). WB-PESQ was used to evaluate the wideband speech perceptual quality, which ranges from -0.5 to 4.5. STOI was used to measure speech intelligibility, which ranges from 0 to 1. CSIG, CBAK, and COVL measure signal distortion, background noise intrusiveness, and overall effect, respectively. All these three mean opinion score (MOS) based metrics range from 1 to 5. For the evaluation on the DNS Challenge dataset, besides the WB-PESQ and the STOI, we further adopted the narrowband PESQ (NB-PESQ) to evaluate the narrowband speech perceptual quality and the scale-invariant signal-to-distortion ratio (SI-SDR) [78] to quantify the distortion between the enhanced and clean speech signals. For all the metrics, higher values indicate better performance.

2) *Comparison with baseline methods*: We evaluated the performance of our proposed MP-SENet separately on the VoiceBank+DEMAND dataset and the large-scale DNS Challenge dataset, with different baseline models used for each of these datasets.

• **Evaluation on the VoiceBank+DEMAND Dataset**: we compared our proposed MP-SENet with several representative time-domain and TF-domain methods. In the time-domain category, we included SEGAN [6], DEMUCS [8], and SE-Conformer [9]. For the TF-domain methods, we selected MetricGAN [26], MetricGAN+ [27], PHASEN [30], and four SOTA methods, namely DPT-FSNet [15], TridentSE [16], DB-AIAT [21], and CMGAN [22]. As depicted in Table I, our proposed MP-SENet outperformed all other speech denoising methods across all metrics, demonstrating its strong denoising capability. It was noteworthy that the time-domain methods generally performed inferior to TF-domain methods, indicating the importance of TF-domain characteristics for the speech denoising task. Within the TF-domain methods, while our proposed MP-SENet and PHASEN both used magnitude and phase spectra as input conditions, MP-SENet exhibited significant improvements of 0.61, 0.60, 0.44, and 0.72 in WB-PESQ, CSIG, CBAK, and COVL scores with a nine times smaller model size

²<https://github.com/kuleshov/audio-super-res>.

³Audio samples of the proposed MP-SENet can be accessed at the official website <https://yxl-0102.github.io/MP-SENet>.

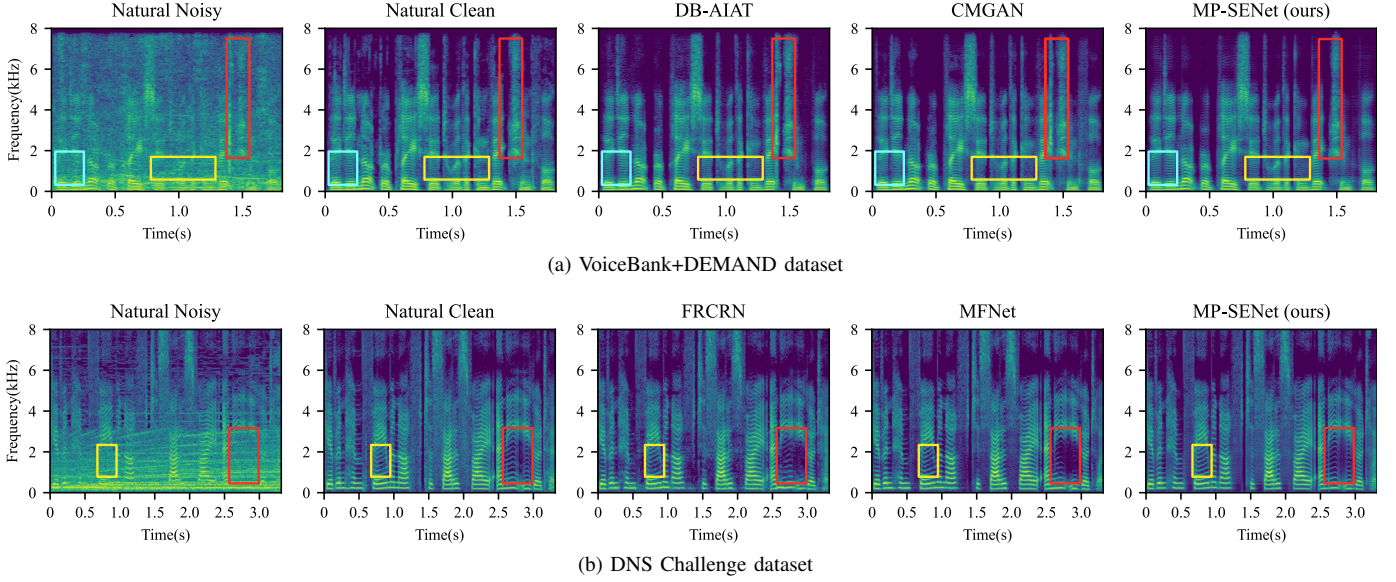


Fig. 3. Spectrogram visualization of the noisy speech, clean speech, and speech waveforms denoised by SOTA baseline methods and our proposed MP-SENet.

TABLE II
EXPERIMENTAL RESULTS FOR SPEECH DENOISING METHODS EVALUATED ON THE DNS CHALLENGE EVALUATION SET *w/o* REVERBERATION.

Method	Year	#Param.	WB-PESQ	NB-PESQ	STOI (%)	SI-SDR (dB)
Noisy	-	-	1.58	2.45	91.52	9.07
FullSubNet [12]	2021	5.64M	2.78	3.31	96.11	17.29
CTSNNet [19]	2021	4.35M	2.94	3.42	96.21	16.69
TaylorSENet [20]	2022	5.40M	3.22	3.59	97.36	19.15
FRCRN [14]	2022	6.90M	3.23	3.60	97.69	19.78
MFNet [36]	2023	-	3.43	3.74	97.78	20.31
MP-SENet	2024	2.26M	3.62	3.92	98.16	21.03

compared to PHASEN. Compared to the other four TF-domain SOTA approaches, our proposed MP-SENet with explicit phase optimization still excelled among these metrics. Especially, the higher WB-PESQ score signified better magnitude spectra generation, indicating that our proposed MP-SENet with explicit phase optimization can further alleviate the compensation effect between magnitude and phase, thereby enhancing the speech perceptual quality.

For more evidence, we visualized the spectrograms of speech waveforms enhanced by DB-AIAT, CMGAN, and our proposed MP-SENet as shown in Fig. 3a. By comparing the contents of boxes with the same color, it is evident that our proposed MP-SENet achieved better denoising performance compared to DB-AIAT and CMGAN, and effectively preserved the harmonic structures of the speech. This further confirmed the effectiveness of explicit phase optimization in promoting the generation of high-fidelity magnitude spectra.

- **Evaluation on the DNS Challenge Dataset:** We further compared our proposed MP-SENet with several TF-domain speech denoising methods on the DNS Challenge dataset, including two complex spectral masking-based methods (i.e., FullSubNet [12] and FRCRN [14]), two decoupling-style methods (i.e., CTSNet [19] and TaylorSENet [20]), and a short-time discrete cosine transform (STDCT) spectrum-

based method (i.e., MFNet). As depicted in Table II, our proposed MP-SENet significantly outperformed these baseline methods in terms of all the metrics with the smallest parameter count. Especially, our MP-SENet demonstrated huge advantages in PESQ metrics, which reflected the improvement in speech perceptual quality. Compared to the STFT spectrum-based methods, our proposed MP-SENet achieved superior performance due to the explicit estimation of the phase. Although MFNet achieved implicit magnitude-phase recovery by predicting a real-valued STDCT spectrum, its performance in waveform generation tasks was still inferior to that of the STFT spectrum, which may be attributed to the fact that an over-complete Fourier basis contributed to improved training stability [79].

We also visualized the spectrograms of speech waveforms enhanced by our proposed MP-SENet and the other two SOTA baselines, including FRCRN and MFNet. As shown in Fig. 3b, we can observe that although the noise distorted the fundamental frequency and the low-order harmonics in the noisy speech, both the baseline models and our MP-SENet can effectively restore them, which could be because the energy of them was relatively high, making them easier to be separated from the noisy components. However, as shown in the yellow and red boxes in the figure, when the higher-order harmonics were disrupted, the baseline models, lacking explicit phase optimization, failed to adequately restore the harmonic structures during enhancement. This further confirmed the importance of accurate phase estimation for the speech denoising task.

3) *SNR-wise Analysis Experiment:* Paliwal *et al.* [5] emphasized the importance of phase information for speech denoising tasks in low-SNR circumstances. Theoretically, The noisy phase \mathbf{Y}_p can be derived as follows [80]:

$$\begin{aligned} \mathbf{Y}_p &= \text{Angle}(\mathbf{X}_m e^{j\mathbf{X}_p} + \mathbf{N}_m e^{j\mathbf{N}_p}) \\ &= \mathbf{X}_p + \text{Angle}\left(1 + \frac{\mathbf{N}_m}{\mathbf{X}_m} e^{j(\mathbf{N}_p - \mathbf{X}_p)}\right), \end{aligned} \quad (13)$$

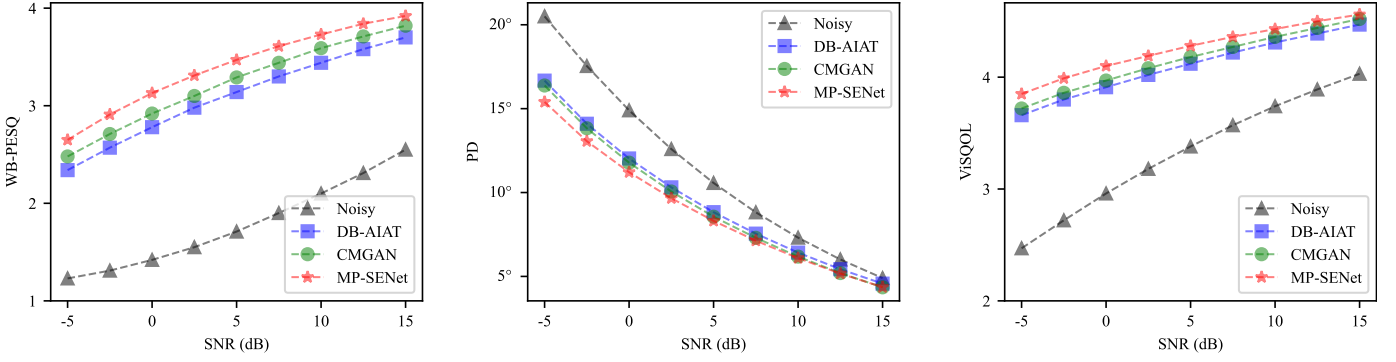


Fig. 4. WB-PESQ, PD, and ViSQOL metrics of the noisy speech and speech waveforms enhanced by DB-AIAT, CMGAN, and our proposed MP-SENet on the re-noised Voice Bank test set with varying SNR conditions.

where $N_m \in \mathbb{R}^{T \times F}$ and $N_p \in \mathbb{R}^{T \times F}$ denote the magnitude and phase spectra of the noise signal $n \in \mathbb{R}^L$, respectively. $\text{Angle}(\cdot)$ calculates the phase spectrum of a complex spectrum. When in a high-SNR circumstance, the clean magnitude spectrum X_m is much larger than the noise magnitude spectrum N_m , the second term in Eq. 13 would approach zero, indicating that the clean phase spectrum X_p can be approximated by the noisy phase spectrum Y_p . So we hypothesized that as the SNR increases, the contribution of phase information to the performance of the SE model would diminish.

To experimentally verify our hypothesis, we conducted an analysis experiment in which we tested the performance of our proposed MP-SENet under different SNR circumstances. Additionally, we selected DB-AIAT and CMGAN, which implicitly enhance the phase, as comparative systems to validate the importance of explicit phase enhancement. With the clean test set of the Voice Bank corpus, we re-noised it with five types of unseen noise (i.e., living room, office space, bus, open area cafeteria, and a public square) under separate SNR conditions, which range from -5 dB to 15 dB with intervals of 2.5 dB. Subsequently, we employed these test sets with varying SNRs to evaluate the performance of DB-AIAT, CMGAN, and our proposed MP-SENet. To quantitatively measure the phase estimation, we used the phase distance (PD) metric [11] which was defined as the average angle difference weighted by target magnitude between the enhanced speech and the target speech in the TF domain. The PD can be formulated using our proposed anti-wrapping function as follows:

$$PD = \frac{360^\circ}{\pi} \sum_{t,f} \frac{X_m[t,f]}{\sum_{t',f'} X_m[t',f']} f_{AW}(X_p[t,f] - \hat{X}_p'[t,f]), \quad (14)$$

where \hat{X}_p' denotes the re-extracted phase spectrum from the enhanced speech \hat{x} . The PD metric ranges from 0° to 180° , the lower the better. Furthermore, since WB-PESQ is a magnitude-related metric that can not reflect the overall speech quality, we adopted the virtual speech quality objective listener (ViSQOL) [81], which utilizes a spectral-temporal measure of similarity between reference and test speech signals to produce a mean opinion score - listening quality objective (MOS-LQO) score. The ViSQOL metric ranges from 1 to 5, the higher the better.

The experimental results were plotted into curves as illustrated in Fig. 4. Firstly, it is obvious that our proposed

MP-SENet outperformed the other two baselines among the three metrics, demonstrating its superiority in refining both magnitude and phase and improving overall speech quality. For the WB-PESQ metric, both the baseline models and MP-SENet exhibited relatively stable improvements compared to the noisy speech across different SNRs. For the PD metric, the enhanced speech showed a significant reduction in phase distance compared to the noisy speech at very low SNRs. However, as the SNR increased, the gap in the PD metric between the enhanced speech and the noisy speech diminished and eventually tended toward zero. At the same time, the difference in the ViSQOL metric between them was also narrowing, which aligned with our hypothesis and highlighted the significance of recovering phase information in low-SNR circumstances. It is noteworthy that, the two baseline models with implicit phase optimization, almost overlapped in PD and ViSQOL. However, our MP-SENet showed apparent improvements in these two metrics in low-SNR conditions, which confirmed the contribution of explicit phase estimation to the overall enhanced speech quality.

4) *Ablation Studies:* To investigate the role of each key component in the MP-SENet, we took the speech denoising task as the ablation example and performed ablation studies on the VoiceBank+DEMAND dataset. We used the WB-PESQ, PD, SI-SDR, and ViSQOL metrics to evaluate the magnitude restoration, phase restoration, distortion of the enhanced waveform, and overall speech quality, respectively. The ablation results are presented in Table. III.

We first conducted ablation studies on the model structure. To compare the modeling capability between the Transformers and the Conformers used in CMGAN, we replaced the Transformers in the TF-Transformer blocks with Conformers (denoted as “w/ Conformer”). As a result, the ablation of the Transformers resulted in performance degradation, which may be attributed to the fact that Transformers have better modeling capabilities than Conformers. However, it was noteworthy that our MP-SENet with Conformers still achieved a PESQ of 3.56 with a smaller model size of 1.71M, which significantly outperformed those of CMGAN in Fig. I (PESQ of 3.41 and model size of 1.83M). This parallel comparison highlighted the contribution of our proposed magnitude-phase parallel prediction architecture and explicit phase optimization.

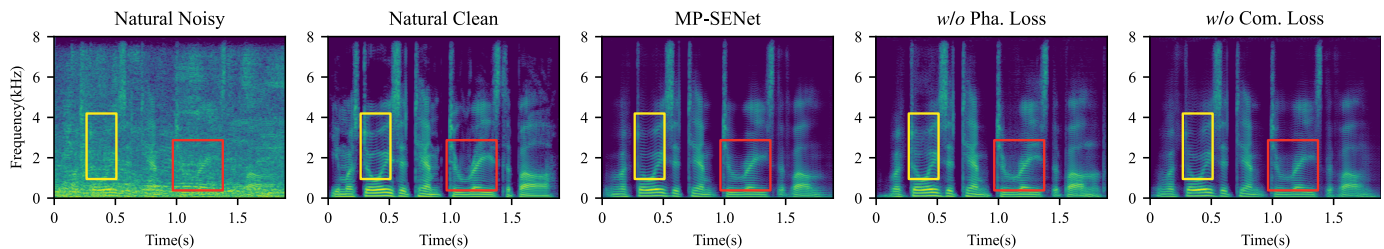


Fig. 5. Spectrogram visualization of the noisy speech, clean speech, and speech waveforms enhanced by our proposed MP-SENet with different combinations of phase optimization approaches on the VoiceBank+DEMAND dataset.

TABLE III
EXPERIMENTAL RESULTS OF THE ABLATION STUDIES ON THE VOICEBANK+DEMAND DATASET.

Method	#Param.	WB-PESQ	PD (°)	SI-SDR (dB)	ViSQOL
MP-SENet	2.26M	3.60	7.28	19.47	4.34
Ablation on Model Structure					
w/ Conformer	1.71M	3.56	7.35	19.38	4.30
Magnitude Only	1.99M	3.44	8.00	18.37	4.28
Complex Only	2.26M	3.58	7.56	18.57	4.28
Ablation on Loss Functions					
w/o Pha. Loss	2.26M	3.52	7.65	19.25	4.32
w/o Com. Loss	2.26M	3.57	7.48	18.60	4.33
w/o Con. Loss	2.26M	3.58	7.26	18.96	4.33
w/o Metric Disc.	2.26M	3.44	7.38	20.05	4.25

Subsequently, we investigated the effect of magnitude-phase parallel modeling. We first ablated the phase decoder to only conduct magnitude enhancement (denoted as “Magnitude Only”), in which the enhanced magnitude was combined with the noisy phase to reconstruct the enhanced speech waveform. As a result, all the metrics significantly degraded, which would be attributed to that without explicit phase prediction, the magnitude would significantly compensate for the noisy phase, resulting in both poor magnitude restoration and phase accuracy. We then re-designed the magnitude mask decoder and the phase decoder to the real-part decoder and imaginary-part decoder, which used complex spectral mapping to directly output the real and imaginary parts of the short-time complex spectrum (denoted as “Complex Only”). Consequently, all the metrics marginally degraded, which demonstrated the superiority of explicit magnitude and phase modeling, and also indicated that complex spectrum-based methods can still achieve exceptional enhancement performance under the guidance of explicit magnitude and phase optimization. Overall, the above experiments demonstrate the effectiveness of the explicit magnitude-phase modeling we proposed.

Furthermore, we conducted ablation studies on the loss functions. To investigate the effects of phase optimization approaches, we conducted ablation studies on the phase spectrum loss (denoted as “w/o Pha. loss”) and complex spectrum loss (denoted as “w/o Com. loss”), which explicitly and implicitly optimized the phase, respectively. Results demonstrated that both of them contributed to the overall performance, but explicit phase optimization played a more pivotal role. We further visualized the spectrograms of the speech waveforms enhanced by these two ablation models and our proposed MP-

SENet as illustrated in Fig 5. It can be clearly observed that after ablating the phase loss, the harmonic structures were significantly distorted, while this impact was slighter when ablating the complex spectrum loss. This further highlighted the importance of explicit phase optimization in alleviating the compensation effect and restoring the harmonic structures. Subsequently, we solely ablated the STFT consistency loss (denoted as *w/o Con. loss*) and the results demonstrated that the STFT consistency did contribute to the model performance. Finally, we removed the metric discriminator (denoted as “w/o Metric disc.”) to assess the effect of the human-perception-related metric loss. Ablation results demonstrated that introducing the metric loss significantly enhanced PESQ and ViSQOL but marginally sacrificed SI-SDR. This suggested that the slight degradation in SI-SDR did not affect the overall enhanced speech quality.

B. Speech Dereverberation

1) *Evaluation metrics*: For the speech dereverberation task, we used the speech quality metrics provided by the REVERB Challenge, including WB-PESQ, cepstral distance (CD), log-likelihood ratio (LLR), frequency-weighted segmental signal-to-noise ratio (FWSegSNR), and signal-to-reverberation modulation energy ratio (SRMR). Since the ‘RealData’ of the REVERB Challenge evaluation set lacked corresponding clean reference speech, only the SRMR metric was employed to evaluate the model’s dereverberation performance in real scenarios. Lower CD and LLR, as well as higher FWSegSNR and SRMR, indicate better performance.

2) *Comparison with baseline methods*: We compared our proposed MP-SENet with a signal processing-based method, i.e., WPE [43], as well as four DNN-based methods, including the magnitude mapping-based UNet [46], and three complex spectrum-based methods (i.e., CMGAN [22], SkipConvGAN [13], and DCN [47]). The objective evaluation results are presented in Table. IV. It’s obvious that DNN-based methods exhibited more robust dereverberation capabilities compared to the signal processing-based WPE. Among the DNN-based methods, UNet which only enhanced the magnitude while retaining the reverberant phase apparently lagged behind other phase-aware methods. Compared to the complex spectrum-based methods, our proposed MP-SENet achieved optimal results in most metrics on the ‘SimData’, with suboptimal results for the remaining LLR and SRMR metrics. However, the MP-SENet achieved the optimal SRMR on the ‘RealData’,

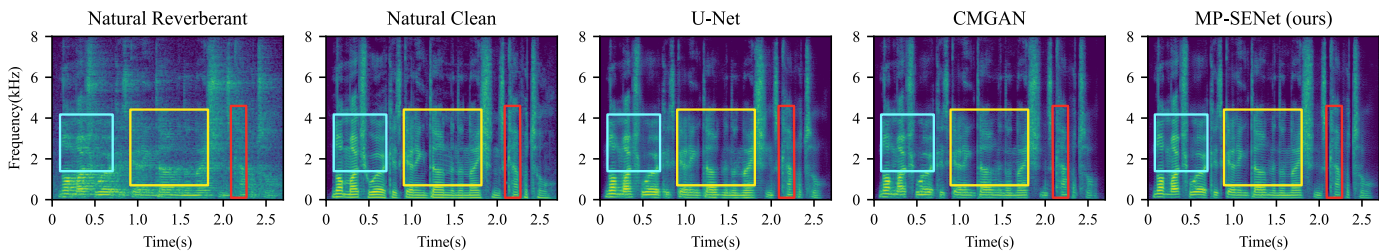


Fig. 6. Spectrogram visualization of the reverberant speech, clean speech, and speech waveforms dereverberated by U-Net, CMGAN, and our proposed MP-SENNet on the REVERB Challenge ‘SimData’ set.

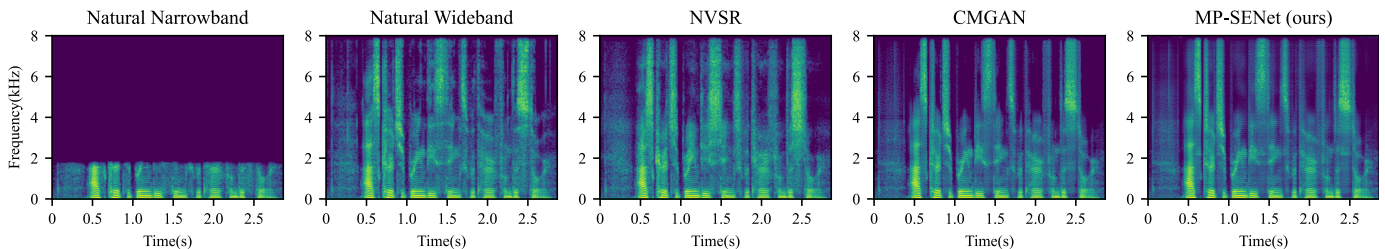


Fig. 7. Spectrogram visualization of the narrowband speech, wideband speech, and speech waveforms extended by NVSR, CMGAN, and our proposed MP-SENNet on the VCTK dataset with the source sampling rate of 4 kHz and target sampling rate of 16 kHz.

TABLE IV

EXPERIMENTAL RESULTS FOR SPEECH DEREVERBERATION METHODS EVALUATED ON THE REVERB CHALLENGE EVALUATION SET.

Method	Year	SimData						RealData
		PESQ	CD	LLR	SNR _{fw}	SRMR	SRMR	
Reverberant	-	1.50	3.97	0.58	3.62	3.69		3.18
WPE [43]	2010	1.72	3.75	0.51	4.90	4.22		3.98
UNet [46]	2018	-	2.50	0.40	10.70	4.88		5.58
SkipConvGAN [13]	2022	2.91	2.32	0.23	11.90	5.89		6.36
CMGAN [22]	2023	-	2.25	0.31	11.74	5.47		6.55
DCN [47]	2024	2.94	2.00	0.23	13.33	5.27		6.48
MP-SENNet	2024	2.97	1.97	0.24	14.07	5.51		6.67

indicating that it exhibited better robustness in real reverberant environments. To intuitively compare the dereverberation performance of different methods, we also visualized the spectrograms of speech waveforms dereverberated by UNet, CMGAN, and our proposed MP-SENNet as shown in Fig. 6. By comparing the contents of boxes with the same color, We can observe that UNet, due to its utilization of the reverberant phase, inevitably suffered from magnitude distortion during iSTFT. Although CMGAN partially restored the phase information through implicit phase optimization, the imprecise phase estimation still induced the compensation effect, resulting in poor restoration of the harmonic structures. Nevertheless, our proposed MP-SENNet, utilizing explicit phase modeling and optimization, effectively restored continuous and intact harmonic structures from the smeared spectra, further improving the quality of the enhanced speech.

C. Speech Bandwidth Extension

1) *Evaluation metrics*: For the speech BWE task, we first adopted the WB-PESQ and ViSQOL to evaluate the magnitude-related speech perceptual quality and the overall speech quality, respectively. Additionally, log-spectral distance

TABLE V

EXPERIMENTAL RESULTS FOR SPEECH BANDWIDTH EXTENSION METHODS ON THE VCTK DATASET. “*” DENOTES THE REPRODUCED RESULTS.

Methods	Year	8 kHz → 16 kHz			4 kHz → 16 kHz		
		WB-PESQ	LSD	ViSQOL	WB-PESQ	LSD	ViSQOL
AudioUNet [7]	2017	3.68	1.32	-	3.39	1.40	-
TFNet [52]	2018	3.72	0.99	-	3.48	1.22	-
AECNN [10]	2021	3.91	0.88	-	3.64	0.95	-
NVSR [51]	2022	3.56	0.79	4.52	2.40	0.95	4.11
CMGAN* [22]	2023	4.24	0.81	4.70	3.73	1.04	4.35
MP-SENNet	2024	4.28	0.66	4.72	3.78	0.81	4.42

(LSD) [82] was utilized to measure the logarithmic distance between two magnitude spectra in dB. lower LSD values indicate better performance.

2) *Comparison with baseline methods*: We compared MP-SENNet with two time-domain methods (i.e., AudioUNet [7] and AECNN [10]), two TF-domain methods (i.e., NVSR [51] and CMGAN [22]), and a hybrid-domain method (i.e., TFNet [52]). The experimental results are presented in Table. V. It is evident that our proposed MP-SENNet achieved the optimal performance with the source sampling rates of both 4 kHz and 8 kHz. Overall, time-domain methods still appeared to be slightly inferior to TF-domain methods, highlighting the importance of low-frequency TF-domain features for recovering high-frequency components. Within the TF-domain methods, we directly visualized the spectrograms of speech waveforms extended by NVSR, CMGAN, and our proposed MP-SENNet with a subsampling rate of 4. Combining the results from Table. V and the spectrograms illustrated in Fig. 7, we can observe that NVSR restored wideband waveforms from the extended mel-spectrogram using a neural vocoder, resulting in higher energy in the high-frequency part compared to real wideband speech, thus leading to the collapse of the PESQ metric. Additionally, since it did not

utilize the low-frequency phase information, its performance in recovering high-frequency harmonics was relatively poorer. Compared to CMGAN, our proposed MP-SENet exhibited a significant improvement in LSD. By comparing the high-frequency components of their spectrograms, we can observe that the MP-SENet demonstrated a strong capability to restore high-frequency harmonics, which confirms the importance of phase information in speech BWE tasks. It is worth noting that MP-SENet adopted the same model structure as the previous speech denoising and dereverberation tasks for the speech BWE task, which achieved an effective extension of narrowband magnitude spectrum through a high-frequency mask, demonstrating the versatility of our proposed parallel magnitude-phase enhancement architecture.

VI. CONCLUSION

In this paper, we proposed a TF-domain monaural SE model, named MP-SENet, which explicitly enhanced both magnitude and phase spectra in parallel. The MP-SENet comprised an encoder-decoder architecture, where the encoder encoded distorted magnitude and phase spectra to TF-domain representations, and the parallel magnitude mask decoder and phase decoder output the enhanced magnitude and phase spectra, respectively. The encoder and decoder are bridged by TF-Transformers to capture time and frequency dependencies. The major contribution of the MP-SENet lay in the explicit modeling and optimization of the phase spectrum, which provided a new solution for the current SE framework, and further elevated the SE performance to a new level. Experimental results demonstrated that our proposed MP-SENet achieved SOTA SE performance across speech denoising, dereverberation, and bandwidth extension tasks with a unified framework. Spectrogram visualizations demonstrated the powerful harmonic restoration capability of our MP-SENet, indicating that explicit phase estimation can effectively alleviate the compensation effect. Further reducing the model parameters and improving the inference efficiency will be the focus of our future work.

REFERENCES

- [1] Y.-X. Lu, Y. Ai, and Z.-H. Ling, "MP-SENet: A speech enhancement model with parallel denoising of magnitude and phase spectra," in *Proc. Interspeech*, 2023, pp. 3834–3838.
- [2] F. Wengler, H. Erdogan, S. Watanabe, E. Vincent, J. Le Roux, J. R. Hershey, and B. Schuller, "Speech enhancement with LSTM recurrent neural networks and its application to noise-robust ASR," in *International Conference on Latent Variable Analysis and Signal Separation (LVA/ICA)*, 2015, pp. 91–99.
- [3] J. L. Desjardins and K. A. Doherty, "The effect of hearing aid noise reduction on listening effort in hearing-impaired adults," *Ear and hearing*, vol. 35, no. 6, pp. 600–610, 2014.
- [4] D. Wang and J. Lim, "The unimportance of phase in speech enhancement," *IEEE/ACM Transactions on Acoustics, Speech, and Signal Processing*, vol. 30, no. 4, pp. 679–681, 1982.
- [5] K. Paliwal, K. Wójcicki, and B. Shannon, "The importance of phase in speech enhancement," *Speech Communication*, vol. 53, no. 4, pp. 465–494, 2011.
- [6] S. Pascual, A. Bonafonte, and J. Serrà, "SEGAN: Speech enhancement generative adversarial network," in *Proc. Interspeech*, 2017, pp. 3642–3646.
- [7] V. Kuleshov, S. Z. Enam, and S. Ermon, "Audio super-resolution using neural nets," in *Proc. ICLR (Workshop Track)*, 2017.
- [8] A. Défossez, G. Synnaeve, and Y. Adi, "Real time speech enhancement in the waveform domain," in *Proc. Interspeech*, 2020, pp. 3291–3295.
- [9] E. Kim and H. Seo, "SE-Conformer: Time-domain speech enhancement using conformer," in *Proc. Interspeech*, 2021, pp. 2736–2740.
- [10] H. Wang and D. Wang, "Towards robust speech super-resolution," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 2058–2066, 2021.
- [11] H.-S. Choi, J.-H. Kim, J. Huh, A. Kim, J.-W. Ha, and K. Lee, "Phase-aware speech enhancement with deep complex U-Net," in *Proc. ICLR*, 2019.
- [12] X. Hao, X. Su, R. Horaud, and X. Li, "FullSubNet: A full-band and sub-band fusion model for real-time single-channel speech enhancement," in *Proc. ICASSP*, 2021, pp. 6633–6637.
- [13] V. Kothapally and J. H. Hansen, "SkipConvGAN: Monaural speech dereverberation using generative adversarial networks via complex time-frequency masking," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 30, pp. 1600–1613, 2022.
- [14] S. Zhao, B. Ma, K. N. Watcharasupat, and W.-S. Gan, "FRCRN: Boosting feature representation using frequency recurrence for monaural speech enhancement," in *Proc. ICASSP*, 2022, pp. 9281–9285.
- [15] F. Dang, H. Chen, and P. Zhang, "DPT-FSNet: Dual-path transformer based full-band and sub-band fusion network for speech enhancement," in *Proc. ICASSP*, 2022, pp. 6857–6861.
- [16] D. Yin, Z. Zhao, C. Tang, Z. Xiong, and C. Luo, "TridentSE: Guiding speech enhancement with 32 global tokens," in *Proc. Interspeech*, 2023, pp. 3839–3843.
- [17] Z.-Q. Wang, P. Wang, and D. Wang, "Complex spectral mapping for single- and multi-channel speech enhancement and robust ASR," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 1778–1787, 2020.
- [18] Z.-Q. Wang, S. Cornell, S. Choi, Y. Lee, B.-Y. Kim, and S. Watanabe, "TF-GridNet: Making time-frequency domain models great again for monaural speaker separation," in *Proc. ICASSP*, 2023, pp. 1–5.
- [19] A. Li, W. Liu, C. Zheng, C. Fan, and X. Li, "Two heads are better than one: A two-stage complex spectral mapping approach for monaural speech enhancement," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 29, pp. 1829–1843, 2021.
- [20] A. Li, S. You, G. Yu, C. Zheng, and X. Li, "Taylor, can you hear me now? a taylor-unfolding framework for monaural speech enhancement," in *Proc. IJCAI*, 2022, pp. 4193–4200.
- [21] G. Yu, A. Li, C. Zheng, Y. Guo, Y. Wang, and H. Wang, "Dual-branch attention-in-attention transformer for single-channel speech enhancement," in *Proc. ICASSP*, 2022, pp. 7847–7851.
- [22] S. Abdulatif, R. Cao, and B. Yang, "CMGAN: Conformer-based metric-gan for monaural speech enhancement," *arXiv preprint arXiv:2209.11112*, 2022.
- [23] Z.-Q. Wang, G. Wichern, and J. Le Roux, "On the compensation between magnitude and phase in speech separation," *IEEE Signal Processing Letters*, vol. 28, pp. 2018–2022, 2021.
- [24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [25] Y. Ai and Z.-H. Ling, "Neural speech phase prediction based on parallel estimation architecture and anti-wrapping losses," in *Proc. ICASSP*, 2023.
- [26] 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, pp. 2031–2041.
- [27] 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," in *Proc. Interspeech*, 2021, pp. 201–205.
- [28] 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, no. 1, pp. 7–19, 2014.
- [29] Y. Ai, J.-X. Zhang, L. Chen, and Z.-H. Ling, "DNN-based spectral enhancement for neural waveform generators with low-bit quantization," in *Proc. ICASSP*, 2019, pp. 7025–7029.
- [30] D. Yin, C. Luo, Z. Xiong, and W. Zeng, "PHASEN: A phase-and-harmonics-aware speech enhancement network," in *Proc. AAAI*, vol. 34, no. 05, 2020, pp. 9458–9465.
- [31] M. Berouti, R. Schwartz, and J. Makhoul, "Enhancement of speech corrupted by acoustic noise," in *Proc. ICASSP*, vol. 4, 1979, pp. 208–211.
- [32] J. Lim and A. Oppenheim, "All-pole modeling of degraded speech," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 26, no. 3, pp. 197–210, 1978.
- [33] Y. Ephraim, "Statistical-model-based speech enhancement systems," *Proceedings of the IEEE*, vol. 80, no. 10, pp. 1526–1555, 1992.

- [34] M. Dendrinos, S. Bakamidis, and G. Carayannis, "Speech enhancement from noise: A regenerative approach," *Speech Communication*, vol. 10, no. 1, pp. 45–57, 1991.
- [35] Y. Ephraim and H. L. Van Trees, "A signal subspace approach for speech enhancement," *IEEE Transactions on Speech and Audio Processing*, vol. 3, no. 4, pp. 251–266, 1995.
- [36] L. Liu, H. Guan, J. Ma, W. Dai, G. Wang, and S. Ding, "A mask free neural network for monaural speech enhancement," in *Proc. Interspeech*, 2023, pp. 2468–2472.
- [37] G. Hu and D. Wang, "Speech segregation based on pitch tracking and amplitude modulation," in *Proc. WASPAA*, 2001, pp. 79–82.
- [38] S. Srinivasan, N. Roman, and D. Wang, "Binary and ratio time-frequency masks for robust speech recognition," *Speech Communication*, vol. 48, no. 11, pp. 1486–1501, 2006.
- [39] A. Narayanan and D. Wang, "Ideal ratio mask estimation using deep neural networks for robust speech recognition," in *Proc. ICASSP*, 2013, pp. 7092–7096.
- [40] Y. Wang, A. Narayanan, and D. Wang, "On training targets for supervised speech separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 12, pp. 1849–1858, 2014.
- [41] H. Erdogan, J. R. Hershey, S. Watanabe, and J. Le Roux, "Phase-sensitive and recognition-boosted speech separation using deep recurrent neural networks," in *Proc. ICASSP*, 2015, pp. 708–712.
- [42] D. S. Williamson, Y. Wang, and D. Wang, "Complex ratio masking for monaural speech separation," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 3, pp. 483–492, 2015.
- [43] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B.-H. Juang, "Speech dereverberation based on variance-normalized delayed linear prediction," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 7, pp. 1717–1731, 2010.
- [44] N. Roman and J. Woodruff, "Intelligibility of reverberant noisy speech with ideal binary masking," *The Journal of the Acoustical Society of America*, vol. 130, no. 4, pp. 2153–2161, 2011.
- [45] X. Li, J. Li, and Y. Yan, "Ideal ratio mask estimation using deep neural networks for monaural speech segregation in noisy reverberant conditions," in *Proc. Interspeech*, 2017, pp. 1203–1207.
- [46] O. Ernst, S. E. Chazan, S. Gannot, and J. Goldberger, "Speech dereverberation using fully convolutional networks," in *Proc. EUSIPCO*, 2018, pp. 390–394.
- [47] V. Kothapally and J. H. Hansen, "Monaural speech dereverberation using deformable convolutional networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.
- [48] Z.-H. Ling, S.-Y. Kang, H. Zen, A. Senior, M. Schuster, X.-J. Qian, H. M. Meng, and L. Deng, "Deep learning for acoustic modeling in parametric speech generation: A systematic review of existing techniques and future trends," *IEEE Signal Processing Magazine*, vol. 32, no. 3, pp. 35–52, 2015.
- [49] J. Abel, M. Strake, and T. Fingscheidt, "A simple cepstral domain DNN approach to artificial speech bandwidth extension," in *Proc. ICASSP*, 2018, pp. 5469–5473.
- [50] K. Li and C.-H. Lee, "A deep neural network approach to speech bandwidth expansion," in *Proc. ICASSP*, 2015, pp. 4395–4399.
- [51] H. Liu, W. Choi, X. Liu, Q. Kong, Q. Tian, and D. Wang, "Neural vocoder is all you need for speech super-resolution," in *Proc. Interspeech*, 2022, pp. 4227–4231.
- [52] T. Y. Lim, R. A. Yeh, Y. Xu, M. N. Do, and M. Hasegawa-Johnson, "Time-frequency networks for audio super-resolution," in *Proc. ICASSP*, 2018, pp. 646–650.
- [53] S. Wisdom, J. R. Hershey, K. Wilson, J. Thorpe, M. Chinen, B. Patton, and R. A. Saurous, "Differentiable consistency constraints for improved deep speech enhancement," in *Proc. ICASSP*, 2019, pp. 900–904.
- [54] A. Pandey and D. Wang, "Densely connected neural network with dilated convolutions for real-time speech enhancement in the time domain," in *Proc. ICASSP*, 2020, pp. 6629–6633.
- [55] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Instance normalization: The missing ingredient for fast stylization," *arXiv preprint arXiv:1607.08022*, 2016.
- [56] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proc. ICCV*, 2015, pp. 1026–1034.
- [57] K. Wang, B. He, and W.-P. Zhu, "TSTNN: Two-stage transformer based neural network for speech enhancement in the time domain," in *Proc. ICASSP*, 2021, pp. 7098–7102.
- [58] 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 *Proc. ICASSP*, 2022, pp. 7417–7421.
- [59] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer normalization," *arXiv preprint arXiv:1607.06450*, 2016.
- [60] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proc. AISTATS*, 2011, pp. 315–323.
- [61] M. Sperber, J. Niehues, G. Neubig, S. Stüker, and A. Waibel, "Self-attentional acoustic models," *Proc. Interspeech*, pp. 3723–3727, 2018.
- [62] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in *Proc. CVPR*, 2016, pp. 1874–1883.
- [63] J. Le Roux, N. Ono, and S. Sagayama, "Explicit consistency constraints for STFT spectrograms and their application to phase reconstruction," in *Proc. SAPA*, 2008, pp. 23–28.
- [64] J. Le Roux, H. Kameoka, N. Ono, and S. Sagayama, "Fast signal reconstruction from magnitude STFT spectrogram based on spectrogram consistency," in *Proc. DAFX*, vol. 10, 2010, pp. 397–403.
- [65] J. Benesty, S. Makino, and J. Chen, *Speech enhancement*. Springer Science & Business Media, 2006.
- [66] A. W. Rix, J. G. Beerends, M. P. Hollier, and A. P. Hekstra, "Perceptual evaluation of speech quality (PESQ)-a new method for speech quality assessment of telephone networks and codecs," in *Proc. ICASSP*, vol. 2, 2001, pp. 749–752.
- [67] C. H. Taal, R. C. Hendriks, R. Heusdens, and J. Jensen, "An algorithm for intelligibility prediction of time-frequency weighted noisy speech," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 7, pp. 2125–2136, 2011.
- [68] C. Valentini-Botinhao, X. Wang, S. Takaki, and J. Yamagishi, "Investigating RNN-based speech enhancement methods for noise-robust text-to-speech," in *Proc. SSW*, 2016, pp. 146–152.
- [69] C. K. Reddy, V. Gopal, R. Cutler, E. Beyrami, R. Cheng, H. Dubey, S. Matushevych, R. Aichner, A. Aazami, S. Braun *et al.*, "The INTER-SPEECH 2020 deep noise suppression challenge: Datasets, subjective testing framework, and challenge results," in *Proc. Interspeech*, 2020, pp. 2492–2496.
- [70] C. Veaux, J. Yamagishi, and S. King, "The voice bank corpus: Design, collection and data analysis of a large regional accent speech database," in *Proc. O-COCOSDA/CASLRE*, 2013, pp. 1–4.
- [71] J. Thiemann, N. Ito, and E. Vincent, "The diverse environments multi-channel acoustic noise database (DEMAND): A database of multichannel environmental noise recordings," in *Proc. ICA*, vol. 19, no. 1, 2013, p. 035081.
- [72] K. Kinoshita, M. Delcroix, S. Gannot, E. A. P. Habets, R. Haeb-Umbach, W. Kellermann, V. Leutnant, R. Maas, T. Nakatani, B. Raj *et al.*, "A summary of the REVERB challenge: state-of-the-art and remaining challenges in reverberant speech processing research," *EURASIP Journal on Advances in Signal Processing*, vol. 2016, pp. 1–19, 2016.
- [73] T. Robinson, J. Franssen, D. Pye, J. Foote, and S. Renals, "Wsjcam0: a british english speech corpus for large vocabulary continuous speech recognition," in *Proc. ICASSP*, vol. 1, 1995, pp. 81–84.
- [74] M. Lincoln, I. McCowan, J. Vepa, and H. K. Maganti, "The multi-channel wall street journal audio visual corpus (MC-WSJ-AV): Specification and initial experiments," in *Proc. ASRU*, 2005, pp. 357–362.
- [75] C. Veaux, J. Yamagishi, K. MacDonald *et al.*, "CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit," *University of Edinburgh. The Centre for Speech Technology Research (CSTR)*, 2017.
- [76] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017.
- [77] Y. Hu and P. C. Loizou, "Evaluation of objective quality measures for speech enhancement," *IEEE Transactions on audio, speech, and language processing*, vol. 16, no. 1, pp. 229–238, 2007.
- [78] J. Le Roux, S. Wisdom, H. Erdogan, and J. R. Hershey, "SDR-half-baked or well done?" in *Proc. ICASSP*, 2019, pp. 626–630.
- [79] A. Gritsenko, T. Salimans, R. van den Berg, J. Snoek, and N. Kalchbrenner, "A spectral energy distance for parallel speech synthesis," *Proc. NeurIPS*, vol. 33, pp. 13 062–13 072, 2020.
- [80] N. Zheng and X.-L. Zhang, "Phase-aware speech enhancement based on deep neural networks," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 1, pp. 63–76, 2018.
- [81] M. Chinen, F. S. Lim, J. Skoglund, N. Gureev, F. O’Gorman, and A. Hines, "ViSQOL v3: An open source production ready objective speech and audio metric," in *Proc. QoMEX*, 2020, pp. 1–6.
- [82] A. Gray and J. Markel, "Distance measures for speech processing," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 24, no. 5, pp. 380–391, 1976.