# VOICEFIXER: TOWARD GENERAL SPEECH RESTORATION WITH NEURAL VOCODER

**Haohe Liu**[1,2*]**, Qiuqiang Kong**[1]**, Qiao Tian**[1]**, Yan Zhao**[1]**,**
**DeLiang Wang**[2]**, Chuanzeng Huang**[1]**, Yuxuan Wang**[1]
[1] Speech, Audio and Music Intelligence (SAMI) group, ByteDance
[2] Department of Computer Science and Engineering, The Ohio State Univeresity

## ABSTRACT

Speech restoration aims to remove distortions in speech signals. Prior methods mainly focus on single-task speech restoration (SSR), such as speech enhancement or speech declipping. However, SSR systems only focus on one task and do not address the general speech restoration problem. Previous SSR systems also have limited performance in speech restoration tasks such as speech super-resolution. To overcome those limitations, we propose a general speech restoration (GSR) task that attempts to remove multiple distortions simultaneously. Furthermore, we propose *VoiceFixer*[1], a generative framework to address the GSR tasks. *VoiceFixer* consists of an *analysis stage* and a *synthesis stage* to mimic the speech analysis and comprehension of the human auditory system. We employ a ResUNet to model the analysis module and a neural vocoder to model the synthesis module. We evaluate *VoiceFixer* with additive noise, room reverberation, low-resolution, and clipping distortions. Our baseline GSR model achieves a 0.499 higher mean opinion score (MOS) than the speech enhancement SSR model. *VoiceFixer* further surpasses the GSR baseline model on MOS score by 0.256. In addition, we observe that *VoiceFixer* generalizes well to severely degraded real speech recordings, indicating its potential in restoring old movies and historical speeches. The source code is available at `https://github.com/haoheliu/voicefixer_main`.

## 1 INTRODUCTIONS

Speech restoration is a process to restore degraded speech signals to high-quality speech signals. Speech restoration is an important research topic due to speech distortions are ubiquitous. For example, speech is usually surrounded by background noise, blurred by room reverbs, or recorded by low-quality devices (Godsill et al., 2002). Those distortions degrade the perceptual quality of speech for human listeners. Speech restoration has a wide range of applications such as online meeting (Defossez et al., 2020), hearing aids (Van den Bogaert et al., 2009), and audio editting (Van Winkle, 2008). Still, speech restoration remains a challenging problem due to the large variety of distortions in the world.

Previous works in speech restoration mainly focus on single task restoration (SSR), which deals with only one kind of distortion at a time. For example, speech enhancement (Loizou, 2007), speech dereveberation (Naylor & Gaubitch, 2010), speech super-resolution (Kuleshov et al., 2017), or speech declipping (Záviška et al., 2020). However, there can be several different distortions that distort speech signals simultaneously in the real world. Previous SSR systems usually oversimplify the speech distortion types (Kashani et al., 2019; Lin et al., 2021; Kuleshov et al., 2017; Birnbaum
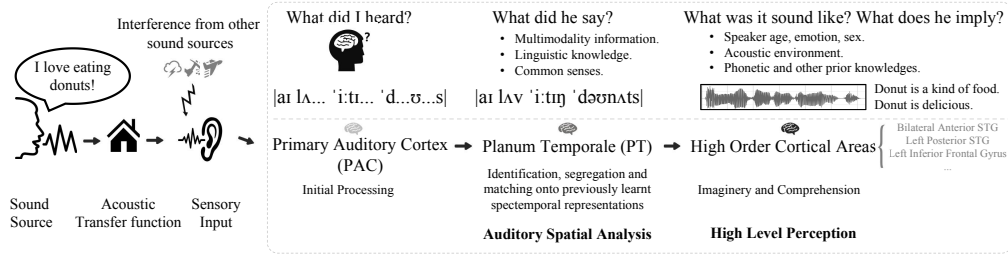
---

*Work done while interning at ByteDance.
[1]Restoration samples can be found in https://haoheliu.github.io/demopage-voicefixer/

et al., 2019). The mismatch between the training data used in SSR and the testing data from the real world degrades the speech restoration performance. Furthermore, previous methods mainly apply one-stage systems to map from restored speech to high-quality speech. However, those one-stage systems do not perform well on generative tasks such as speech super-resolution to generate high-resolution speech from low-resolution speech (Sulun & Davies, 2020; Kuleshov et al., 2017; Lin et al., 2021; Lee & Han, 2021).

To address those problems, we propose a new task called general speech restoration (GSR) aims at restoring multiple distortions in a single model. A numerous studies (Cutler et al., 2021; Cauchi et al., 2014; Han et al., 2015) have reported the benefits of joint training multiple speech restoration tasks. Nevertheless, performing GSR using one-stage systems suffer from the problems in each SSR method. Furthermore, we propose a two-stage system called *VoiceFixer* to address those problems.



**Figure 1:** The neural and cognitive model of how human brain understand and restore distorted speech.

Our proposed *VoiceFixer* is a two-stage speech restoration system that consists of an analysis stage and a synthesis stage. The design of VoiceFixer is motivated by the biological mechanisms of human hearing when a human is restoring distorted speech (Kennedy-Higgins, 2019). Intuitively, if a man tries to identify a strongly distorted voice, his brain can recover the degraded speech by utilizing both the distorted speech signal and his prior knowledge of the speech of the language. VoiceFixer be modeled by a two-stage system, an *auditory scene analysis* stage (Bregman, 1994), and a *high-level comprehension/synthesis* stage (Griffiths & Warren, 2002) as shown in Figure 1. In the analysis stage, the *primary auditory cortex* (*PAC*) processes sound information into acoustic representations. In the analysis stage, the *planum temporale* (*PT*) is hypothesised to be functionally connected with a cortical area that would perform the high-level perception tasks (Griffiths & Warren, 2002; Kennedy-Higgins, 2019). We model the analysis stage with a deep residual UNet and model the synthesis stage with a convolutional vocoder trained with adversarial losses. One advantage of the two-stage *VoiceFixer* is that the analysis and synthesis stages can be trained separately and there are more data to train each stage separately than training the one-stage system. Two-stage methods have also been successfully applied to the speech synthesis task (Wang et al., 2016; Ren et al., 2019; Lin et al., 2021) where acoustic models and vocoders are trained separately.

To our best knowledge, *VoiceFixer* is the speech super-resolution model that successfully restores low-resolution speech of a wide range of sampling rates from 2kHz to 44.1kHz compared to previous works trained and tested on constant sampling rates (Lim et al., 2018; Wang & Wang, 2021; Lee & Han, 2021). *VoiceFixer* is the first model that jointly performs speech enhancement, speech dereverberation, speech super-resolution, and speech declipping in a unified model.

The rest of this paper is organized as follows. Section 2 introduces the formulations of speech distortions. Section 3 introduces the design of *VoiceFixer*. Section 5 discusses the evaluation results. Section 6 concludes this work and discusses future directions. Appendixes introduces related works and shows speech restoration demos.

## 2 PROBLEM FORMULATION

We denote the segment of a speech signal as $s \in \mathbb{R}^L$, where $L$ is the samples number in the segment. We model the distortion process of the speech signal as a function $d(\cdot)$. The degraded speech $x \in \mathbb{R}^L$ can be written as:

$$x = d(s). \tag{1}$$

Speech restoration is a task to restore high-quality speech $\hat{s}$ from $x$:

$$\hat{s} = f(\boldsymbol{x}) \tag{2}$$

where $f(\cdot)$ is the restoration function and can be viewed as a reverse process of $d(\cdot)$. The target is to restore $\hat{s}$ from the observed speech $\boldsymbol{x}$ so that $\hat{s}$ is as close to $\boldsymbol{s}$ as possible. Recently, several deep learning based one-stage methods has been proposed to model $f(\cdot)$ such as fully connected neural networks, recurrent neural networks, and convolutional neural networks introduced in Appendix A.2.

**Distortion simulation** is a an important stage to simulate distorted speech $\boldsymbol{x}$ to build speech restoration systems. Several previous works simulate distortions in a sequential order (Vincent et al., 2017; Cauchi et al., 2014; Tan et al., 2020; Zhao et al., 2019). We denote the distortion $d(\cdot)$ as a composition function:

$$d(\boldsymbol{x}) = d_1 \circ d_2 \circ ...d_Q(\boldsymbol{x}), d_q \in D, q = 1, 2, ..., Q \tag{3}$$

where $\circ$ stands for function composition and $Q$ is the number of distortions to consist $d(\cdot)$. Set $D = \{d_v(\cdot)\}_{v=1}^V$ is the set of different types of distortions where $V$ is the distortion types number. Equation (3) describes the procedure of compounding different distortions from $D$ in a sequential order. We introduce four speech distortion types as follows.

**Additive noise** is one of the most common distortion type and can be modeled by the linear addition between speech $\boldsymbol{s}$ and noise $\boldsymbol{n} \in \mathbb{R}^L$ (Xu et al., 2014; Benesty et al., 2006):

$$d_{\mathrm{noise}}(\boldsymbol{s}) = \boldsymbol{s} + \mathbf{n}. \tag{4}$$

For example, $\boldsymbol{n}$ can be white noise, background sounds, or any non-speech sounds.

**Reverberation** is caused by the reflections of speech in a room. Reverberation leads to speech signals sound distant and blurred. Reverberation can be modeled by convolving speech signals with Room Impulse Response Filters (RIRs) $\boldsymbol{r}$ (Han et al., 2015; Williamson & Wang, 2017):

$$d_{\mathrm{rev}}(\boldsymbol{s}) = \boldsymbol{s} * \boldsymbol{r} \tag{5}$$

where $*$ stands for a convolution operation.

**Low-resolution** distortions refer to audio recordings that are resampled or stored in low sampling rates. There are several causes for low-resolution distortions, such as microphones have low responses in the high frequency, or audio recordings are compressed to low sampling rates. We follow (Wang & Wang, 2021) to perform low-resolution distortion but add more diverse filter type (Sulun & Davies, 2020). For an original sampling rate $o$ and a lower sampling rate $u$, we first convolve $\boldsymbol{s}$ with a low pass filter $\boldsymbol{h}$ to avoid the aliasing phenomenon. Then, the signal is resampled to a sampling rate of $u$:

$$d_{\mathrm{low\_res}}(\boldsymbol{s}) = \mathrm{Resample}(\boldsymbol{s} * \boldsymbol{h}, o, u), \tag{6}$$

**Clipping** distortions refer to the magnitude of audio recordings that are clipped usually caused by low-quality microphones. Clipping can be modeled by restricting signal amplitudes within $[-\eta, +\eta]$:

$$d_{\mathrm{clip}}(\boldsymbol{s}) = \max(\min(\boldsymbol{s}, \eta), -\eta), \eta \in [0, 1]. \tag{7}$$
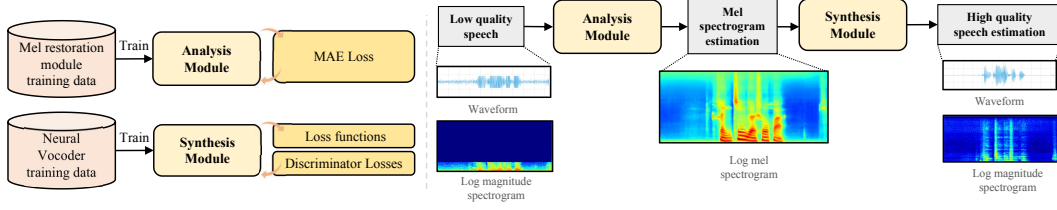
In the frequency domain, clipping effect will produce strong harmonic component in high frequencies and will degrade speech intelligibility.

## 3 METHODOLOGY

### 3.1 ONE-STAGE SPEECH RESTORATION MODELS

Previous deep learning based speech restoration models are usually in one-stage. That is, a models predict restored speech $\hat{s}$ from input $\boldsymbol{x}$ directly:

$$f : \boldsymbol{x} \to \hat{s}. \tag{8}$$

**Figure 2:** Overview of the proposed VoiceFixer model.

The mapping $f(\cdot)$ can be modeled by time domain speech restoration systems such as one-dimensional convolutional neural networks or time-frequency domain systems such as mask-based methods:

$$\hat{\boldsymbol{S}} = (F_{\text{sp}}(|\boldsymbol{X}|\,;\theta) \odot |\boldsymbol{X}|)e^{j\angle\boldsymbol{X}}. \tag{9}$$

where $\boldsymbol{X}$ is the short-time Fourier transform (STFT) of $\boldsymbol{x}$. The STFT $\boldsymbol{X}$ has a shape of $T \times F$ where $T$ is the frames number and $F$ is the frequency bins number. The output of the mask estimation function $F(\cdot;\theta)$ is multiplied with the magnitude of spectrogram $|\boldsymbol{X}|$ to produce the target spectrogram estimation $\hat{\boldsymbol{S}}$. Then, inverse short-time Fourier transform (iSTFT) is applied on $\hat{\boldsymbol{S}}$ to obtain $\hat{s}$. The one-stage speech restoration models are optimized by minimizing the Mean Absolute Error (MAE) between the estimated spectrogram $\hat{\boldsymbol{S}}$ and the target spectrogram $\boldsymbol{S}$ using:

$$L = \left\|\,|\hat{\boldsymbol{S}}| - |\boldsymbol{S}|\,\right\|_1 \tag{10}$$

Previous one-stage models usually build systems on a high-dimensional feature such as time samples and the STFT spectrogram. However, Kuo & Sloan (2005); Trunk (1979) points out that the sparsity in high-dimensional feature will lead to exponential growth in search space. The model can take effect on the high-dimensional feature under the premise of enlarging the model capacity but may also fail in challenging tasks. So, it would be beneficial if we could build a system on a more delicate, low-dimensional feature.

## 4   VOICEFIXER

In this work, we propose *VoiceFixer*, a two-stage speech restoration framework. Multi-stage methods have achieved state-of-the-art performance in the speech synthesis task (Jarrett et al., 2009; Takahama et al., 2019; Zhao et al., 2019; Tan et al., 2020). In speech restoration, our proposed *Voicefixer* breaks the conventional one-stage system into a two-stage system:

$$f : \boldsymbol{x} \mapsto \boldsymbol{z}, \tag{11}$$

$$g : \boldsymbol{z} \mapsto \hat{\boldsymbol{s}}. \tag{12}$$

Equation (11) shows that the first stage of *VoiceFixer* is an analysis stage where a distorted speech $\boldsymbol{x}$ is mapped into a representation $\boldsymbol{z}$. Equation (12) shows that the second stage of *VoiceFixer* is a synthesis stage to synthesis $\boldsymbol{z}$ to restored speech $\hat{\boldsymbol{s}}$. The two-stage *VoiceFixer* mimics the human perception of sounds as described in Section 1 (Kennedy-Higgins, 2019).
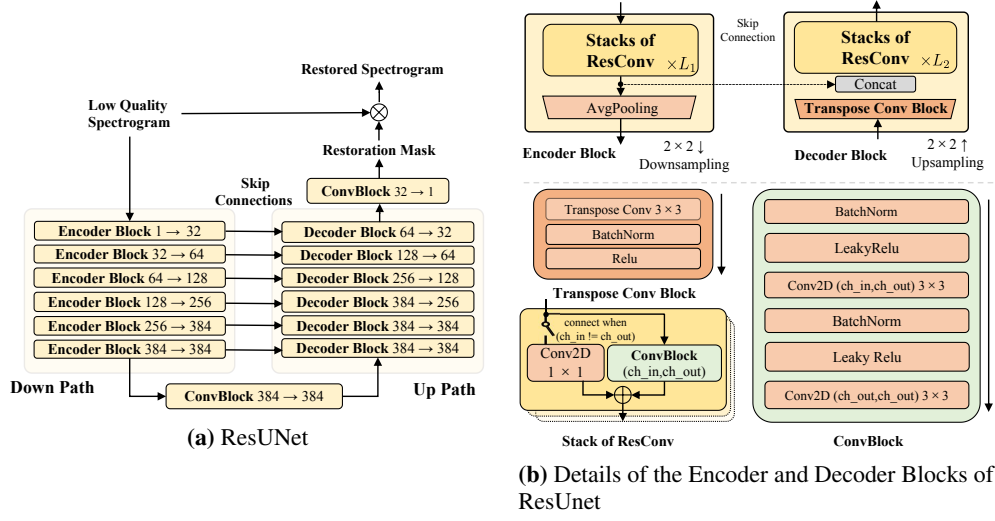
### 4.1   ANALYSIS MODULE

The goal of the analysis module is to predict the intermediate representation $\boldsymbol{z}$ where $\boldsymbol{z}$ can be used to recover the restored speech signals. We choose mel spectrogram as the intermediate representation because mel spectrogram has been widely used as inputs in speech synthesis. The objective of the analysis module becomes to restore mel spectrograms as close as target mel spectrograms as possible. The mel restoration process can be written as equation 13

$$\hat{\boldsymbol{S}}_{\text{mel}} = f_{\text{mel}}(\boldsymbol{X}_{\text{mel}};\alpha) \odot \boldsymbol{X}_{\text{mel}}, \tag{13}$$

where $\boldsymbol{X}_{\text{mel}}$ is the mel spectrogram of $\boldsymbol{x}$ and is calculated by $\boldsymbol{X}_{\text{mel}} = \boldsymbol{X}\boldsymbol{W}$ where $\boldsymbol{W}$ is a set of mel filter banks and has a shape of $F \times F'$. The frequency dimension $F'$ of mel spectrogram is

**(a)** ResUNet

**(b)** Details of the Encoder and Decoder Blocks of ResUnet

**Figure 3:** The architecture of ResUnet, which output have the same size as input.

usually much smaller than the frequency dimension $F$ of STFT. Mel spectrogram is regarded as a way for feature dimension reduction and has been widely used in speech synthesis (Shen et al., 2018; Ren et al., 2019; Kong et al., 2019; Narayanan & Wang, 2013). The mapping $f_{mel}(\cdot; \alpha)$ is the mel restoration mask estimation module parameterized by $\alpha$. The output of $f_{mel}$ is multiplied by $X_{mel}$ to predict the target mel spectrogram.
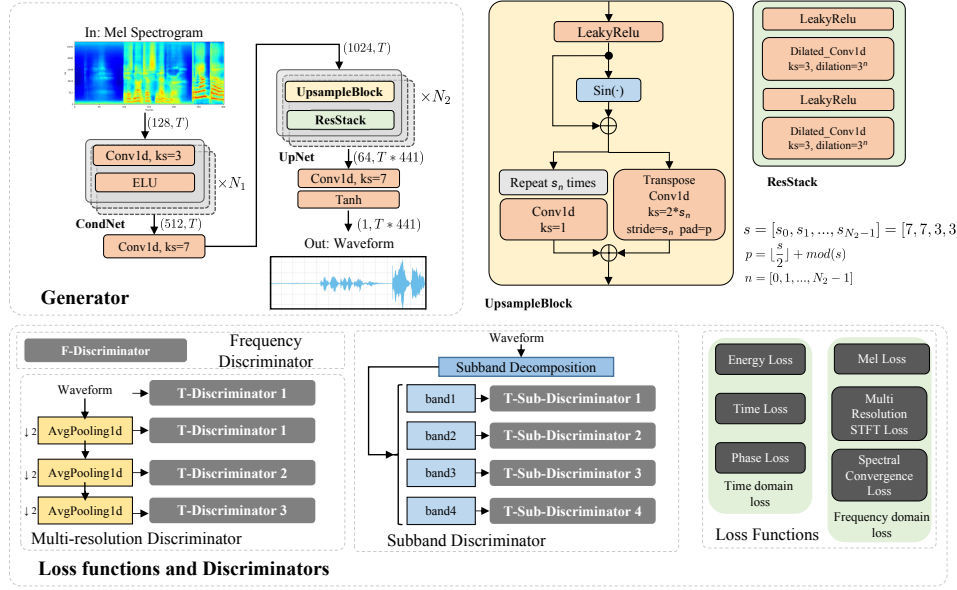
We propose to use ResUNet (Kong et al., 2021a) which is an improvement to UNets (Jansson et al., 2017) to build the analysis module as shown in Figure 3a. The ResUNet consists of several encoder and decoder blocks modeled by convolutional layers. There are skip connections between encoder and decoder blocks with the same level. Figure 3b shows the details of the encoder and decoder block. Both encoder and decoder block share the same *Stack of ResConv* structure, which is a series of residual convolutions (*ResConv*). Each convolutional layer in ResConv consists of a batch normalization (BN) (Ioffe & Szegedy, 2015), a leaky ReLU activation (Xu et al., 2015), and a linear convolutional operation. The encoder block applies average pooling after *Stack of ResConv* for downsampling. The decoder layers apply transpose convolution after *Stack of ResConv* for upsampling. In addition to ResUNet, we implement fully connected Deep Neural Network (DNN) (Ciregan et al., 2012; Szegedy et al., 2013), and Bidirectional GRU (BiGRU) (Chung et al., 2014) as analysis modules for comparison. The *DNN* consists of six fully connected layers. The *BiGRU* have similar structures with *DNN* except for substituting the last two layers in DNN into bi-directional GRU layers.

The details of these three models are discussed in B.1. We will refer to *ResUNet* as *UNet* later for abbreviation. We optimize the analysis module using the $L_1$ norm between the estimated mel spectrogram $\hat{S}_{mel}$ and the target mel spectrogram $S_{mel}$ for training:

$$L_{ana} = \left\| \hat{S}_{mel} - S_{mel} \right\|_1 \tag{14}$$

## 4.2 SYNTHESIS MODULE

There are two advantages of using neural vocoders in the synthesis stage of *VoiceFixer*. First, neural vocoders contain prior knowledge on the structural distribution of speech signals which is important to restore distorted speech (Kennedy-Higgins, 2019). Our vocoder is trained on a large amount of speech data with thousands of speakers and provides sufficient prior knowledge of speech signals. This amount of data is more than conventional SSR methods trained with limited numbers of speakers. Second, building the speech restoration system on the mel spectrogram is regarded as a way of feature dimension reduction compared with processing on the STFT spectrogram because $F' < F$. Dimension reduction can help to remove redundant features, reduce computational costs, simplify the original tasks, and achieve better results (Chao et al., 2019; Chiou & Chen, 2013).

**Figure 4:** The architecture and training scheme of *TFGAN*, which generator is later used as vocoder. The generator takes mel spectrogram as input and upsampled it into waveform. Both the output waveform and its STFT spectrogram are used for the calculation of loss functions. We employ both time and frequency discriminator for discriminative training.

The synthesis module is a neural vocoder that synthesis the mel spectrogram into waveform: equation 15.

$$\hat{s} = g(\hat{X}_{\text{mel}}; \beta), \tag{15}$$

where $g(\cdot; \beta)$ stands for the vocoder model parameterized by $\beta$. We implement the synthesis module using *TFGAN* (Tian et al., 2020), a recently proposed non-autoregressive vocoder.

Figure 4 shows the detailed architecture of *TFGAN*, in which the input mel spectrogram $\hat{X}_{\text{mel}}$ will first pass through a condition network *CondNet*. *CondNet* contains $N_1$ one dimensional convolution (*Conv1d*) layers with exponential linear unit (*ELU*) activations (Clevert et al., 2015). Then, the *UpNet* shown in Figure 4 contains a stack of *UpsampleBlocks* and *ResStacks*. In *UpNet*, the input is upsampled $N_2$ times with upsampling ratios of $s_0, s_1, ...,$ and $s_{N_2-1}$ by each *UpsampleBlock*. The *UpsampleBlock* shows that the input is passed through a *LeakyRelu* activation and is fed into a sinusoidal functions added with the original input to remove periodic artifacts in breathing part of speech (Tian et al., 2020). Then, the output is bifurcated into two branches for upsampling. One branch repeats the samples $s_n$ times followed by a one dimensional convolution. Another branch uses a stride $s_n$ transpose convolution. The output of the repeat and transpose convolution branches are added together as the output of *UpsampleBlock*. *ResStacks* module contains two dilated convolution layers with a *LeakyRelu* non-linearity. The exponentially growing dilation in *ResStack* enable the model to capture long range dependencies. The *TFGAN* in our synthesis model applies $N_2 = 4$. After four *UpsampleBlock* blocks with upsampling ratios $[7, 7, 3, 3]$, each frame of the mel spectrogram is transformed into a sequence with 441 samples corresponding to 10 ms of audio with a sampling rate of 44,100 kHz.

The training criteria of the synthesis module consists of frequency domain losses $L_F$, time domain losses $L_T$, and weighted minimax discriminator losses $L_D$:

$$L_{\text{syn}} = L_F + L_T + \lambda_D L_D, \tag{16}$$

The frequency domain losses $L_F$ is the combination of a mel loss $L_{\text{mel}}$ and multi-resolution STFT losses:

$$L_F(\hat{s}, s) = \lambda_{\text{mel}} L_{\text{mel}}(\hat{s}, s) + \sum_{k=1}^{K_F} (\lambda_{\text{sc}} L_{\text{sc}}^{(k)}(\hat{s}, s) + \lambda_{\text{mag}} L_{\text{mag}}^{(k)}(\hat{s}, s)) \tag{17}$$

where $L_{\text{sc}}$ and $L_{\text{mag}}$ are the spectrogram $L_1$ loss calculated in the linear and log scale respectively. There are $K_F$ different window sizes ranging from 64 to 4096 to calculate $L_{\text{sc}}$ and $L_{\text{mag}}$ so that the trained vocoder is tolerant over phase mismatch Yamamoto et al. (2020); Juvela et al. (2019); Wang et al. (2019). Table 2 in Appendix B.2 shows the detailed configurations.

Time domain losses are complementary to frequency domain losses to address artifacts such as periodic artifacts (Tian et al., 2020). Time domain losses include segment loss $L_{\text{time}}^{(k)}$, energy loss $L_{\text{energy}}^{(k)}$ and phase loss $L_{\text{phase}}^{(k)}$:

$$L_T(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \sum_{k=1}^{K_T} (\lambda_{\text{energy}} L_{\text{energy}}^{(k)}(\hat{\boldsymbol{s}}, \boldsymbol{s}) + \lambda_{\text{phase}} L_{\text{phase}}^{(k)}(\hat{\boldsymbol{s}}, \boldsymbol{s}) + \lambda_{\text{time}} L_{\text{time}}^{(k)}(\hat{\boldsymbol{s}}, \boldsymbol{s})) \tag{18}$$

where $L_{\text{time}}^{(k)}$, energy loss $L_{\text{energy}}^{(k)}$ and phase loss $L_{\text{phase}}^{(k)}$ are described in equation (24 - 26) of Appendix B.2. There are $K_T$ different window sizes ranging from 1 to 960 described to calculate time domain loss with different resolutions. The details of window sizes are shown in Table 3 of Appendix B.2. The energy loss and phase loss have the advantage of alleviating artificial sounds.

Discriminative training is an effective way to optimize the training of neural vocoders (Kong et al., 2020; Kumar et al., 2019). The discriminator we used includes a multi-resolution time discriminator $D_T$, a subband discriminator $D_{T\_sub}$, and frequency discriminator $D_F$:

$$D(\boldsymbol{s}) = D_{T\_sub}(\boldsymbol{s}) + D_F(\boldsymbol{s}) + \sum_{r=1}^{\text{subbands}} D_T^{(r)}(\boldsymbol{s}) \tag{19}$$

$$L_D(\boldsymbol{s}, \hat{\boldsymbol{s}}) = \min_g \max_D (\mathbb{E}_{\boldsymbol{s}}(\log(D(\boldsymbol{s}))) + \mathbb{E}_{\hat{\boldsymbol{s}}}(\log(1 - D(\hat{\boldsymbol{s}})))). \tag{20}$$

The multi-resolution $D_T$ takes signals from four time resolutions after average pooling as input. The subband discriminator $D_{T\_sub}$ performs subband decomposition (Liu et al., 2020) on the waveform, producing four subband signals which are fed into four *T-discrminators* respectively. Frequency discriminator $D_F$ takes the linear spectrogram as input and output real or fake labels. The bottom part of Figure 4 shows the architecture of *T-discriminators* and *F-discriminator*. Appendix B.2 describes the details of the discriminators.

## 5 EXPERIMENTS

### 5.1 DATASETS AND EVALUATION METRICS

**Train sets** We use *VCTK* (Yamagishi et al., 2019) dataset to train the analysis module and , *AISHELL-3* (Shi et al., 2020), and *VCTK* (Yamagishi et al., 2019), *HQ-TTS* (van Niekerk et al., 2017; Sodimana et al., 2018; Guevara-Rukoz et al., 2020) datasets to train the synthesis module. We call the noise datasets used for training *VD-Noise*. To simulate the reverberations, we simulate a set of Room Impulse Response filter to create a *RIR-44k* dataset [2]. The details of those datasets and the simulation of RIR filters are discussed in appendix C.1.

**Test sets**[3] We simulate the additive noise dataset called *DENOISE* using the *VCTK-Demand* (Valentini-Botinhao et al., 2017) dataset. We call our simulated speech super resolution, declipping, and dereverberation evaluationd datasets as *SR*, *DECLI*, and *DEREV*. In addition, we simulate a *ALL-GSR* dataset containing all distortions.

**Evaluation Metrics** The metrics we use include Log-spectral distance (LSD) (Erell & Weintraub, 1990; Lim et al., 2018; Wang & Wang, 2021), Wide Band Perceptual Evaluation of Speech Quality (PESQ-wb) (Rix et al., 2001; Macartney & Weyde, 2018), Structural Similarity (SSIM) (Wang et al., 2004; Sheng et al., 2019), and Scale-Invariant Signal to Noise Ratio (SiSNR) (Záviška et al., 2019a; Zhao et al., 2021). We use Mean Opinion Scores (MOS) to subjevtively evaluate different systems. Details of the metrics are described in Appendix C.4.

---

[2]https://zenodo.org/record/5528124
[3]https://zenodo.org/record/5528144

## 5.2 Distortion Simulation

For the SSR task, we perform only one type of distortion for evaluation. For the GSR task, we first assume that $D = \{d_{\text{noise}}, d_{\text{rev}}, d_{\text{low\_res}}, d_{\text{clip}}\}$ because those distortions are the most common distortions in natural environment (Sulun & Davies, 2020; Ribas et al., 2016). Second, we assume that $Q \leq 4$ in equation 3 which means each distortion in $D$ is performed at most one time. Then, we simulate the distortion process with a specific order $d_{\text{low\_res}}$, $d_{\text{rev}}$, $d_{\text{clip}}$, and $d_{\text{noise}}$. We simulate reverberation, clipping, low-resolution, and noise with random probabilities and random parameters.

## 5.3 Baselines

Table 5 in appendix D concludes the experiments we conduct. We implement several SSR and GSR systems using previous restoration models. For the GSR model, we train a ResUNet model called *USR-UNet* with all distortions. For the SSR model, we implement an *Enh-UNet* with additive noise distortion, a *Derev-UNet* with reverberation distortion, a *SR-UNet* with low-resolution distortion, and a *Declip-UNet* for declipping with declipping distortion. For the SR task, we implement two state-of-the-art models *NuWave* (Lee & Han, 2021) and *SEANet* (Li et al., 2021) for comparison. For Declipping, we implement a state-of-the-art synthesis-based method *SSPADE* (Kitić et al., 2015; Záviška et al., 2019a) using the toolbox[4] provided by Záviška et al. (2020). To explore the impact of parameter number on the performance of mel restoration model, we setup ResUNets with two different sizes. *UNet-S* and *UNet* have one and four *ResConv* blocks in each encoder and decoder block, respectively.
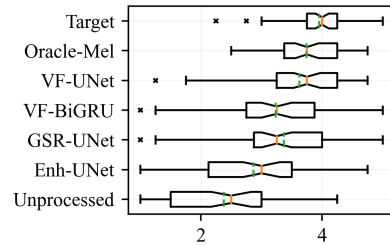
## 5.4 Result Analysis

**Neural Vocoder** To evaluate the performance of the neural vocoder, we provide three baselines. The *Target* system indicates using the perfect $s$ for evaluation. The *Unprocessed* system indicates using distorted speech $x$ for evaluation. The *Oracle-Mel* system indicates using the mel spectrogram of perfect $s$ as input to a Vocoder. The *Oracle-Mel* system shows the performance of the Vocoder. Table 1 shows that the *Oracle-Mel* system achieves a MOS score of 3.74 which is close to the *Target* MOS of 3.95, indicating that the Vocoder performs well in the synthesis task.

| Models | PESQ | LSD | SiSPNR | SSIM | MOS |
|--------|------|-----|--------|------|-----|
| Unprocessed | 1.94 | 2.00 | 7.20 | 0.64 | 2.38 |
| Oracle-Mel | 2.52 | 0.91 | 11.73 | 0.74 | 3.74 |
| Target | 4.64 | 0.01 | 110.55 | 1.00 | 3.95 |
| GSR-UNet | **2.67** | 1.01 | **12.19** | **0.79** | 3.37 |
| Enh-UNet | 2.33 | 1.98 | 9.65 | 0.65 | 2.87 |
| Derev-UNet | 1.97 | 1.81 | 8.50 | 0.59 | / |
| VF-DNN | 1.55 | 1.18 | 10.13 | 0.68 | / |
| VF-BiGRU | 1.92 | 1.02 | 10.98 | 0.71 | 3.24 |
| VF-UNet-S | 2.01 | 1.02 | 11.09 | 0.71 | / |
| VF-UNet | 2.05 | **1.01** | 11.14 | 0.71 | **3.62** |

**Table 1:** Average PESQ, LSD, SiSPNR, SSIM and MOS Score on the general speech restoration testset that contain all kinds of random distortions.
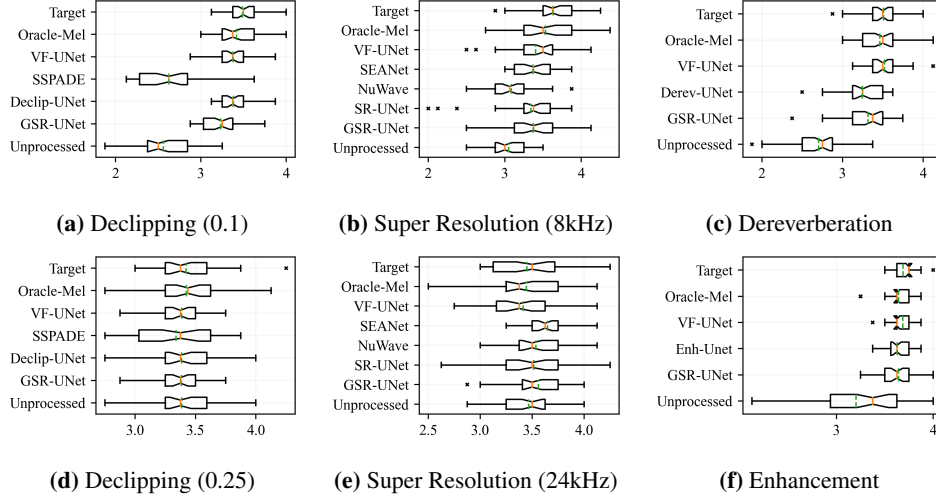
**Figure 5:** Box plot of the MOS scores on general speech restoration task. Red solid line and green dashed line represent median and mean value.



**General Speech Restorations** Table. 1 shows the evaluation results of testset *ALL-GSR*. Figure 5 shows the box plot of the MOS scores of these systems. The *GSR-UNet* outperform the two SSR models, *Enh-UNet* and *Derev-UNet* by a large margin. The *GSR-UNet* model surpasses *Enh-UNet* model by 0.5 on MOS score. This suggests GSR model is more powerful over SSR model on this testset. We abbreviate *VoiceFixer* as *VF*. We observe that the *VF-UNet* model achieves the highest MOS score and LSD score, which is 0.256 higher than the MOS score of *GSR-UNet*. This result indicates that *VoiceFixer* is better than previous models on overall quality. Also, we notice that the MOS score of *VF-UNet* is only 0.11 lower than the *Oracle-Mel*, which indicate the analysis module works well. The MOS score of *Oracle-Mel* is only 0.21 lower than *Target*, meaning *Oracle-Mel* sound very close to the groundtruth. Among the *VoiceFixer* analysis models, the *UNet* front-end achieves the best performance. The *VF-BiGRU* model achieves similar subjective metrics with the

---

[4]https://rajmic.github.io/declipping2020/

*VF-UNet* model but has much lower MOS scores. This phenomenon shows that the improvement in subjective metrics in VoiceFixer is not always consistent with objective evaluation results.



**(a)** Declipping (0.1)  **(b)** Super Resolution (8kHz)  **(c)** Dereverberation

**(d)** Declipping (0.25)  **(e)** Super Resolution (24kHz)  **(f)** Enhancement

**Figure 6:** Box plot of the MOS scores on speech super resolution, declipping, dereveberation and enhancement test sets.

**Speech Super Resolution** Table 6 in appendix D.1 shows the evaluation results on the super resolution testset *SR*. For the 2kHz, 4kHz, and 8kHz to 44.1kHz super-resolution tasks, *VF-UNet* achieved a significantly higher LSD, SiSPNR and SSIM scores than other models. The LSD value of *VF-UNet* in 2k sample rate is still higher than the 8k sample rate score of *GSR-UNet*, *SR-UNet*, *NuWave*, and *SEANet*. This demonstrate the strong performance of *VoiceFixer* on low sample rate cases. The *VF-BiGRU* model outperform *VF-UNet-S* model on average scores for its better performance on low upsample-ratio cases. MOS box plot in Figure 6b shows that *VF-UNet* achieves the highest performance on 8kHz to 44.1kHz super resolution test set. Figure 6e shows 24kHz and 44.1kHz speech have almost the same perceptual quality.

**Speech Enhancement** We evaluate the speech enhancement performance on the *DENOISE* testset and show result in table. 7 in appendix D.1. We found that *GSR-UNet* preserve more details in high-frequency part and have better the PESQ and SiSPNR value than the Enhancement only model *Enh-UNet*. We suppose that's because the data augmentation and joint performing super-resolution can increase the model's generalization and inpainting ability (Hao et al., 2020). The PESQ score of *VF-UNet* reachs to 2.43, higher than *SEGAN*, *WaveUNet*, and model trained with weakly labeled data in Kong et al. (2021b). The MOS evaluation in Figure 6f on speech enhancement task also proves that the result of *VF-UNet* sound on par with previous speech enhancement models.

**Speech Declipping and dereveberation** Table. 9 and table. 8 in appendix D.1 share similar trends on the speech declipping and speech dereveberation performance. In both task, the SSR model *Derev-UNet* and *Declip-UNet* achieved the highest score. The performance of *GSR-UNet* is slightly lower, but it's acceptable considering that *GSR-UNet* is a model that doesn't need extra training for each task. *SSPADE* shows better performance on SiSNR metrics but its PESQ and STOI are lower, especially in the 0.1 threshold case. The MOS score in Figure 6d shows that the clipping effect in the 0.25 threshold case is not apparent enough to perceive, leading to a high MOS score across all methods. While in Figure 6a, the 0.1 threshold clipping test set both *Dec-UNet* and *VF-UNet* achieves the highest objective score. On dereverberation test set *DEREV*, *VF-UNet* achieve the highest MOS score 3.52.

## 6 CONCLUSION

In conclusion, *VoiceFixer* is a effective architecture for speech restoration. It achieves a leading performance across all tasks. The result also shows model trained in a GSR way can achieve a comparable or even better performance than the SSR model.

## REFERENCES

Jacob Benesty, Shoji Makino, and Jingdong Chen. *Speech enhancement*. Springer Science & Business Media, 2006.

Fanhu Bie, Dong Wang, Jun Wang, and Thomas Fang Zheng. Detection and reconstruction of clipped speech for speaker recognition. *Speech Communication*, pp. 218–231, 2015.

Sawyer Birnbaum, Volodymyr Kuleshov, Zayd Enam, Pang Wei Koh, and Stefano Ermon. Temporal film: Capturing long-range sequence dependencies with feature-wise modulations. *arXiv preprint arXiv:1909.06628*, 2019.

Albert S Bregman. *Auditory scene analysis: The perceptual organization of sound*. MIT press, 1994.

Benjamin Cauchi, Ina Kodrasi, Robert Rehr, Stephan Gerlach, Ante Jukic, Timo Gerkmann, Simon Doclo, and Stefan Goetze. Joint dereverberation and noise reduction using beamforming and a single-channel speech enhancement scheme. In *Proc. REVERB Challenge Workshop*, pp. 1–8, 2014.

Guoqing Chao, Yuan Luo, and Weiping Ding. Recent advances in supervised dimension reduction: A survey. *Machine Learning and Knowledge Extraction*, pp. 341–358, 2019.

Bo-Chang Chiou and Chia-Ping Chen. Feature space dimension reduction in speech emotion recognition using support vector machine. In *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference*, pp. 1–6. IEEE, 2013.

Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.

Dan Ciregan, Ueli Meier, and Jürgen Schmidhuber. Multi-column deep neural networks for image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3642–3649, 2012.

Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.

Ross Cutler, Ando Saabas, Tanel Parnamaa, Markus Loide, Sten Sootla, Marju Purin, Hannes Gamper, Sebastian Braun, Karsten Sorensen, Robert Aichner, et al. Acoustic echo cancellation challenge. In *INTERSPEECH*, 2021.

Alexandre Defossez, Gabriel Synnaeve, and Yossi Adi. Real time speech enhancement in the waveform domain. *arXiv preprint arXiv:2006.12847*, 2020.

Per Ekstrand. Bandwidth extension of audio signals by spectral band replication. In *Proceedings of the IEEE Benelux Workshop on Model Based Processing and Coding of Audio*. Citeseer, 2002.

Adoram Erell and Mitch Weintraub. Estimation using log-spectral-distance criterion for noise-robust speech recognition. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 853–856, 1990.

William Fong and Simon Godsill. Monte carlo smoothing for non-linearly distorted signals. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 3997–4000, 2001.

Simon Godsill, Peter Rayner, and Olivier Cappé. Digital audio restoration. In *Applications of digital signal processing to audio and acoustics*, pp. 133–194. Springer, 2002.

Timothy D Griffiths and Jason D Warren. The planum temporale as a computational hub. *Trends in neurosciences*, pp. 348–353, 2002.

Adriana Guevara-Rukoz, Isin Demirsahin, Fei He, Shan-Hui Cathy Chu, Supheakmungkol Sarin, Knot Pipatsrisawat, Alexander Gutkin, Alena Butryna, and Oddur Kjartansson. Crowdsourcing Latin American Spanish for Low-Resource Text-to-Speech. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pp. 6504–6513, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL `https://www.aclweb.org/anthology/2020.lrec-1.801`.

Archit Gupta, Brendan Shillingford, Yannis Assael, and Thomas C Walters. Speech bandwidth extension with wavenet. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pp. 205–208. IEEE, 2019.

Kun Han, Yuxuan Wang, DeLiang Wang, William S Woods, Ivo Merks, and Tao Zhang. Learning spectral mapping for speech dereverberation and denoising. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 982–992, 2015.

Peter SK Hansen. *Signal subspace methods for speech enhancement*. PhD thesis, Citeseer, 1997.

Xiang Hao, Xiangdong Su, Shixue Wen, Zhiyu Wang, Yiqian Pan, Feilong Bao, and Wei Chen. Masking and inpainting: A two-stage speech enhancement approach for low snr and nonstationary noise. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 6959–6963, 2020.

Yanxin Hu, Yun Liu, Shubo Lv, Mengtao Xing, Shimin Zhang, Yihui Fu, Jian Wu, Bihong Zhang, and Lei Xie. DCCRN: Deep complex convolution recurrent network for phase-aware speech enhancement. *arXiv preprint arXiv:2008.00264*, 2020.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pp. 448–456. PMLR, 2015.

Bernd Iser and Gerhard Schmidt. Neural networks versus codebooks in an application for bandwidth extension of speech signals. In *the 8th European Conference on Speech Communication and Technology*, 2003.

Andreas Jansson, Eric Humphrey, Nicola Montecchio, Rachel Bittner, Aparna Kumar, and Tillman Weyde. Singing voice separation with deep U-Net convolutional networks. 2017.

Kevin Jarrett, Koray Kavukcuoglu, Marc'Aurelio Ranzato, and Yann LeCun. What is the best multistage architecture for object recognition? In *IEEE 12th International Conference on Computer Vision*, pp. 2146–2153. IEEE, 2009.

Lauri Juvela, Bajibabu Bollepalli, Junichi Yamagishi, and Paavo Alku. GELP: GAN-Excited linear prediction for speech synthesis from mel-spectrogram. *arXiv preprint arXiv:1904.03976*, 2019.

Thomas Kailath. Lectures on Wiener and Kalman filtering. In *Lectures on Wiener and Kalman Filtering*, pp. 1–143. Springer, 1981.

Hamidreza Baradaran Kashani, Ata Jodeiri, Mohammad Mohsen Goodarzi, and Shabnam Gholamdokht Firooz. Image to image translation based on convolutional neural network approach for speech declipping. *arXiv preprint arXiv:1910.12116*, 2019.

Hideki Kawahara. Straight, exploitation of the other aspect of vocoder: Perceptually isomorphic decomposition of speech sounds. *Acoustical Science and Technology*, pp. 349–353, 2006.

Dan Kennedy-Higgins. *Neural and cognitive mechanisms affecting perceptual adaptation to distorted speech*. PhD thesis, University College London, 2019.

Keisuke Kinoshita, Marc Delcroix, Takuya Yoshioka, Tomohiro Nakatani, Emanuel Habets, Reinhold Haeb-Umbach, Volker Leutnant, Armin Sehr, Walter Kellermann, and Roland Maas. The REVERB challenge: A common evaluation framework for dereverberation and recognition of reverberant speech. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pp. 1–4. IEEE, 2013.

Sran Kitić, Nancy Bertin, and Rémi Gribonval. Sparsity and cosparsity for audio declipping: a flexible non-convex approach. In *International Conference on Latent Variable Analysis and Signal Separation*, pp. 243–250. Springer, 2015.

Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis. *arXiv preprint arXiv:2010.05646*, 2020.

Qiuqiang Kong, Yong Xu, Iwona Sobieraj, Wenwu Wang, and Mark D Plumbley. Sound event detection and time–frequency segmentation from weakly labelled data. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 777–787, 2019.

Qiuqiang Kong, Yin Cao, Haohe Liu, Keunwoo Choi, and Yuxuan Wang. Decoupling magnitude and phase estimation with deep resunet for music source separation. In *The International Society for Music Information Retrieval*, 2021a.

Qiuqiang Kong, Haohe Liu, Xingjian Du, Li Chen, Rui Xia, and Yuxuan Wang. Speech enhancement with weakly labelled data from audioset. *arXiv preprint arXiv:2102.09971*, 2021b.

Juho Kontio, Laura Laaksonen, and Paavo Alku. Neural network-based artificial bandwidth expansion of speech. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 873–881, 2007.

Volodymyr Kuleshov, S Zayd Enam, and Stefano Ermon. Audio super resolution using neural networks. *arXiv preprint arXiv:1708.00853*, 2017.

Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron Courville. Melgan: Generative adversarial networks for conditional waveform synthesis. *arXiv preprint arXiv:1910.06711*, 2019.

Rithesh Kumar, Kundan Kumar, Vicki Anand, Yoshua Bengio, and Aaron Courville. NU-GAN: High resolution neural upsampling with gan. *arXiv preprint arXiv:2010.11362*, 2020.

Frances Y Kuo and Ian H Sloan. Lifting the curse of dimensionality. *Notices of the American Mathematical Society*, pp. 1320–1328, 2005.

Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R Hershey. SDR–half-baked or well done? In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 626–630, 2019.

Katia Lebart, Jean-Marc Boucher, and Philip N Denbigh. A new method based on spectral subtraction for speech dereverberation. *Acta Acustica united with Acustica*, pp. 359–366, 2001.

Junhyeok Lee and Seungu Han. Nu-wave: A diffusion probabilistic model for neural audio upsampling. *arXiv preprint arXiv:2104.02321*, 2021.

Kehuang Li and Chin-Hui Lee. A deep neural network approach to speech bandwidth expansion. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 4395–4399, 2015.

Yunpeng Li, Marco Tagliasacchi, Oleg Rybakov, Victor Ungureanu, and Dominik Roblek. Real-time speech frequency bandwidth extension. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 691–695, 2021.

Teck Yian Lim, Raymond A Yeh, Yijia Xu, Minh N Do, and Mark Hasegawa-Johnson. Time-frequency networks for audio super-resolution. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 646–650, 2018.

Ju Lin, Yun Wang, Kaustubh Kalgaonkar, Gil Keren, Didi Zhang, and Christian Fuegen. A two-stage approach to speech bandwidth extension. *INTERSPEECH*, pp. 1689–1693, 2021.

Haohe Liu, Lei Xie, Jian Wu, and Geng Yang. Channel-wise subband input for better voice and accompaniment separation on high resolution music. *arXiv preprint arXiv:2008.05216*, 2020.

QG Liu, B Champagne, and KC Ho. On the use of a modified fast affine projection algorithm in sub-bands for acoustic echo cancelation. In *IEEE Digital Signal Processing Workshop Proceedings*, pp. 354–357. IEEE, 1996.

Philipos C Loizou. *Speech enhancement: theory and practice*. CRC press, 2007.

Yi Luo and Nima Mesgarani. Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 1256–1266, 2019.

Craig Macartney and Tillman Weyde. Improved speech enhancement with the Wave-U-Net. *arXiv preprint arXiv:1811.11307*, 2018.

Wolfgang Mack and Emanuël AP Habets. Declipping speech using deep filtering. In *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, pp. 200–204. IEEE, 2019.

Rainer Martin. Spectral subtraction based on minimum statistics. *Power*, 1994.

Annamaria Mesaros, Toni Heittola, and Tuomas Virtanen. A multi-device dataset for urban acoustic scene classification. *arXiv preprint arXiv:1807.09840*, 2018.

Masanori Morise, Fumiya Yokomori, and Kenji Ozawa. World: a vocoder-based high-quality speech synthesis system for real-time applications. *IEICE Transactions on Information and Systems*, pp. 1877–1884, 2016.

Anna K Nábělek, Tomasz R Letowski, and Frances M Tucker. Reverberant overlap-and self-masking in consonant identification. *The Journal of the Acoustical Society of America*, pp. 1259–1265, 1989.

Yoshihisa Nakatoh, Mineo Tsushima, and Takeshi Norimatsu. Generation of broadband speech from narrowband speech based on linear mapping. *Electronics and Communications in Japan (Part II: Electronics)*, pp. 44–53, 2002.

Arun Narayanan and DeLiang Wang. Ideal ratio mask estimation using deep neural networks for robust speech recognition. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 7092–7096, 2013.

Patrick A Naylor and Nikolay D Gaubitch. *Speech dereverberation*. Springer Science & Business Media, 2010.

Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.

Santiago Pascual, Antonio Bonafonte, and Joan Serra. SEGAN: Speech enhancement generative adversarial network. *arXiv preprint arXiv:1703.09452*, 2017.

Wei Ping, Kainan Peng, Kexin Zhao, and Zhao Song. Waveflow: A compact flow-based model for raw audio. In *International Conference on Machine Learning*, pp. 7706–7716. PMLR, 2020.

Adam Polyak, Lior Wolf, Yossi Adi, Ori Kabeli, and Yaniv Taigman. High fidelity speech re-generation with application to speech enhancement. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 7143–7147, 2021.

Ryan Prenger, Rafael Valle, and Bryan Catanzaro. Waveglow: A flow-based generative network for speech synthesis. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 3617–3621, 2019.

Joko Radic and Nikola Rozic. Reconstruction of the samples corrupted with impulse noise in multi-carrier systems. In *IEEE Wireless Communications and Networking Conference*, pp. 1–5. IEEE, 2009.

Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. FastSpeech: Fast, robust and controllable text to speech. *arXiv preprint arXiv:1905.09263*, 2019.

Lucas Rencker, Francis Bach, Wenwu Wang, and Mark D Plumbley. Sparse recovery and dictionary learning from nonlinear compressive measurements. *IEEE Transactions on Signal Processing*, pp. 5659–5670, 2019.

Dayana Ribas, Emmanuel Vincent, and José Ramón Calvo. A study of speech distortion conditions in real scenarios for speech processing applications. In *2016 IEEE Spoken Language Technology Workshop (SLT)*, pp. 13–20. IEEE, 2016.

Antony W Rix, John G Beerends, Michael P Hollier, and Andries P Hekstra. Perceptual evaluation of speech quality (pesq)-a new method for speech quality assessment of telephone networks and codecs. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 749–752, 2001.

Tara N Sainath, Oriol Vinyals, Andrew Senior, and Haşim Sak. Convolutional, long short-term memory, fully connected deep neural networks. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 4580–4584, 2015.

Boaz Schwartz, Sharon Gannot, and Emanuël AP Habets. Online speech dereverberation using kalman filter and em algorithm. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 394–406, 2014.

Jonathan Shen, Ruoming Pang, Ron J Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, Rj Skerrv-Ryan, et al. Natural TTS synthesis by conditioning wavenet on Mel spectrogram predictions. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 4779–4783, 2018.

Leyuan Sheng, Dong-Yan Huang, and Evgeniy N Pavlovskiy. High-quality speech synthesis using super-resolution mel-spectrogram. *arXiv preprint arXiv:1912.01167*, 2019.

Yao Shi, Hui Bu, Xin Xu, Shaoji Zhang, and Ming Li. Aishell-3: A multi-speaker mandarin tts corpus and the baselines. *arXiv preprint arXiv:2010.11567*, 2020.

Xiaofeng Shu, Yehang Zhu, Yanjie Chen, Li Chen, Haohe Liu, Chuanzeng Huang, and Yuxuan Wang. Joint echo cancellation and noise suppression based on cascaded magnitude and complex mask estimation. *arXiv preprint arXiv:2107.09298*, 2021.

Brett Y Smolenski and Ravi P Ramachandran. Usable speech processing: A filterless approach in the presence of interference. *IEEE Circuits and Systems Magazine*, pp. 8–22, 2011.

Keshan Sodimana, Knot Pipatsrisawat, Linne Ha, Martin Jansche, Oddur Kjartansson, Pasindu De Silva, and Supheakmungkol Sarin. A Step-by-Step Process for Building TTS Voices Using Open Source Data and Framework for Bangla, Javanese, Khmer, Nepali, Sinhala, and Sundanese. In *Proc. The 6th Intl. Workshop on Spoken Language Technologies for Under-Resourced Languages*, pp. 66–70, Gurugram, India, August 2018. URL http://dx.doi.org/10.21437/SLTU.2018-14.

Serkan Sulun and Matthew EP Davies. On filter generalization for music bandwidth extension using deep neural networks. *IEEE Journal of Selected Topics in Signal Processing*, pp. 132–142, 2020.

Christian Szegedy, Alexander Toshev, and Dumitru Erhan. Deep neural networks for object detection. 2013.

Shusuke Takahama, Yusuke Kurose, Yusuke Mukuta, Hiroyuki Abe, Masashi Fukayama, Akihiko Yoshizawa, Masanobu Kitagawa, and Tatsuya Harada. Multi-stage pathological image classification using semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10702–10711, 2019.

Ke Tan, Yong Xu, Shi-Xiong Zhang, Meng Yu, and Dong Yu. Audio-visual speech separation and dereverberation with a two-stage multimodal network. *IEEE Journal of Selected Topics in Signal Processing*, pp. 542–553, 2020.

Qiao Tian, Yi Chen, Zewang Zhang, Heng Lu, Linghui Chen, Lei Xie, and Shan Liu. TFGAN: Time and frequency domain based generative adversarial network for high-fidelity speech synthesis. *arXiv preprint arXiv:2011.12206*, 2020.

Gerard V Trunk. A problem of dimensionality: A simple example. *IEEE Transactions on pattern analysis and machine intelligence*, (3):306–307, 1979.

Cassia Valentini-Botinhao et al. Noisy speech database for training speech enhancement algorithms and TTS models. 2017.

Jean-Marc Valin and Jan Skoglund. Lpcnet: Improving neural speech synthesis through linear prediction. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 5891–5895, 2019.

Tim Van den Bogaert, Simon Doclo, Jan Wouters, and Marc Moonen. Speech enhancement with multichannel wiener filter techniques in multimicrophone binaural hearing aids. *The Journal of the Acoustical Society of America*, pp. 360–371, 2009.

Daniel van Niekerk, Charl van Heerden, Marelie Davel, Neil Kleynhans, Oddur Kjartansson, Martin Jansche, and Linne Ha. Rapid development of TTS corpora for four South African languages. In *INTERSPEECH*, pp. 2178–2182, Stockholm, Sweden, 2017.

Charles Van Winkle. Audio analysis and spectral restoration workflows using adobe audition. In *Audio Engineering Society Conference: 33rd International Conference: Audio Forensics-Theory and Practice*. Audio Engineering Society, 2008.

Emmanuel Vincent, Shinji Watanabe, Aditya Arie Nugraha, Jon Barker, and Ricard Marxer. An analysis of environment, microphone and data simulation mismatches in robust speech recognition. *Computer Speech & Language*, pp. 535–557, 2017.

Heming Wang and DeLiang Wang. Towards robust speech super-resolution. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2021.

Hong Wang and Fumitada Itakura. Dereverberation of speech signals based on sub-band envelope estimation. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, pp. 3576–3583, 1991.

Wenfu Wang, Shuang Xu, and Bo Xu. First step towards end-to-end parametric tts synthesis: Generating spectral parameters with neural attention. In *INTERSPEECH*, pp. 2243–2247, 2016.

Xin Wang, Shinji Takaki, and Junichi Yamagishi. Neural source-filter-based waveform model for statistical parametric speech synthesis. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 5916–5920, 2019.

Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, pp. 600–612, 2004.

Donald S Williamson and DeLiang Wang. Time-frequency masking in the complex domain for speech dereverberation and denoising. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 1492–1501, 2017.

Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.

Yong Xu, Jun Du, Li-Rong Dai, and Chin-Hui Lee. A regression approach to speech enhancement based on deep neural networks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 7–19, 2014.

Junichi Yamagishi, Christophe Veaux, Kirsten MacDonald, et al. Cstr vctk corpus: English multi-speaker corpus for cstr voice cloning toolkit (version 0.92). 2019.

Ryuichi Yamamoto, Eunwoo Song, and Jae-Min Kim. Parallel WaveGAN: A fast waveform generation model based on generative adversarial networks with multi-resolution spectrogram. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 6199–6203, 2020.

Chengzhu Yu, Heng Lu, Na Hu, Meng Yu, Chao Weng, Kun Xu, Peng Liu, Deyi Tuo, Shiyin Kang, Guangzhi Lei, et al. Durian: Duration informed attention network for multimodal synthesis. *arXiv preprint arXiv:1909.01700*, 2019.

Pavel Záviška, Pavel Rajmic, Ondřej Mokrỳ, and Zdeněk Průša. A proper version of synthesis-based sparse audio declipper. In *Proceedings of the IEEE Conference on Acoustics, Speech, and Signal Processing*, pp. 591–595, 2019a.

Pavel Záviška, Pavel Rajmic, and Jiří Schimmel. Psychoacoustically motivated audio declipping based on weighted l 1 minimization. In *2019 42nd International Conference on Telecommunications and Signal Processing*, pp. 338–342. IEEE, 2019b.

Pavel Záviška, Pavel Rajmic, Alexey Ozerov, and Lucas Rencker. A survey and an extensive evaluation of popular audio declipping methods. *IEEE Journal of Selected Topics in Signal Processing*, pp. 5–24, 2020.

Mengchen Zhao, Xiujuan Yao, Jing Wang, Yi Yan, Xiang Gao, and Yanan Fan. Single-channel blind source separation of spatial aliasing signal based on stacked-lstm. *Sensors*, pp. 4844, 2021.

Yan Zhao, Zhong-Qiu Wang, and DeLiang Wang. Two-stage deep learning for noisy-reverberant speech enhancement. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, pp. 53–62, 2019.

# A    APPENDIX A

## A.1    SPEECH DISTORTIONS

For the distortion types, firstly, the noise signal $\mathbf{n}$ from another sound source can interference with the original speech, degrading its intelligibility. To address this problem, Researchers proposed to conduct Speech Enhancement (Benesty et al., 2006) to remove the undesired noise. Secondly, even in an absolutely quiet environment, speech can interference with itself too. That's because the signal microphone receives is not only the original speech, but also the one reflected by other objects or the echo generated by hardware. The former is related to the problem of Speech Dereveberation (Naylor & Gaubitch, 2010) and the latter can be handled by Acoustic Echo Cancelation (AEC) (Liu et al., 1996). Third, speech distortions caused by hardware deficiencies are also common since not every recording device is professional and ideal. If the recording device has a low response in the high-frequency part, its recording will result in a loss in the higher band, making the speech sound less clear and sparkle. In this case, Band Width Extension (Ekstrand, 2002) can be used for the generation of high-frequency signals. Also, Audio Super Resolution (SR) (Kuleshov et al., 2017) can handle a similar situation when the speech is originally stored at a low sample rate. For example, the Public Switched Telephone Network (PSTN) (Li & Lee, 2015), which provides infrastructure and services for public telecommunication, still uses an 8kHz sample rate with 8-bits resolution to transmit signal from one end to another. In this case, if users need better quality, they can use SR methods to generate the upsampled high-resolution signal. Other methods like click and pop removal (Smolenski & Ramachandran, 2011), corrupted samples replacement (Radic & Rozic, 2009), and declipping (Fong & Godsill, 2001) are also common restoration algorithms for low quality recordings.

## A.2    RELATED WORKS

### A.2.1    SPEECH RESTORATION TASKS

**Audio Super Resolution (SR)** A lot of early study (Nakatoh et al., 2002; Iser & Schmidt, 2003; Kontio et al., 2007) break SR into spectral envelop estimation and excitation generation from low-resolution part. At that time, the direct mapping from the low-resolution part to the high-resolution feature is not widely explored since the dimension of the high-resolution part is relatively high. Later, Deep Neural Network (Li & Lee, 2015; Kuleshov et al., 2017) is introduced to perform SR using spectral mapping on missing high-frequency parts. These approaches show better subjective quality comparing with traditional methods. Later, to increase the modeling capacity, TFilm (Birnbaum et al., 2019) is proposed to modeling affine transformation among each time

block using LSTM. WaveNet also shows effectiveness in extending the bandwidth of a band-limited speech (Gupta et al., 2019). To utilize the information both from the time and frequency domain, (Wang & Wang, 2021) proposed a time-frequency loss that can yield a balanced performance both on time domain and frequency domain metrics. Recently, NU-GAN (Kumar et al., 2020) and NU-Wave (Lee & Han, 2021) pushed the SR's target sample rate to high fidelity, up to 44.1k and 48k.

Although employing Deep Neural Network in BWE task show promising results, the capability of these methods are still restricted and cannot generatlize well on mismatch data. For example, previous approaches (Kuleshov et al., 2017; Gupta et al., 2019) usually train models with a fix setting, i.e., a fix initial sample rate and target sample rate, the same in test data. While in real world applications, speech bandwidth is not usually constant, in which case previous approach usually will fails. Also, since the high-low quality speech pair is impossible to collect. Almost all of BWE model mannuly simulate low quality audio with lowpass filter or bandpass filter during training. In this case, many works based on DNN tend to suffer from overfitting on a specific kind of lowpass filter. As mentioned in (Sulun & Davies, 2020), when the kind of filter used during training and testing differ, the performance can fall considerably. To alleviate filter overfitting, (Sulun & Davies, 2020) proposed to train model with multiple kinds of lowpass filters. By performing data augmentation in this way, unseen filter can be handled properly.

**Speech Declipping (DEC)** The methods for speech declipping can be categorized as supervised methods and unsupervised methods. The unsupervised, or blind methods usually perform declipping based on some generic regularization and assumption of what natural audio should look like, such as ASPADE (Kitić et al., 2015), Dictionary Learning (Rencker et al., 2019), and Psychoacoustically motivated l1 minimization (Záviška et al., 2019b). The supervised models, mostly based on DNN (Bie et al., 2015; Mack & Habets, 2019), are usually trained on clipped and target data pair with backpropagation algorithm. For example, Kashani et al. (2019) treat the declipping as an image to image translation problem and utilize the UNet to do the spectral mapping. By comparison, most of the state-of-the-art methods are unsupervised (Záviška et al., 2020). One of the reasons for their better performance is that they are usually designed to work on all kinds of audio, while the supervised model mainly specialized on the type of their training data, which is usually in limited types However, Záviška et al. (2020) believes supervised model still have the potential for better declipping performance and it will be a promising research topic.

**Speech Enhancement (SE)** Many methods have been proposed in the literature. Classical methods are efficient and effective on stationary noise, such as spectral substraction (Martin, 1994), Wiener and Kalman filtering (Kailath, 1981), and subspace methods (Hansen, 1997) By comparison, deep learning based model such as CLDNN (Sainath et al., 2015), Conv-TasNet (Luo & Mesgarani, 2019) shows higher subjective score and more robustness on complex cases. Recently, new schemes for the training of SE model have emerged. SEGAN (Pascual et al., 2017) tried a generative way to train SE model. DCCRN (Hu et al., 2020) employ the full complex network to perform enhancement. Kong et al. (2021b) achieved a speech enhancement model using only weakly labeled data. And Polyak et al. (2021) realize a enhancement model using a regeneration approach.

**Speech Dereveberation** Many works have been proposed to address this problem. Some of the early methods, such as Inverse filtering (Naylor & Gaubitch, 2010) and sub-band envelope estimation (Wang & Itakura, 1991), aiming at deconvolve the reverberate signal by estimating an inverse filter. But actually, the inverse filter is hard to do the precise estimate and is not robust to the change of RIR filter. Other techniques, spectral substraction (Lebart et al., 2001), is based on an important overlap-masking (Nábělek et al., 1989) effect of reverberation. Schwartz et al. (2014) perform dereverberation using Kalman-Filter and EM algorithm. Recently, deep learning based dereverberation methods have emerged as the state of the art. Han et al. (2015) use full connected Deep Neural Network (DNN) to learn a spectral mapping from reverberate speech to clean speech. In Williamson & Wang (2017), similar to the masking-based SE, authors proposed to do time-frequency mask estimation to perform dereverberation.

### A.2.2 JOINT RESTORATION AND SYNTHETIC RESTORATION

**Joint Restoration** Many works have adopted the joint restoration scheme to increase the model performance. To make the Acoustic Echo Cancellation (AEC) result sound cleaner, MC-TCN (Shu et al., 2021) proposed to jointly perform AEC and Noise suppression at the same time and achieved a mean opinion score of 4.41, outperforming the baseline of INTERSPEECH2021 AEC

Challenge (Cutler et al., 2021) by 0.54. What's more, in the REVERB challenge (Kinoshita et al., 2013), the test set has both reverberation and noise. So the methods (Cauchi et al., 2014) in this challenge need to both perform denoising and dereverberation. Later, in Han et al. (2015), the authors proposed to perform dereverberation and denoising within a single DNN and substantially outperform related methods regarding quality and intelligibility. However, previous joint processing usually involved only two sub-tasks, which are usually denoising and the main task. In our work, we tried to joint performing four or more tasks so that to achieve general restoration. Although the ordering for distortions is complex, the restoration process can all be formulated as multiplying an estimation $\hat{E}$ to the original signal.

**Synthetic Restoration** Directly estimate the source signal from the input mixture is hard sometimes especially when the source signal to noise ratio (SNR) is low. Some works adopted a regeneration approach for the quality enhancement of speech. In Polyak et al. (2021), the authors utilize an ASR model, a pitch extraction model, and a loudness model to extract semantic level information from the speaker. Then they used these features in an encode-decoder network to do the regeneration of speech. To maintain the consistency of speaker characteristics. It used an auxiliary identity network to compute the identity feature. Its result is shown to have better quality and intelligibility. Besides the restoration task, Text-to-Speech (TTS), or Speech Synthesis, is another heated research area. Similar to synthetic speech restoration, which regenerates speech from distorted speech, TTS can be treated as the regeneration of speech from texts.
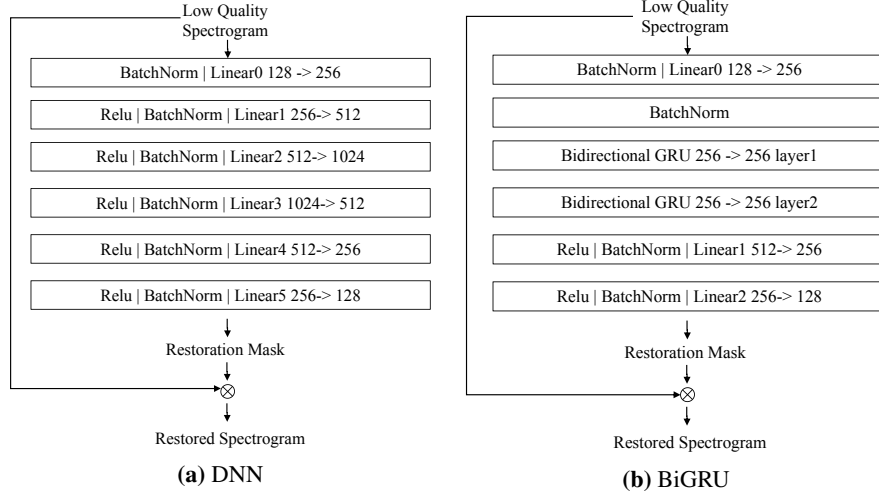
### A.2.3 NEURAL VOCODER

Vocoder, which can directly map the encoded speech feature to the waveform, is an important component in various speech synthesis tasks. The most widely used input feature for vocoder is Mel spectrogram. In recent years, since the emergence of WaveNet (Oord et al., 2016), Neural Network based Vocoder start to demonstrat clear advantages over traditional parametric Vocoder (Morise et al., 2016; Kawahara, 2006). Comparing with traditional methods, the synthesis quality of WaveNet is more closer to human voice. Later, WaveRNN (Yu et al., 2019) is proposed to model the waveform with a single GRU. In this way, WaveRNN has much lower complexity comparing with WaveNet. To get more efficiency, LPCNet (Valin & Skoglund, 2019) combines linear prediction with RNN, which significantly improves inference speed. However, the autoregressive nature of these models and extremely deep structure make their inference speed relatively slow and hard to speed up by parallel. To address this problem, non-autoregressive models like WaveGlow (Prenger et al., 2019) and WaveFlow (Ping et al., 2020) were proposed, which obviously has faster inference speed. Afterward, non-autoregressive GAN based models such as Yamamoto et al. (2020); Kumar et al. (2019) push the synthesis quality to a comparable level with auto-regressive models. Recently, TFGAN (Tian et al., 2020) demonstrated strong capability in vocoding. Directed by frequency discriminator and multiple time-domain loss, it learns waveform information both in the time domain and frequency domain. As a result, the synthesis quality of TFGAN is more natural and less metallic comparing with other GAN-based non-auto-regressive models. In this work, we realize a universal vocoder based on TFGAN, which can reconstruct waveform from Mel spectrogram with perfect perceptual quality. We open-source the pretrained vocoder for the convenience of later research and reproduction of this work.

## B APPENDIX B

### B.1 ANALYSIS MODULE ARCHITECTURE

The DNN and BiGRU we used are shown in Figure 7. the DNN we use in this module is a six layers fully connected network with batch normalization and relu activations. The DNN accept each time step of the low-quality spectrogram as the input feature and output the restoration mask. Similarly, for the BiGRU model, we substitute some layers in DNN to a two-layer bidirectional GRU to capture the time dependency between time steps. To increase the modeling capacity of BiGRU, we expanded the input dimension of GRU to twice of the mel frequency dimension with full connected networks.

The detailed architecture of ResUNet is shown in Figure 3a. In the downpath, the input low quality mel spectrogram will go through 6 encoder blocks, which includes *Stack of ResConv* and a 2 by 2 average pooling, where *Stack of ResConv* is a stack of $L_1$ *ResConv*, in which the outputs of

**(a)** DNN           **(b)** BiGRU

**Figure 7:** The architecture of DNN and Bi-RU we used

*ConvBlock* and the residual convolution are added together as the output. *ConvBlock* is a typical two layers convolution with *BatchNorm* and *LeakyRelu* activation functions. The kernel size of residual convolution and the convolution in *ConvBlock* is $1 \times 1$ and $3 \times 3$. Correspondingly, the decoder blocks have the symmetric structure of the encoder blocks. It first performs a transpose convolution with a 2 by 2 stride and $3 \times 3$ kernels, which result is concatenate with the output of encoder at the same level to form the input of decoder. The *ConvStack* in decoder also contain $L_2$ layers of *res-convblock*. The output of the final decoder block is passed to a final *ConvBlock* to fit the output channel.

## B.2 SYNTHESIS MODULE DETAILS

As shown in table. 3, we use 7 kinds of STFT resolutions and 4 kinds of time resolution during the calculation of $L_F$ and $L_T$. So $K_F = 7$ in equation 17 and $K_T = 4$ in equation 18.

| $k$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| win-length | 4096 | 2048 | 1024 | 512 | 256 | 128 | 64 |
| hop-length | 2048 | 1024 | 512 | 256 | 128 | 64 | 32 |
| fft-size | 8192 | 4096 | 2048 | 1024 | 512 | 256 | 128 |

**Table 2:** STFT setup for different $k$ in $L_F$.

| $k$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| frame-length | 1 | 240 | 480 | 960 |
| hop-length | 1 | 120 | 240 | 480 |

**Table 3:** Windowing setup for different $k$ in $L_T$.

The mel loss $L_{\text{mel}}$, spectral convergence loss $L_{\text{sc}}$, STFT magnitude loss $L_{\text{mag}}$, time loss $L_{\text{time}}$, energy loss $L_{\text{energy}}$, and phase loss $L_{\text{phase}}$ are defined in equation 21, equation 22, equation 23, equation 24, equation 25, and equation 26. The function $v(\cdot)$ is the windowing function that divide time sample into $w$ windows and compute mean value within each window, $v(s)_{1 \times w} = (\text{mean}(\boldsymbol{s}_0), \text{mean}(\boldsymbol{s}_1), ..., \text{mean}(\boldsymbol{s}_w))$. Each $s_w$ stand for windowed $s$. $\Delta$ stand for first difference.

$$L_{\text{mel}}(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \left\| \hat{\boldsymbol{S}}_{\text{mel}} - \boldsymbol{S}_{\text{mel}} \right\|_2 \tag{21}$$

$$L_{\text{sc}}(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \frac{\left\| |\hat{\boldsymbol{S}}| - |\boldsymbol{S}| \right\|_F}{\left\| |\hat{S}| \right\|_F} \tag{22}$$

$$L_{\text{mag}}(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \left\| \log(|\hat{\boldsymbol{S}}|) - \log(|\boldsymbol{S}|) \right\|_1, \tag{23}$$

$$L_{\text{time}}(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \left\| v(\hat{\boldsymbol{s}}) - v(\boldsymbol{s}) \right\|_1, \tag{24}$$

$$L_{\text{energy}}(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \left\| v(\hat{\boldsymbol{s}_w}^2) - v(\boldsymbol{s}_w^2) \right\|_1, \tag{25}$$

$$L_{\text{phase}}(\hat{\boldsymbol{s}}, \boldsymbol{s}) = \left\| \Delta v(\hat{\boldsymbol{s}_w}^2) - \Delta v(\boldsymbol{s}_w^2) \right\|_1,\tag{26}$$

Table 9 and table 8 show the structure of frequency and time domain discriminator. The sub-band discriminator $D_{T\_sub}$ and multi-resolution time discriminator $D_T^{(r)}(s)$ use the structure of *T-discriminator*, which is a stack of one dimensional convolution with grouping and large kernal size. The frequency discriminator $D_F$ use the similar module *ResConv* similar to *ResUNet* shown in Figure 3b.

| T-discriminator |
|---|
| Conv1d(1, 128, kernal_size=16), LeakyRelu(0.2) |
| Conv1d(128, 128, kernal_size=41, stride=4, padding=20, groups=8), LeakyRelu(0.2) |
| Conv1d(128, 128, kernal_size=41, stride=4, padding=20, groups=16), LeakyRelu(0.2) |
| Conv1d(128, 128, kernal_size=41, stride=4, padding=20, groups=32), LeakyRelu(0.2) |
| Conv1d(128, 1, kernal_size=3, stride=1, padding=1), LeakyRelu(0.2) |

**Figure 8:** The structure of *T-discriminator*.

| F-discriminator |
|---|
| Conv2d(1,32,kernal_size=(3,3)) |
| ResConv(32, 32, stride=1,kernal_size=(3,3)) |
| ResConv(32, 32, stride=1,kernal_size=(3,3)) |
| ResConv(32, 64, stride=2,kernal_size=(3,3)) |
| ResConv(64, 64, stride=1,kernal_size=(3,3)) |
| ResConv(64, 32, stride=2,kernal_size=(3,3)) |
| ResConv(32, 32, stride=1,kernal_size=(3,3)) |
| ResConv(32, 32, stride=2,kernal_size=(3,3)) |
| ResConv(32, 32, stride=1,kernal_size=(3,3)) |

**Figure 9:** The structure of *F-discriminator*.

### B.3 TRAINING DETAILS

Except for the training of vocoder, all other models use the same training setup. We use Adam optimizer with $\beta_1 = 0.5, \beta_2 = 0.999$ and a 3e-4 learning rate. We treat the first 1000 steps as the warmup phase, during which the learning rate grows linearly from 0 to 3e-4. We decay the learning rate by 0.9 every 400 hours of training data. We perform an evaluation every 200 hours of training data. If we observe three consecutive evaluations with no improvement, we will interrupt the experiment.

For all the STFT and iSTFT, we use hanning windows with a window length of 2048 and a hop length of 441. As all the audio we use is at 44.1kHz sample rate, the corresponding spectrogram size using this setting will be $T \times 1025$, where T is the dimension of time frames. For mel spectrogram, the dimension of the linear spectrogram is transformed into $T \times 128$ with 128 mel filterbanks.

For the training of synthesis module, we setting up the $\lambda_D$ to $\lambda_{\text{time}}$ value in equation 16, equation 17, and equation 18 as $\lambda_D = 4.0, \lambda_{\text{mel}} = 50, \lambda_{\text{sc}} = 5.0, \lambda_{\text{mag}} = 5.0, \lambda_{\text{energy}} = 100.0, \lambda_{\text{phase}} = 100.0$, and $\lambda_{\text{time}} = 200.0$

## C APPENDIX C

### C.1 DATASETS PREPARATION

**Clean Speech** CSTR VCTK Corpus (Yamagishi et al., 2019) is a multi-speaker English Corpus containing 110 speaker with different accents. We split it into a training part *VCTK-Train* and a testing part *VCTK-Test*. The version of VCTK we used is 0.92. To follow the data preparation strategy of Lee & Han (2021), only the *mic1* microphone data is used for experiments, and *p280* and *p315* are omitted for the technicial issue they have. For the remaining 108 speakers, the last 8 speakers, *p360,p361,p362,p363,p364,p374,p376,s5* are splitted as testset *vctk-test*. Within the other 100 speakers, *p232* and *p257* are omitted because they are used later in the testset *vd-test-noisy*, remaining 98 speakers are defined as *VCTK-Train*. Except for the training of *NuWave*, all the utterance are resampled at 44.1kHz sample rate. *AISHELL-3* is an open source Hi-Fi Mandarin Speech Corpus, containing 88035 utterance with a total duration of 85 hours. The Dataset is originally sampled at 44.1kHz. *hq-tts* dataset contains 191 hours of clean speech data collected from a serial of dataset (van Niekerk et al., 2017; Sodimana et al., 2018; Guevara-Rukoz et al., 2020) collected from *openslr.org*. Here We includes the details of *HQ-TTS* includeing each subset's URL and language type.

**Noise Data** Noise data is used for the simulation of a noisy environment. One of the noise dataset we use come from *VCTK-Demand* (VD) (Valentini-Botinhao et al., 2017), a widely used corpus for

**Table 4:** The components of HQ-TTS dataset.

| URL | Languages | URL | Languages |
|---|---|---|---|
| http://www.openslr.org/32/ | Afrikaans, Sesotho, Setswana and isiXhosa | http://www.openslr.org/70/ | Nigerian English |
| http://www.openslr.org/37/ | Bangladesh Bengali and Indian Bengali | http://www.openslr.org/71/ | Chilean Spanish |
| http://www.openslr.org/41/ | Javanese | http://www.openslr.org/72/ | Colombian Spanish |
| http://www.openslr.org/42/ | Khmer | http://www.openslr.org/73/ | Peruvian Spanish |
| http://www.openslr.org/43/ | Nepali | http://www.openslr.org/74/ | Puerto Rico Spanish |
| http://www.openslr.org/44/ | Sundanese | http://www.openslr.org/75/ | Venezuelan Spanish |
| http://www.openslr.org/61/ | Spanish | http://www.openslr.org/76/ | Basque |
| http://www.openslr.org/63/ | Malayalam | http://www.openslr.org/77/ | Galician |
| http://www.openslr.org/64/ | Marathi | http://www.openslr.org/78/ | Gujarati |
| http://www.openslr.org/65/ | Tamil | http://www.openslr.org/79/ | Kannada |
| http://www.openslr.org/66/ | Telugu | http://www.openslr.org/80/ | Gujarati |
| http://www.openslr.org/69/ | Catalan | | |

Speech Enhancement and noise-robust TTS training and evaluation. This dataset contains a training part *VD-Train* and a testing part *VD-Test*, in which both contain two noisy set *VD-Train-Noisy*, *VD-Test-Noisy* and two clean speech set *VD-Train-Clean*, *VD-Test-Clean*. To obtain the noise data from this dataset, we minus each noisy data from *VD-Train-Noisy* with its corresponding clean part in *VD-Train-Clean* to get the final training noise dataset *VD-Noise*. The noise data was also resampled to 44.1kHz. Another noise dataset we adopt is the TUT Urban Acoustic Scenes 2018 dataset (Mesaros et al., 2018), which is originally used for the acoustic scene classification task of DCASE 2018 Challenge. The dataset contains 89 hours of high-quality recording from 10 acoustic scenes such as airport and shopping mall. The total amount of audio is divided into development *DCASE-Dev* and evaluation *DCASE-Eval* partitions. Both of them contain audio from all cities and all acoustic scenes.

**Room Impulse Response** We randomly simulated a collection of Room Impulse Response filters to simulate the 44.1kHz speech room reverberation using a open source tool [5]. The meters of height, width and length of the room is sampled randomly in a uniform distribution $X \sim U(1, 12)$. The placement of the microphone is then randomly selected within the room space. For the placement of sound source, we first determined the distance between the microphone and sound source, which is randomly sampled in a Gaussian distribution $X \sim N(\mu, \sigma^2), \mu = 2, \sigma = 4$. If the sampled value is negative or greater than five meters, we will sample the distance again until it meets the requirement. After determined the distance between the microphone and sound source, the placement of the sound source is randomly selected on the sphere centered at the microphone. The RT60 value we choose come from the uniform distribution $X \sim U(0.05, 1.0)$. For the pickup pattern of the microphone, we randomly choose from types omnidirectional and cardioid. Finally, we simulated 43239 filters following this scheme, in which we randomly split out 5000 filters as test set *RIR-Test* and named other 38239 filters as *RIR-Train*.

## C.2    TRAINING DATA SIMULATION ALGORITHM

We describe this simulation process in algorithm. 1. $X = \{\boldsymbol{x}^{(0)}, \boldsymbol{x}^{(1)}, ..., \boldsymbol{x}^{(i)}\}$, $N = \{\boldsymbol{n}^{(0)}, \boldsymbol{n}^{(1)}, ..., \boldsymbol{n}^{(i)}\}$, and $R = \{\boldsymbol{r}(0), \boldsymbol{r}(1), ..., \boldsymbol{r}^{(i)}\}$ are the speech dataset, noise dataset, and room Impulse dataset. We use several helper function to describe this algorithm. $\mathrm{randomFilterType}(\cdot)$ is a function that randomly select a type of filter within *Butterworth, Chebyshev, Bessel, Ellipic*. $\mathrm{Resample}(\boldsymbol{x}_1, \boldsymbol{s}_0, \boldsymbol{s}_1)$ is a resampling function that resample the one dimensional signal $x_1$ from a original samplerate $s_0$ to the target $s_1$ samplerate. $\mathrm{buildFilter}(t, c, o)$ is a filter design function that return a type $t$ filter with cutoff frequency $c$ and order $o$. $\max(\cdot), \min(\cdot), and \mathrm{abs}(\cdot)$ is the element wise maximum, minimum, and absolute value function. $\mathrm{mean}(\cdot)$ is the mean value of the input.

We first select a speech utterance $\boldsymbol{x}$, a segment of noise $\boldsymbol{n}$ and a rir filter $\boldsymbol{r}$ randomly from the dataset. Then with $p_1$ probability, we add the reverberate effect using $\boldsymbol{r}$. And with $p_2$ probability, we add clipping effect with a clipping ratio $\eta$, which is sampled in a uniform distrubution $U(\eta_{low}, \eta_{high})$. To simulate the low-resolution data, we first randomly sample the filter type $t$, the cutoff frequency $c$ and order $o$ from {*Butterworth, Chebyshev, Bessel, Ellipic*}, the uniform distribution $U(C_{low}, C_{high})$ and $U(O_{low}, O_{high})$. Secondly we perform convolution using the type $t$ order $o$ lowpass filter at a cutoff frequency $c$. Finally the filtered data will be resampled twice, one is resample to $c * 2$ samplerate and another is resample back to 44.1kHz. We also perform the same lowpass filtering to the noise signal randomly. This operation is necessary because, if not, the model will overfit the

---

[5]https://github.com/sunits/rir_simulator_python

pattern that the bandwidth of noise signal is always different from speech. In this case, the model will fail to remove noise when the bandwidth of noise and speech are similar. For the simulation of noisy environment, we randomly add the noise $n$ into the speech signal $x'$ using a random snr $s \sim U(S_{low}, S_{high})$. To fit the model with all energy level, we randomly conduct a $q \sim U(Q_{low}, Q_{high})$ scaling to the input and target data pair.

In our work, we choose the following parameters to perform this algorithm, $p_1 = 0.25$, $p_2 = 0.25$, $p_3 = 0.5$, $\eta_{low} = 0.06$, $\eta_{high} = 0.9$, $C_{low} = 750$, $C_{high} = 22050$, $O_{low} = 2$, $O_{high} = 10$, $S_{low} = -5$, $S_{high} = 40$, $Q_{low} = 0.3$, $Q_{high} = 1.0$.

---

**Algorithm 1:** Add high quality speech $x$ with random distortions

**In:** $x \leftarrow \mathbb{X}; n \leftarrow \mathbb{N}; r \leftarrow \mathbb{R}$

1    $x' = x$
2    with $p_1$ probability:
3        $x' = x * r$ ;                    `/* Convolute with RIR filter */`
4    with $p_2$ probability:
5        $\theta = S(U(\Theta_{low}, \Theta_{high}))$ ;           `/* Choose clipping ratio */`
6        $x' = \max(\min(x', \theta), -\theta)$ ;        `/* Hard clipping */`
7    with $p_3$ probability:
8        $t = \text{randomFilterType}()$ ;
9        $c \sim U(C_{low}, C_{high}); o \sim U(O_{low}, O_{high})$ ;   `/* Random cutoff and order */`
10       $x' = x' * \text{buildFilter}(t, c, o)$ ;          `/* Low pass filtering */`
11       $x' = \text{Resample}(\text{Resample}(x', 44100, c * 2), c * 2, 44100)$ ;   `/* Resample */`
12       with $p_4$ probability:
13          $n = n * \text{buildFilter}(t, c, o)$ ;  `/* Low pass filtering on noise */`
14          $n = \text{Resample}(\text{Resample}(n, c * 2), 44100)$ ;       `/* Resample */`
15    with $p_5$ probability:
16        $s \sim U(S_{low}, S_{high}); q \sim U(Q_{low}, Q_{high})$ ;   `/* Random SNR and scale */`
17        $n = \frac{n}{\text{mean}(\text{abs}(n))/\text{mean}(\text{abs}(x'))}$ ; `/* Normalize the energy of noise */`
18        $x' = (x' + \frac{n}{10^{s/20}})$ ;             `/* Add noise */`
19    $x' = x' \cdot q$ ;                        `/* Scaling */`
20    $x = x \cdot q$ ;                          `/* Scaling */`
    **Out:** The randomly distorted speech $x'$ and its target $x$

---

### C.3  TEST SET SIMULATIONS

Testing data is crucial for the evaluation on each kind of distortion. The Testing data we use either come from publicly available testset or simulated by ourself.

**Super Resolution** The simulation of *SR* test set follows the work of (Kuleshov et al., 2017; Wang & Wang, 2021; Lim et al., 2018). The low-resolution and target data pair is obtained by transform 44.1kHz sample rate utterances in target speech data *VCTK-Test* to $R_{ls}$ sample rate. To achieve that, we first convolve the speech data with a order 8 Chebyshev type I lowpass filter with the $\frac{R_{ls}}{2}$ cutoff frequency. Then we subsample the signal to $R_{ls}$ sample rate using polyphase filtering. In this work, to test the performance on different samplerate settings, $R_{ls}$ is choosen at 2k, 4k, 8k, 16k and 24k. So there are five testsets for SR task, we denote them as *VCTK-4k*, *VCTK-4k*, *VCTK-8k*, *VCTK-16k*, and *VCTK-24k*.

**Enhancement** For the enhancement task, we adopt the open-sourced *VD-Test-Noisy* described in C.1 as the test set *DENOISE*. This test set contains 824 utterances from a female speaker and a male speaker. The type of noise data comprises a domestic noise (living room), an office noise (office space), noise in transport scene (bus), and two street noises (open area cafeteria and a public square). The test set is simulated at four SNR levels, which are 17.5 dB, 12.5 dB, 7.5 dB, and 2.5 dB. The original data of *VD-Test-Noisy* is sampled at 48kHz. We downsample it to 44.1kHz to fit our experiments.

**Dereverberation** The testset for dereverberation, *DEREV*, is simulated using *VCTK-Test* and *RIR-Test*. For each utterance in *VCTK-Test*, we first randomly select a rir from *RIR-Test*, then we calculate

the convolution between the rir and utterance to build the reverberate speech. Finally we build 2937 reverberate and target data pair in this testset.

**Declipping** *DECLI*, the evaluation set for declipping, is also constructed based on *VCTK-Test*. We perform clipping on *VCTK-Test* following the equation in section 2 and choose 0.25,0.1 as the two setups for the clipping ratio. This result in two declipping testset with different levels, each containing 2937 clipped speech and targets.

**Comprehensive Evaluation** In order to evaluate the model's performance on Universal Speech Restoration, we simulate a testset *ALL-GSR* comprising of speech with all kinds of distortion. The clean speech and noise data used to build *ALL-GSR* is *VCTK-Test* and *DCASE-EVAL*. The simulation procedure of *ALL-GSR* is almost the same to the training data simulation described in 5.2. Totally 501 three seconds long utterances were simulated in this testset.

**MOS Evaluation** We select a small portion from the testsets to carry out MOS evaluation for each one. In *SR*, *DECLI*, and *DEREV*, we select 38 atterances out for human ratings. In *DENOISE* and *ALL-GSR*, we randomly choose 42 and 51 atterances.

## C.4 EVALUATION METRICS

**Log-spectral distance** LSD is a commonly used metrics on the evaluation of super resolution performance (Kumar et al., 2020; Lee & Han, 2021; Wang & Wang, 2021). For target signal $s$ and output estimate $\hat{s}$, LSD can be computed as equation 27, where $S(f,t)$ and $\hat{S}(f,t)$ is the magnitude spectrogram of $s$ and $\hat{s}$.

$$\text{LSD}(s,\hat{s}) = \frac{1}{T}\Sigma_{t=1}^{T}\sqrt{\frac{1}{F}\Sigma_{f=1}^{F}\log_{10}(\frac{S(f,t)^2}{\hat{S}(f,t)^2})^2} \tag{27}$$

**Perceptual Evaluation of Speech Quality** PESQ is widely used in Speech Enhancement literature as their evaluation metrics (Pascual et al., 2017; Hu et al., 2020). It was originally developed to model the subjective test commonly used in telecommunication. PESQ provides a score ranging from -0.5 to 4.5 and the higher the score, the better quality a speech has. In our work, we used an open-sourced implementation of PESQ to compute these metrics. Since PESQ only works on a 16000 sample rate, we performed a 16k downsampling to the output 44.1k audio before evaluation.

**Structural Similarity** SSIM (Wang et al., 2004) is a metrics in image super-resolution. It addresses the shortcoming of pixel-level metrics by taking the image texture into account. We match the implementation of SSIM in (Wang et al., 2004) with ours and compute SSIM as equation 28, where $\mu_S$ and $\sigma_S$ is the mean and standard deviation of $S$. $\text{Cov}_{S\hat{S}}$ is the Covariance of $S$ and $\hat{S}$. $\epsilon_1 = 0.01$ and $\epsilon_2 = 0.02$ are two constant used to avoid zero division. Similarity is measured within the $K$ 7*7 blocks divided from $S$ and $\hat{S}$.

$$\text{SSIM}(s,\hat{s}) = \Sigma_{k=1}^{K}(\frac{(2\mu_{S_k}\mu_{\hat{S}_k} + \epsilon_1)(2\text{Cov}_{S_k\hat{S}_k} + \epsilon_2)}{(\mu_{S_k}^2 + \mu_{\hat{S}_k}^2 + \epsilon_1)(\sigma_{S_k}^2 + \sigma_{\hat{S}_k}^2 + \epsilon_2)}) \tag{28}$$

**Scale-invariant Spectral to Noise Ratio** SiSPNR is a spectral metrics similar to Scale Invariant Signal to Noise Ratio (Si-SNR) (Le Roux et al., 2019). They have the similar idea except Si-SPNR is computed on the magnitude spectrogram. Given the target spectrogram $S$ and estimation $\hat{S}$ the computation of Si-SPNR can be formulated as

$$SiSPNR = 10 * \log_{10}\frac{\|\hat{S}_{target}\|^2}{\|e_{noise}\|^2} \tag{29}$$

where $\hat{S}_{target} = \frac{<\hat{S}S>S}{\|\hat{S}\|^2}$. The scale invariant is guranteed by mean normalization of estimated and target spectrogram. SNR is widely used in literature (Fong & Godsill, 2001) to compare the energy of a signal to its background noise. The higher SiSPNR or SNR indicate less discrepancy between the estimation and target.

**Table 5:** The experiments we performed, and trainsets and testsets we used to during training and evaluation. We use check mark and cross to denote whether a model use the framework of *VoiceFixer* and whether it's trained to perform a single task (SSR) or multiple tasks (GSR). *VF* here stand for *VoiceFixer*.

| Name | VoiceFixer | SSR | GSR | TrainSets | TestSets |
|---|---|---|---|---|---|
| Unprocessed | | | | | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| Oracle-Mel | | | | | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| Vocoder-TFGAN | | | | VCTK-Train; HQ-TTS; AISHELL-3 | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| Enh-UNet | | ✓ | | VCTK-Train; VD-Noise; | DENOISE; ALL-GSR; |
| Derev-UNet | | ✓ | | VCTK-Train; RIR-Train; | DEREV |
| SR-UNet | | ✓ | | VCTK-Train; | SR |
| Declip-UNet | | ✓ | | VCTK-Train; | DECLI |
| NuWave | | ✓ | | VCTK-Train; | SR |
| SEANet | | ✓ | | VCTK-Train; | SR |
| SSPADE | | ✓ | | | DECLI |
| GSR-UNet | | | ✓ | VCTK-Train; VD-Noise; RIR-Train; | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| VF-DNN | ✓ | | ✓ | VCTK-Train; VD-Noise; RIR-Train; | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| VF-BiGRU | ✓ | | ✓ | VCTK-Train; VD-Noise; RIR-Train; | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| VF-UNet-S | ✓ | | ✓ | VCTK-Train; VD-Noise; RIR-Train; | DENOISE; DEREV; SR; DECLI; ALL-GSR; |
| VF-UNet | ✓ | | ✓ | VCTK-Train; VD-Noise; RIR-Train; | DENOISE; DEREV; SR; DECLI; ALL-GSR; |

**Table 6:** Evaluation result on speech super resolution test set *SR*, which contain five kinds of samplerate settings. The metrics is calculated at a target sample rate of 44.1kHz

| TRAINING SCHEME | | REGRESSION BASED MODELS | | | | VOICEFIXER MODELS | | | | OTHERS | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SampleRate Up Ratio | Metrics | GSR-UNet | SR-UNet | NuWave | SEANet | VF-DNN | VF-BiGRU | VF-UNet-S | VF-UNet | Unprocessed | Oracle-Mel | Target |
| 2kHz 22.1 | LSD | 1.34 | 1.19 | 1.41 | 1.33 | 1.18 | 1.08 | 1.08 | **1.05** | 3.13 | 0.89 | / |
| | SiSPNR | 11.03 | 10.89 | 9.19 | 9.78 | 10.67 | 11.84 | 11.65 | **12.10** | 9.18 | 13.65 | / |
| | SSIM | 0.75 | 0.77 | 0.73 | 0.72 | 0.75 | 0.77 | 0.78 | **0.78** | 0.68 | 0.85 | / |
| 4kHz 11.0 | LSD | 1.27 | 1.18 | 1.35 | 1.24 | 1.15 | 1.03 | 1.04 | **1.02** | 2.97 | 0.89 | / |
| | SiSPNR | 11.48 | 11.10 | 9.65 | 10.58 | 11.07 | 12.27 | 11.98 | **12.41** | 9.52 | 13.65 | / |
| | SSIM | 0.77 | 0.78 | 0.76 | 0.72 | 0.75 | 0.79 | 0.79 | **0.79** | 0.71 | 0.85 | / |
| 8kHz 5.5 | LSD | 1.21 | 1.11 | 1.24 | 1.20 | 1.06 | 0.99 | 1.01 | **0.99** | 2.70 | 0.89 | / |
| | SiSPNR | 12.07 | 11.82 | 10.73 | 11.11 | 11.94 | 12.68 | 12.34 | **12.74** | 9.93 | 13.65 | / |
| | SSIM | 0.81 | **0.82** | 0.80 | 0.74 | 0.78 | 0.81 | 0.81 | 0.81 | 0.76 | 0.85 | / |
| | MOS | 3.37 | 3.34 | 3.09 | 3.37 | / | / | / | **3.40** | 3.05 | 3.53 | 3.63 |
| 16kHz 2.8 | LSD | 1.10 | 0.99 | 1.18 | 1.16 | 1.01 | **0.94** | 0.96 | 0.94 | 2.32 | 0.89 | / |
| | SiSPNR | 13.02 | 13.01 | 11.54 | 11.90 | 12.37 | 13.14 | 12.70 | **13.14** | 10.08 | 13.65 | / |
| | SSIM | 0.85 | **0.88** | 0.81 | 0.75 | 0.82 | 0.82 | 0.82 | 0.82 | 0.83 | 0.85 | / |
| 24kHz 1.8 | LSD | 0.97 | **0.91** | 1.12 | 1.15 | 0.93 | 0.91 | 0.94 | 0.92 | 1.91 | 0.89 | / |
| | SiSPNR | **13.96** | 13.81 | 11.63 | 12.58 | 13.21 | 13.38 | 12.86 | 13.38 | 10.40 | 13.65 | / |
| | SSIM | 0.87 | **0.91** | 0.81 | 0.75 | 0.84 | 0.83 | 0.83 | 0.84 | 0.89 | 0.85 | / |
| | MOS | 3.56 | 3.52 | 3.54 | **3.65** | / | / | / | 3.41 | 3.47 | 3.44 | 3.45 |
| Average Score | LSD | 1.18 | 1.07 | 1.26 | 1.21 | 1.07 | 0.99 | 1.01 | **0.98** | 2.61 | 0.89 | / |
| | SiSPNR | 12.31 | 12.13 | 10.55 | 11.19 | 11.85 | 12.66 | 12.31 | **12.75** | 9.82 | 13.65 | / |
| | SSIM | 0.81 | **0.83** | 0.79 | 0.74 | 0.79 | 0.80 | 0.81 | 0.81 | 0.77 | 0.85 | / |

# D   APPENDIX D

## D.1   EVALUATION RESULTS

**Table 7:** Evaluation result on speech enhancement test set *DENOISE*

| Models | SiSNR | PESQ | SiSPNR | MOS |
|---|---|---|---|---|
| Unprocessed | 8.40 | 1.97 | 9.78 | 3.20 |
| Oracle-Mel | -17.52 | 2.85 | 12.84 | 3.64 |
| Target | / | / | / | 3.69 |
| SEGAN * | / | 2.16 | / | / |
| Wave-U-Net * | / | 2.40 | / | / |
| Wiener * | / | 2.22 | / | / |
| Weakly Labelled * | / | 2.28 | / | / |
| GSR-UNet | 16.42 | **2.82** | **12.25** | 3.64 |
| Enh-UNet | **17.58** | 2.71 | 11.82 | 3.63 |
| VF-DNN | / | 1.71 | 10.93 | / |
| VF-BiGRU | / | 2.29 | 11.72 | / |
| VF-UNet-S | / | 2.33 | 11.19 | / |
| VF-UNet | / | 2.43 | 11.71 | 3.69 |

**Table 8:** Evaluation result on speech dere-verberation test set *DEREV*

| Models | PESQ | SiSPNR | MOS |
|---|---|---|---|
| Unprocessed | 1.99 | 14.58 | 2.70 |
| Oracle-Mel | 2.36 | 13.65 | 3.46 |
| Target | / | / | 3.51 |
| GSR-UNet | 2.35 | 14.10 | 3.32 |
| Rev-UNet | **2.49** | **14.99** | 3.25 |
| VF-DNN | 1.41 | 11.70 | / |
| VF-BiGRU | 1.69 | 13.00 | / |
| VF-UNet-S | 1.78 | 12.80 | / |
| VF-UNet | 1.86 | 13.21 | 3.52 |

**Table 9:** Evaluation result on speech declipping test set *DECLI*

| Clipping Level | 0.25 | | | | 0.10 | | | | Average | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Models | SiSNR | STOI | PESQ | MOS | SiSNR | STOI | PESQ | MOS | SiSNR | STOI | PESQ | MOS |
| Unprocessed | 9.60 | 0.95 | 2.38 | 2.56 | 4.00 | 0.89 | 1.51 | 2.72 | 6.80 | 0.92 | 1.95 | 2.64 |
| Oracle-Mel | -19.94 | 0.81 | 2.36 | 3.44 | -19.94 | 0.81 | 2.36 | 3.42 | -19.94 | 0.81 | 2.36 | 3.43 |
| Target | / | / | / | 3.42 | / | / | / | 3.49 | / | / | / | 3.46 |
| GSR-UNet | 11.01 | 0.97 | 3.54 | 3.38 | 7.47 | 0.94 | 2.89 | 3.23 | 9.24 | 0.95 | 3.21 | 3.31 |
| Declip-UNet | 12.45 | **0.99** | **3.98** | 3.38 | 8.43 | **0.96** | **3.40** | 3.38 | 10.44 | **0.98** | 3.69 | 3.38 |
| SSPADE | **17.43** | 0.98 | 3.55 | 3.34 | **10.31** | 0.92 | 2.12 | 2.63 | **13.87** | 0.95 | 2.84 | 2.98 |
| VF-DNN | / | 0.76 | 1.72 | / | / | 0.72 | 1.48 | / | / | 0.74 | 1.60 | / |
| VF-BiGRU | / | 0.81 | 2.09 | / | / | 0.79 | 1.82 | / | / | 0.80 | 1.95 | / |
| VF-UNet-S | / | 0.82 | 2.13 | / | / | 0.80 | 1.85 | / | / | 0.81 | 1.99 | / |
| VF-UNet | / | 0.82 | 2.21 | 3.38 | / | 0.80 | 1.93 | 3.38 | / | 0.81 | 2.07 | 3.38 |

## D.2 METRICS CALCULATE ON ANALYSIS MODULE

In this section, we report the Mel Spectrogram Restoration score on different testset. They are calculated for the evaluation of mel restoration performance using the output of mel restoration module and target spectrogram. We calculate the LSD, SiSPNR, and SSIM value on each setup. The *Unprocessed* column is calculated using the target and unprocessed mel spectrogram. And the *Oracle-Mel* column is calculated between the target spectrogram and itself.

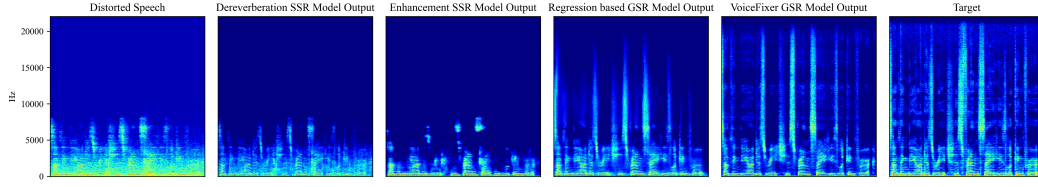**Table 10:** The Performance of Mel Spectrogram Restroation on *DENOISE*, *DEREV*, and *ALL-GSR* testsets

| | DENOISE | | | DEREV | | | ALL-GSR | | |
|---|---|---|---|---|---|---|---|---|---|
| Models | LSD | SiSPNR | SSIM | LSD | SiSPNR | SSIM | LSD | SiSPNR | SSIM |
| Unprocessed | 1.31 | -1.41 | 0.57 | 0.84 | 10.02 | 0.63 | 1.65 | -3.90 | 0.47 |
| VF-DNN | 0.76 | 7.61 | 0.69 | 0.93 | 8.86 | 0.59 | 0.87 | 6.26 | 0.58 |
| VF-BiGRU | 0.55 | 10.98 | 0.79 | 0.56 | 12.91 | 0.75 | 0.59 | 10.49 | 0.70 |
| VF-UNet-S | 0.52 | 10.29 | 0.82 | 0.47 | 13.61 | **0.82** | 0.55 | 11.08 | 0.75 |
| VF-UNet | **0.46** | **12.27** | **0.84** | **0.46** | **14.89** | 0.82 | **0.53** | **11.36** | **0.76** |

Table.10 shows that on *DENOISE*, *DEREV*, and *ALL-GSR*, all four *VoiceFixer* based models are effective on the restortion of mel spectrogram. Among the front-end restoration models, *VoiceFixer-UNet* is consistantly better than the other three models.

**Table 11:** The Performance of Mel Spectrogram Restroation on *SR* testset

| SampleRate Upsampling Ratio | Metrics | MODELS | | | | | |
|---|---|---|---|---|---|---|---|
| | | VF-DNN | VF-BiGRU | VF-UNet-S | VF-UNet | Unprocessed | Oracle-Mel |
| 2kHz 22.1 | LSD | 0.80 | 0.68 | 0.65 | **0.60** | 2.99 | 0.00 |
| | SiSPNR | 8.02 | 9.62 | 9.82 | **11.32** | 2.54 | 127.43 |
| | SSIM | 0.56 | 0.63 | 0.66 | **0.68** | 0.40 | 1.00 |
| 4kHz 11.0 | LSD | 0.68 | 0.54 | 0.55 | **0.50** | 2.54 | 0.00 |
| | SiSPNR | 9.66 | 12.23 | 11.22 | **12.83** | 3.16 | 127.43 |
| | SSIM | 0.65 | 0.72 | 0.74 | **0.76** | 0.51 | 1.00 |
| 8kHz 5.5 | LSD | 0.51 | **0.40** | 0.46 | 0.42 | 2.02 | 0.00 |
| | SiSPNR | 12.53 | **14.85** | 12.67 | 14.20 | 4.26 | 127.43 |
| | SSIM | 0.77 | 0.82 | 0.83 | **0.84** | 0.64 | 1.00 |
| 16kHz 2.8 | LSD | 0.43 | **0.26** | 0.37 | 0.33 | 1.53 | 0.00 |
| | SiSPNR | 13.62 | **19.00** | 14.07 | 16.13 | 5.64 | 127.43 |
| | SSIM | 0.83 | 0.91 | 0.90 | **0.91** | 0.77 | 1.00 |
| 24kHz 1.8 | LSD | 0.29 | **0.18** | 0.31 | 0.27 | 1.16 | 0.00 |
| | SiSPNR | 17.94 | **22.16** | 15.53 | 18.59 | 7.40 | 127.43 |
| | SSIM | 0.92 | 0.95 | 0.94 | **0.95** | 0.86 | 1.00 |
| Average | LSD | 0.54 | **0.41** | 0.47 | 0.43 | 2.05 | 0.00 |
| | SiSPNR | 12.35 | **15.57** | 12.66 | 14.61 | 4.60 | 127.43 |
| | SSIM | 0.75 | 0.80 | 0.81 | **0.83** | 0.64 | 1.00 |

Table.11 lists the mel restoration performance on different samplerates. We found that although *VF-BiGRU* have less parameter than *VF-UNet*, it still achieved the highest score on average LSD and SiSPNR. This result shows the recurrent structure is more suitable for this mel spectrogram super resolution task, especially when the initial samplerate is high, such as 8k, 16k, and 24k.

**Figure 10:** Comparison between different restoration mothods. The unprocessed speech is noisy, reveberate and in low-resolution. The leftmost spectrogram is the unprocessed low quality speech and the rightmost figure is the target high quality spectrogram. In the middle from left to right, the figures show results processed by regression based SSR dereveberation model, SSR enhancement model, GSR model and VoiceFixer based GSR model.
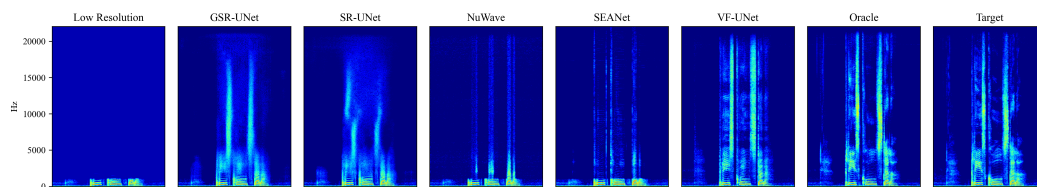
### D.3   DEMOS

In this sections we provide a number of restoration demos using our proposed *VoiceFixer* model. In Figure 12, we provides eight restoration demos using our *VF-UNet* model. All the demos we show are the audio collected from the internet or recorded by ourselves. In each example, the left hand side is the unprocessed spectrogram and the right hand side is the restored one. After restoration, these seriously distorted speech can be revert to a relatively high quality one.

Figure 12b is the speech I recorded myself using Adobe Audition. I set the sample rate of the original recording to 8kHz and manually add the clipping effect after recording. It also contains some low-frequency noise and reverberation introduced by the recording device and environment. Figure 12a is a speech [6] delivered by Amelia Earhart, 1897-1937, appeared in the Library of Congress, United States. It's originally a 5 minutes 43 seconds long audio, in which we only select part of them for this demo. The original version sounds like a mumble because it is in low-resolution with white noise. Figure 12f is an interview in a TV news program. It includes distortions like room reverberation, noise, and low-resolution. Figure 12e is the audio part of a vlog uploaded by a Youtuber. Probably due to the recording device, her speech is deteriorated seriously by noise and the energy of speech in the low-frequency part is also relatively low. Figure 12c is the restoration of a Chinese famous old movie *railroad guerrilla*. Its speech only has limited bandwidth, and part of the frequency information is completely lost. The audio in Figure 12d is selected from a well-known TV series in China, *Romance Three Kindoms*. It's worth noticing that in the original spectrogram, some parts are completely masked off due to the audio compression. Figure 12g is a recording selected from a speech delivered by Sun Yat-sen, 1866-1925. The speech is in extremely low-resolution and includes multiple kinds of unknown distortions. Figure 12h shows the result of a subway broadcasting I record in Shanghai. The low-frequency part of speech is almost lost completely and the reverberation is very serious.

To sum up, all these examples prove the effectiveness of the VoiceFixer model on Universal Restoration. And to our surprise, it can generate will on unseen distortions such as the spectrogram lost in Figure 12c, Figure 12f, and Figure 12d. Also, Figure 12e shows that VoiceFixer is effective for the compensation of low-frequency energy, making speech sound less machinery and distant. Last but not least, despite the abnormal harmonic structure in the low-frequency part in Figure 12g, our proposed model can still repair it into a normal distribution, which proves the advantages of utilizing the prior knowledge of Vocoder.
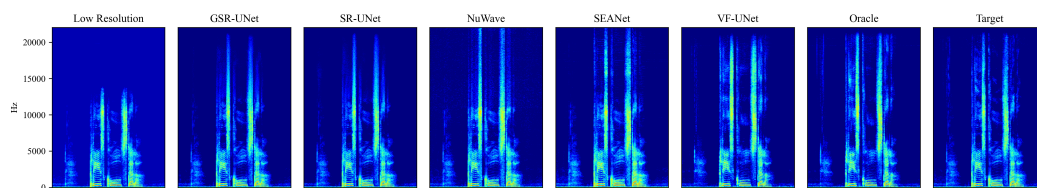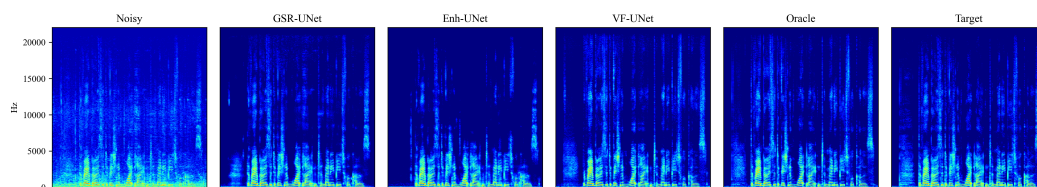
---

[6]https://www.loc.gov/item/afccal000004/

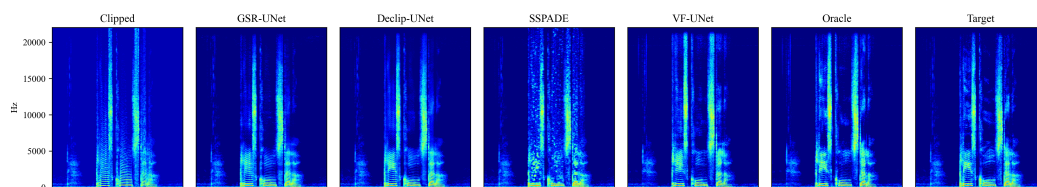**(a)** Speech super resolution results on 2k source samplerate test data.



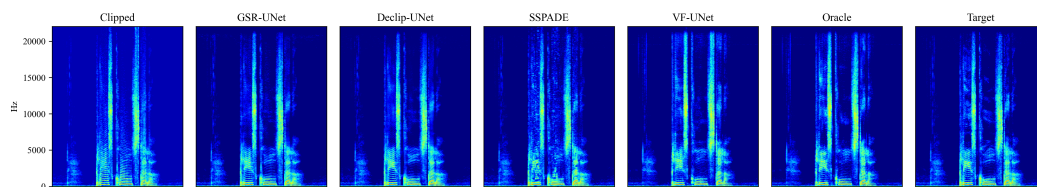**(b)** Speech super resolution results on 8k source samplerate test data.



**(c)** Speech super resolution results on 24k source samplerate test data.
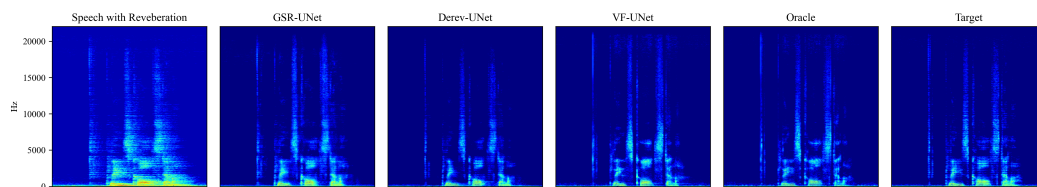


**(d)** Speech Enhancement



**(e)** Speech declipping results on speech with 0.1 clipping threshold
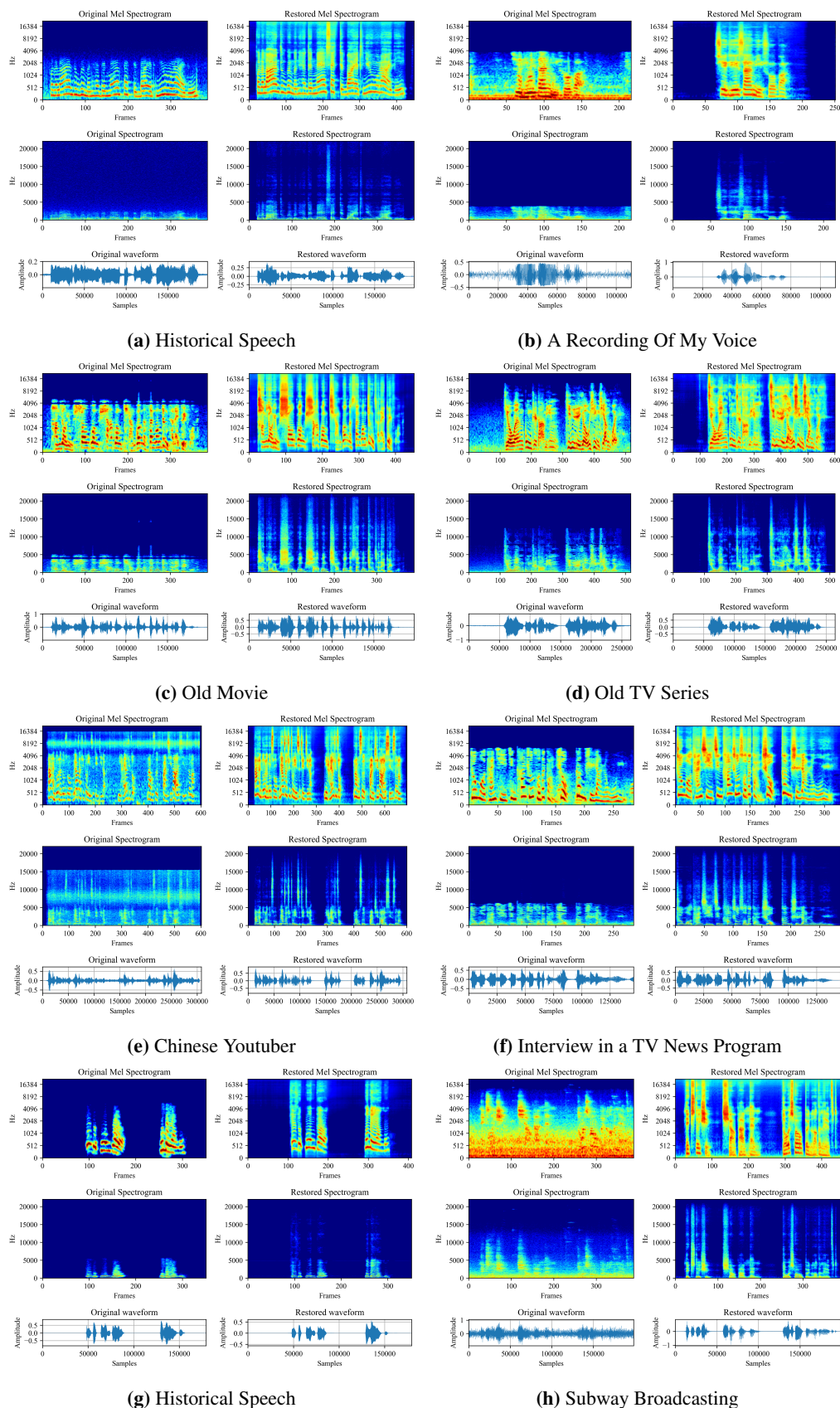


**(f)** Speech declipping results on speech with 0.25 clipping threshold



**(g)** Speech Dereverberation

**Figure 11:** Comparison between different model on four different tasks using simulated data.

**(a)** Historical Speech

**(b)** A Recording Of My Voice

**(c)** Old Movie

**(d)** Old TV Series

**(e)** Chinese Youtuber

**(f)** Interview in a TV News Program

**(g)** Historical Speech

**(h)** Subway Broadcasting

**Figure 12:** Restoration on the data either collected from the internet or recorded by ourselves.