

# FEDSR: FREQUENCY-AWARE ENHANCEMENT FOR DIFFUSION-BASED IMAGE SUPER-RESOLUTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Image super-resolution (ISR) is a classic and challenging problem in low-level vision because the data collection process often introduces complex and unknown degradation patterns. Leveraging powerful generative priors, diffusion-based algorithms have recently established new state-of-the-art ISR performance. Despite the promise, current diffusion-based ISR methods mostly focus on the spatial domain. To bridge this gap, we first experimentally validate that the key to solving the ISR problem lies in addressing the degradation of image amplitude information and high-frequency details. Based on this, we propose a novel *training-free* frequency-aware enhancement framework (**FedSR**) for diffusion-based ISR methods, which consists of two critical components. Firstly, we design the Amplitude Enhancement Module (AEM), which selectively enhances crucial amplitude channels through weighted optimization. Secondly, we introduce the High-Frequency Enhancement Module (HEM) that adaptively masks the skip features to perform high-pass filtering. Through extensive evaluations on both synthetic datasets and real-world image collections, our method demonstrates outstanding performance in reproducing realistic image details without additional tuning. For instance, FedSR improves StableSR across three datasets by **+10.53%** on MUSIQ metric.

## 1 INTRODUCTION

Image super-resolution (ISR) is a fundamental task in low-level vision that aims to reconstruct high-resolution (HR) images from their low-resolution (LR) counterparts. It has widespread applications in areas such as medical imaging (Li et al., 2024; Mao et al., 2023b), satellite imagery (Shermeyer & Etten, 2019; Cornebise et al., 2022), and surveillance systems (Liu et al., 2017; Liang, 2021), where obtaining high-quality images can naturally be subject to hardware limitations and transmission losses. Early ISR algorithms (Dong et al., 2016a; Tai et al., 2017; Chen et al., 2021) attempt to construct synthetic image pairs through simple handcrafted degradation operations (e.g., bicubic downsampling). However, they fail to generalize well in realistic scenarios since real-world LR images typically involve more complex and unknown degradation patterns.

To address this problem, some work (Zhang et al., 2021; Wang et al., 2021) resorts to Generative Adversarial Networks (GAN) (Goodfellow et al., 2014) to enhance visual perception generated by using the adversarial training loss. However, these methods tend to introduce unpleasant visual artifacts because of the instability of adversarial training. Recently, a series of studies (Wang et al., 2023c; Lin et al., 2023; Yu et al., 2024; Wu et al., 2023; Yang et al., 2023) have discovered that incorporating diffusion priors (Rombach et al., 2022) can result in realistic restoration results, achieving state-of-the-art (SOTA) ISR performance. For example, StableSR (Wang et al., 2023c) trains a time-aware encoder to guide Stable Diffusion (Rombach et al., 2022) to achieve promising restoration results; DiffBIR (Lin et al., 2023) employs an IRControlNet trained based on condition images to generate realistic details. Despite the promising results, current diffusion-based ISR methods operate solely in the spatial domain and lack a deep understanding of the frequency domain.

To explore the opportunity to improve diffusion-based ISR models from a frequency perspective, we refer to the following well-established observations: **(1) Loss of High-Frequency Details:** Image degradation often leads to the loss of high-frequency details. **(2) Degradation of Amplitude:** Inspired by tasks such as dehazing (Yu et al., 2022) and deraining (Guo et al., 2022), image degradation can also result in the loss of amplitude information. To systematically validate these phenomena

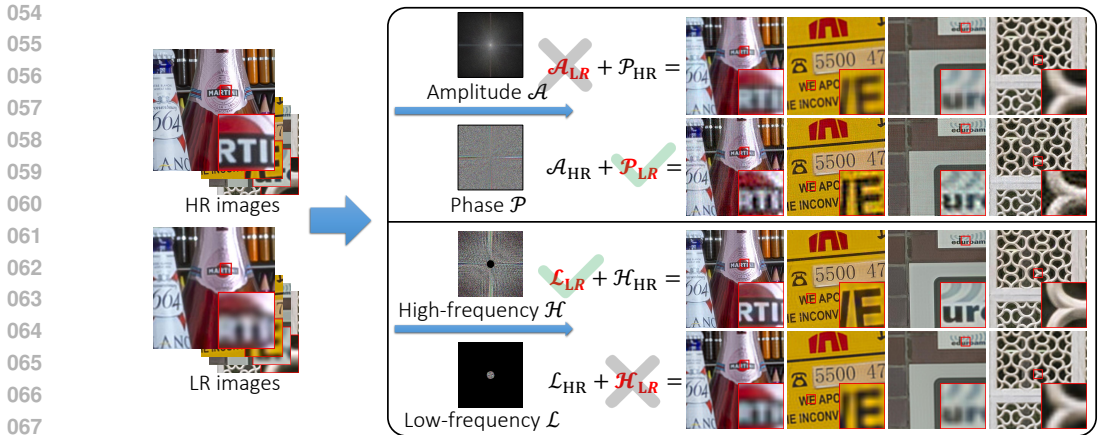


Figure 1: The impact of the real-world degradation on each component. **Top:** We replace the amplitude and phase components of the original HR image ( $\mathcal{A}_{HR}$  and  $\mathcal{P}_{HR}$ ) with the corresponding components from the degraded LR image ( $\mathcal{A}_{LR}$  and  $\mathcal{P}_{LR}$ ). **Bottom:** Similarly, the original low- and high-frequency components  $\mathcal{L}_{HR}$  and  $\mathcal{H}_{HR}$  are replaced with  $\mathcal{L}_{LR}$  and  $\mathcal{H}_{LR}$ .

and facilitate readers’ understanding, we conducted additional experiments (see Figure 1 and Appendix B). From a technical perspective, researchers have explored various frequency-based ISR algorithms. However, early efforts focus on improving traditional model architectures like ResNet (He et al., 2016) and GANs (Goodfellow et al., 2014). Though several recent works (Luo et al., 2023; Wang et al., 2024b; Zhao et al., 2024; Moser et al., 2024) also explore improving diffusion-based ISR, they rely on heavy training processes and handcrafted network structure modifications.

In this paper, we propose a generic and *training-free* Frequency-aware Enhancement framework for Diffusion-based Super-Resolution (dubbed FedSR). Specifically, FedSR encapsulates two key components. **(1) Amplitude Enhancement Module:** To enhance the lost amplitude components, we develop an amplitude enhancement module (AEM) that utilizes a channel-aware mechanism to enhance the amplitude components which convey crucial details. **(2) High-Frequency Enhancement Module:** We further design a high-frequency enhancement module (HEM) that operates on the skip connection features, which employs a spectral modulation method to adaptively enhance the prominent high-frequency information in the skip features. The two modules can be simultaneously integrated into current diffusion-based ISR models, without requiring any further fine-tuning. Through extensive experiments, our FedSR significantly improves state-of-the-art diffusion-based ISR algorithms StableSR, PASD by +10.53%, +10.67% on MUSIQ metric, respectively. These results clearly validate the superiority of our FedSR algorithm in enhancing the amplitude and high-frequency details from a frequency perspective.

The main contributions of our work are as follows: **(A General framework)** We present a **general framework** that is able to improve most diffusion-based SR algorithms without extra training costs. **(Technical Novelty)** Motivated by our empirical findings, we propose a novel channel selection mechanism for enhancing the amplitude information. Also, we develop a new semantic-aware high-pass filtering algorithm that adaptively determines the thresholds by feature inputs. Again, we note that the two modules are **totally training-free**. **(Experiments)** We conduct extensive experiments on three benchmarks, demonstrating that FedSR improves 5 SOTA diffusion-based SR methods, verifying its generality. Moreover, we have no extra training cost, maintaining almost the same complexity parameters.

## 2 RELATED WORK

**Image Super-Resolution (ISR).** Although deep learning-based ISR techniques have gained widespread adoption, most CNN-based methods (Dong et al., 2016a; Lim et al., 2017; Kim et al., 2016; Dong et al., 2016b; Shi et al., 2016) still suffer from the issue of excessive detail smoothing. To address this problem and better enhance visual perception, recent advances (Zhang et al., 2021;

108 Wang et al., 2021; Liang et al., 2021; Chen et al., 2022; Liang et al., 2022; Wang et al., 2024a) using  
 109 the GAN-based models in the field of Real-ISR have explored more complex degradation models for  
 110 adversarial training. For instance, BSRGAN (Zhang et al., 2021) synthesizes more realistic degrada-  
 111 tion by using a random shuffling strategy, and RealeSRGAN (Wang et al., 2021) employs high-order  
 112 degradation modeling techniques. While these methods have made progress in generating more per-  
 113 ceptually realistic details, GAN-based ISR methods often suffer from unstable adversarial training,  
 114 frequently introducing unnatural visual artifacts. In recent years, the powerful Stable Diffusion (SD)  
 115 (Rombach et al., 2022) model has been applied to ISR tasks (Wang et al., 2023c; Lin et al., 2023; Yu  
 116 et al., 2024; Wu et al., 2023; Yang et al., 2023; Wang et al., 2023d; Cui et al., 2024). For instance,  
 117 PASD (Yang et al., 2023) utilizes pixel-aware cross attention to perceive image local structures.  
 118 SUPIR (Yu et al., 2024) develops a trimmed ControlNet (Zhang et al., 2023) and ZeroSFT to reduce  
 119 the model size. Although these methods demonstrate excellent performance in real-world ISR tasks,  
 120 they are limited to operations in the spatial domain and do not thoroughly explore the characteristics  
 121 of the frequency domain. In contrast, we discuss the degradation processes of various frequency  
 122 components and design a training-free method to enhance these degraded components.

123 **Frequency-based Super-Resolution.** Frequency analysis of image processing has been widely  
 124 used in computer vision (Yu et al., 2022; Huang et al., 2024; Yang & Soatto, 2020; Cai et al., 2021;  
 125 Si et al., 2023; Yu et al., 2022; Ji et al., 2021). For super-resolution tasks, many studies improve  
 126 images reconstruction quality by applying frequency domain transformations to comprehensively  
 127 extract feature information from low-resolution images (Guan et al., 2024; Li et al., 2023a; Xu et al.,  
 128 2024; Xie et al., 2021). Some methods enhance performance by constructing frequency domain loss  
 129 functions that focus on recovering frequency information through heavy network training (Zhu et al.,  
 130 2023; Fuoli et al., 2021; Dong et al., 2023; Ji et al., 2021; Wang et al., 2024c; Li et al., 2023b). For  
 131 example, Fuoli et al. (2021) designs Fourier space supervision losses to enhance perceptual quality in  
 132 image super-resolution. Additionally, some methods improve reconstruction quality by separating  
 133 specific components (such as high-frequency components) in the frequency domain (Guan et al.,  
 134 2024; Li et al., 2023a; Xu et al., 2024; Xie et al., 2021; Dai et al., 2024; Yang et al., 2022a; Jiang  
 135 et al., 2023). Appendix A.1 lists the effects of different frequency components on image quality  
 136 for other computer vision tasks. Although these existing frequency domain-based ISR methods  
 137 significantly improve performance, they have two main drawbacks: first, they rely on frequency  
 138 domain loss functions to achieve realistic outcomes with heavy training; second, they typically focus  
 139 only on certain specific components in the frequency domain. In contrast, our training-free FedSR  
 140 systematically analyzes the degradation process from the perspective of image modeling and then  
 141 enhances these degraded components.

### 142 3 BACKGROUND AND PRELIMINARIES

#### 143 3.1 DIFFUSION MODELS FOR IMAGE SUPER-RESOLUTION

144 Diffusion models, such as DDPM (Ho et al., 2020) and LDM (Rombach et al., 2022), are a class of  
 145 latent variable models, which primarily consist of a diffusion process and a denoising process. In the  
 146 diffusion process, Gaussian noise is gradually added at each time step  $t$  according to a predefined  
 147 variance schedule denoted as  $\beta_1, \dots, \beta_t$ , via a Markov chain. It eventually results in a random noise  
 148 distribution, which is defined as,  
 149

$$150 q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}\left(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathcal{I}\right). \quad (1)$$

151 In the denoising process, given the noisy input  $\mathbf{x}_t$ , the model outputs the clean data  $\mathbf{x}_{t-1}$  before  
 152 noise is added, which is represented as,

$$153 p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)). \quad (2)$$

154 Here,  $\boldsymbol{\mu}_\theta$  and  $\boldsymbol{\Sigma}_\theta$  are determined by the denoising model. Current diffusion-based generative models  
 155 (Ho et al., 2020; Rombach et al., 2022) are implemented using a U-Net (Ronneberger et al., 2015)  
 156 architecture to remove noise from data samples, which consists of a contracting path for downsam-  
 157 pling and an expansive path for upsampling. Each upsampling block concatenates both the backbone and  
 158 skip features in the skip connections.  
 159

160 To ensure that diffusion-based generative models meet the requirements for image quality and fi-  
 161 delity in ISR tasks, existing methods typically utilize LR images to guide model generation. First,

the LR image is used as a conditional input and transformed into an embedding through the image encoder. Then, these embeddings are fused with the U-Net using a cross-attention mechanism or a custom control module to guide the generation of HR images. Through iterative diffusion and reverse processes, these models effectively capture complex image features, enhancing the capability to recover realistic details.

### 3.2 FOURIER FREQUENCY DOMAIN TRANSFORMATION

The Fast Fourier Transform (FFT) is widely applied in low-level vision tasks, transforming images from the spatial domain to the Fourier domain, denoted as,

$$\mathcal{F}(\mathbf{x})(u, v) = \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \mathbf{x}(h, w) e^{-j2\pi(\frac{h}{H}u + \frac{w}{W}v)}. \quad (3)$$

Its inverse function (IFFT) is formulated as,

$$\mathcal{G}(\mathbf{f})(h, w) = \frac{1}{UV} \cdot \sum_{u=0}^{U-1} \sum_{v=0}^{V-1} \mathbf{f}(u, v) e^{-j2\pi(\frac{u}{U}h + \frac{v}{V}w)}, \quad (4)$$

where  $j$  is the imaginary unit;  $e$  is Euler's number, which can be formulated as  $e^{j\theta} = \cos\theta + j\sin\theta$ .  $\mathcal{F}(\cdot)$  and  $\mathcal{G}(\cdot)$  are 2D Fourier transform and inverse 2D Fourier transform, respectively. The frequency features  $\mathcal{F}(\mathbf{x})$  in Eq. (3) and  $\mathbf{f}$  in Eq. (4) are both tensors in complex domain, expressed as  $\mathcal{F}(\mathbf{x}) = \mathcal{R}(\mathbf{x}) + j\mathcal{I}(\mathbf{x})$ , where  $\mathcal{R}(\mathbf{x})$  and  $\mathcal{I}(\mathbf{x})$  are the real parts and imaginary parts, respectively.

In this paper, we explore two decomposition methods in the frequency domain, and the related analysis refers to Appendix A. The first is composition-based decomposition, which separates frequency into the amplitude  $\mathcal{A}$  and phase  $\mathcal{P}$ , represented as,

$$\begin{aligned} \mathcal{A}(\mathbf{x})(u, v) &= \sqrt{\mathcal{R}^2(\mathbf{x})(u, v) + \mathcal{I}^2(\mathbf{x})(u, v)}, \\ \mathcal{P}(\mathbf{x})(u, v) &= \arctan\left[\frac{\mathcal{I}(\mathbf{x})(u, v)}{\mathcal{R}(\mathbf{x})(u, v)}\right]. \end{aligned} \quad (5)$$

The other method is distance-based decomposition, where we divide the frequency information into high-frequency and low-frequency parts based on their distance from the frequency center.

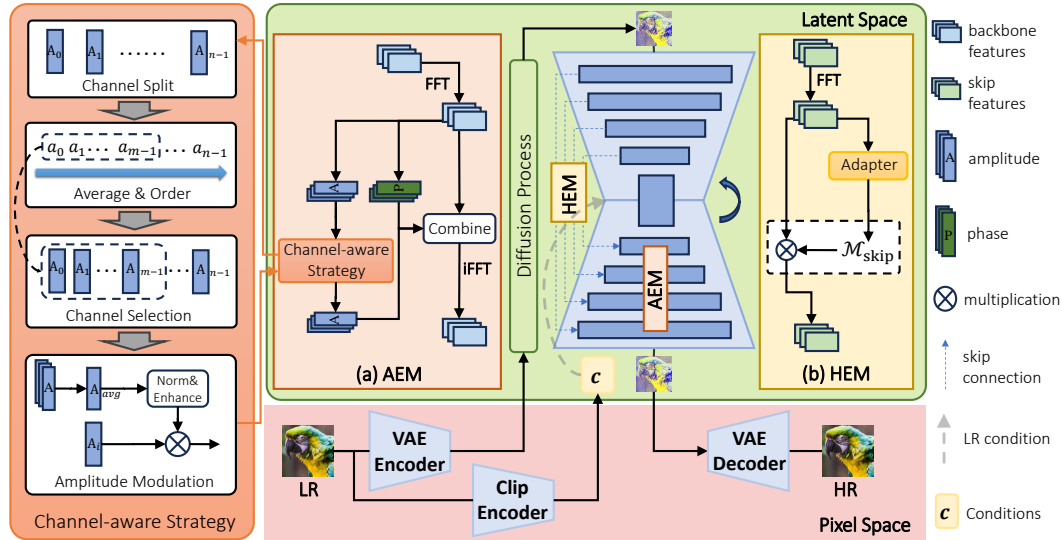


Figure 2: The overview of FedSR, which has two modules, (a) AEM: a channel-aware amplitude enhancement module which selectively enhances crucial amplitude channels through reweighting strategy; (b) HEM: a high-frequency enhancement module which utilizes adaptive masking.

## 4 THE PROPOSED FRAMEWORK

In this section, we describe our novel FedSR framework in detail. Essentially, FedSR comprises two key components: an amplitude modulation module that operates on the backbone features (Section 4.1), and a spectral modulation module designed for adaptive enhancement of high-frequency components (Section 4.2). Importantly, both modules are post-hoc adjustments to diffusion models, requiring no additional heavy tuning. Notably, FedSR can be seamlessly integrated as plugins into any off-the-shelf diffusion-based ISR models. The overall architecture is shown in Figure 2.

### 4.1 CHANNEL-AWARE AMPLITUDE ENHANCEMENT

In our preliminary experiments, we demonstrate the phenomenon of amplitude degradation in LR images. Thus, during the training process, the ISR model would actively learn the amplitude signal features from HR images. However, due to the black-box nature of DNNs, enhancing amplitude features cannot be directly achieved by simply following traditional image processing conclusions and requires further exploration; see Appendix D.1 for further discussion.

**Analysis of Amplitude Channels.** Inspired by some studies (Hu et al., 2018; Zhao et al., 2019), which enhance model performance by adjusting the importance of different channel features in convolutional neural networks (CNNs), we hypothesize that the amplitude features of various channels in the U-Net backbone network which contains convolutional layers, may also convey information of varying significance. To validate this hypothesis, we transform the image features generated by Stable Diffusion (Rombach et al., 2022) into the frequency domain and then select the channels according to their amplitude values. Figure 3 (a) illustrates the reconstructed images using different channels at different sampling steps. Our observations reveal that amplitude channels with lower amplitude values convey crucial details of the image, while channels with higher amplitude values result in disorganized and chaotic images, indicating that these channels have learned meaningless signals. This underscores the importance of emphasizing significant amplitude features during the sampling of ISR models to enhance image quality.

Based on this finding, we develop a simple yet effective channel-aware Amplitude Enhancement Module (AEM) aimed at selectively modulating the amplitude information in the backbone network by identifying channels with rich information, thereby improving the overall visual quality of the images. Technically, the AEM first transforms the U-Net backbone features before the concatenation of skip connections in upsampling blocks into the frequency domain. Then it extracts the amplitude components with Eq. (5) as the optimization target. Subsequently, we design the aforementioned channel-aware strategy, which consists of the following four steps.

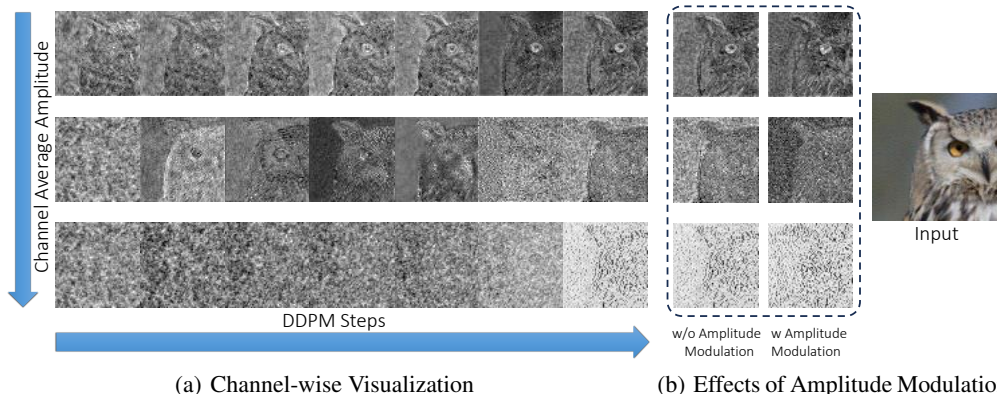


Figure 3: Average amplitude on features and effect of amplitude modulation. (a) During the generation process of the diffusion model, lower average amplitude in the channel leads to clearer generated images. (b) The application of our amplitude modulation further enhances feature clarity.

270 **a) Channel Separation.** Inspired by SENet’s (Hu et al., 2018) different processing of features  
 271 across varying channels, we split the amplitude component along the channel dimension and then  
 272 obtain the amplitude set with all channels, denoted as  $\mathcal{S}_A = \{\mathcal{A}(\mathbf{x}_{\text{bone}})_i\}_{i=1}^C$ , to separate various  
 273 pieces of information, where  $\mathbf{x}_{\text{bone}}$  is the backbone features;  $C$  is the number of amplitude channels.  
 274

275 **b) Average Amplitude Value Ranking.** Recall that channels with lower average amplitude gener-  
 276 ally exhibit clearer details in Figure 3 (a). Therefore, we compute the average value of the ampli-  
 277 tude component for each channel, formulated as  $a_i = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W \mathcal{A}(\mathbf{x}_{\text{bone}})_i^{(h,w)}$ , where  
 278  $\mathcal{A}(\mathbf{x}_{\text{bone}})_i = (\mathcal{A}(\mathbf{x}_{\text{bone}})_i^{(h,w)})_{H \times W}$  with  $H$  as the height of the feature and  $W$  as the width of the  
 279 feature. Then we rank the amplitude of each channel  $\mathcal{A}(\mathbf{x}_{\text{bone}})_i$  in  $\mathcal{S}_A$  in ascending order based on  
 280 the average value  $a_i$ , to identify channels with richer detailed information.  
 281

282 **c) Channel Selection.** Based on the ranking results, we select channels with lower ampli-  
 283 tude values that contain abundant information and then combine them into a subset  $\mathcal{S}_S =$   
 284  $\{\mathcal{A}(\mathbf{x}_{\text{bone}})_i | a_i \leq a_{\min} + P_s \times (a_{\max} - a_{\min})\}$  for subsequent amplitude modulation. Here,  $P_s$  is the  
 285 selection thresholds. To achieve better results of AEM, we design a supplementary experiment ex-  
 286 haustively testing  $P_s$  between 0 and 1 on the validation set, shown in Section 5.3.  
 287

288 **d) Amplitude Modulation.** To amplify the impact of these selected amplitude channels, we apply  
 289 amplitude reweight at the final step of sampling. Specifically, to align the amplitude components  
 290 with their size characteristics, we first compute the average amplitude  $\bar{\mathcal{A}} = \frac{1}{C} \sum_{i=1}^C \mathcal{A}(\mathbf{x}_{\text{bone}})_i$   
 291 along the channels, followed by linear normalization to construct a factor map  
 292

$$293 \mathcal{M}_{\text{bone}} = 1 - P_a \cdot \frac{\bar{\mathcal{A}} - \bar{\mathcal{A}}_{\min}}{\bar{\mathcal{A}}_{\max} - \bar{\mathcal{A}}_{\min}}, \quad (6)$$

294 where  $\bar{\mathcal{A}}_{\min}$  and  $\bar{\mathcal{A}}_{\max}$  means the minimum and maximum of  $\bar{\mathcal{A}}$ ;  $P_a$  is a positive linearization pa-  
 295 rameter. We then multiply the factor map  $\mathcal{M}_{\text{bone}}$  with the channel features in the subset  $\mathcal{S}_S$  one by  
 296 one, formulated as follows,  
 297

$$298 \mathcal{A}(\mathbf{x}_{\text{bone}})_i' = \begin{cases} \mathcal{A}(\mathbf{x}_{\text{bone}})_i \odot \mathcal{M}_{\text{bone}}, & \text{if } i \in \mathcal{S}_S; \\ \mathcal{A}(\mathbf{x}_{\text{bone}})_i, & \text{otherwise.} \end{cases} \quad (7)$$

299 To apply the enhanced amplitude, we use the modulated amplitude and the original phase compo-  
 300 nents to combine into the frequency domain by  $\mathcal{F}(\mathbf{x}_{\text{bone}})' = \mathcal{R}(\mathbf{x}_{\text{bone}})' + j\mathcal{I}(\mathbf{x}_{\text{bone}})'$ , and further  
 301 transfer to the spatial domain by the inverse Fourier transformation  $\mathbf{x}'_{\text{bone}} = \mathcal{G}(\mathcal{F}(\mathbf{x}_{\text{bone}})')$ .  
 302

303 There is a subtle point worth deeper discussion. At first glance, it might seem counterintuitive that  
 304 reducing the amplitude values in Eq. (6) and Eq. (7) would enhance image super-resolution (ISR).  
 305 However, our further experiments (see Appendix D.1) indicate that the logic behind traditional image  
 306 processing may differ from that of diffusion networks. Increasing the amplitude of the original  
 307 image signal typically affects image contrast and brightness. In contrast, within FedSR, reducing the  
 308 amplitude of the deep feature signals results in clearer detail. We speculate the reason might be that  
 309 diffusion models are prone to highlight channels with smaller amplitude values. Further exploration  
 310 of this behavior requires deeper theoretical insights from diffusion models in the frequency domain,  
 311 which we leave for future work.  
 312

## 313 4.2 ADAPTIVE MASKING FOR HIGH-FREQUENCY ENHANCEMENT

314 Next, we discuss our modification to diffusion models to enhance ISR performance from the per-  
 315 spective of high-frequency details. Inspired by FreeU (Si et al., 2023), we know that the skip connec-  
 316 tions in U-Net blocks can transmit high-frequency, information-rich features to deeper layers of the  
 317 network, thereby preserving more comprehensive image information. Note that FreeU is designed  
 318 for text-to-image tasks which only applies two constant scaling transformations to low-frequency  
 319 features on all layers to achieve high-pass filtering. However, for the diffusion-based ISR problems,  
 320 the features on different U-Net layers convey various semantic information. Therefore, considering  
 321 the varying richness of information, we propose a high-frequency enhancement module (HEM) with  
 322 adaptive masking, which can be divided into the following two steps; see Figure 2 (b).  
 323

Table 1: Quantitative comparison with SOTA methods on the synthetic benchmark DIV2K-Val (Agustsson & Timofte, 2017). **Bold** and  $\Delta$  represent the improvement and the performance boost brought by FedSR, respectively. **Red** and **blue** colors represent the best and second-best performance.  $\downarrow$  represents the smaller the better, while  $\uparrow$  represents the opposite.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
<b>StableSR</b> (IJCV2024)	23.26	0.5644	0.3119	0.6771	65.91	4.742	0.4208
<b>FedSR+StableSR</b>	22.59	0.553	0.3711	0.7275	<b>71.48</b>	<b>4.112</b>	0.4914
$\Delta$ StableSR	-0.67	-0.0114	+0.0592	<b>+0.0504</b>	<b>+5.57</b>	<b>-0.63</b>	<b>+0.0706</b>
<b>SUPIR</b> (CVPR2024)	22.14	0.5180	0.3930	0.7130	63.60	5.705	0.5533
<b>FedSR+SUPIR</b>	20.98	0.4847	0.4125	<b>0.7368</b>	66.57	5.437	<b>0.5841</b>
$\Delta$ SUPIR	-1.16	-0.0333	+0.0195	<b>+0.0238</b>	<b>+2.97</b>	<b>-0.268</b>	<b>+0.0308</b>
<b>SeeSR</b> (CVPR2024)	23.67	0.5978	0.3200	0.6940	68.72	4.806	0.5044
<b>FedSR+SeeSR</b>	23.56	0.6101	0.3401	0.6893	70.21	<b>4.598</b>	0.5184
$\Delta$ SeeSR	-0.11	+0.0123	+0.0201	-0.0047	<b>+1.49</b>	<b>-0.208</b>	<b>+0.0140</b>
<b>PASD</b> (ECCV2024)	24.16	0.6099	0.3705	0.5848	61.85	5.169	0.4028
<b>FedSR+PASD</b>	24.25	0.6213	0.3644	0.5948	65.72	4.904	0.4223
$\Delta$ PASD	+0.09	+0.0114	-0.0061	<b>+0.0100</b>	<b>+3.87</b>	<b>-0.265</b>	<b>+0.0195</b>
<b>DiffBIR</b> (Arxiv2023)	23.14	0.5370	0.3667	0.7301	69.90	4.991	0.5675
<b>FedSR+DiffBIR</b>	22.38	0.5222	0.4236	<b>0.7382</b>	<b>73.04</b>	4.729	<b>0.5838</b>
$\Delta$ DiffBIR	-0.76	-0.0148	+0.0569	<b>+0.0081</b>	<b>+3.14</b>	<b>-0.262</b>	<b>+0.0163</b>

**a) Adaptive Mask Construction.** To accurately filter and dynamically enhance the high-frequency components in the skip features, we construct an adaptive high-frequency mask  $\mathcal{M}_{\text{skip}}$ . Considering that lower-level and smaller-scale features often contain less image detailed information, the mask adjusts the enhancement factor based on scale adaptively, to better adapt to the frequency structure of features at different levels, formulated as,

$$\mathcal{M}_{\text{skip}}(r) = \begin{cases} 1 + \left(\frac{S - S_{\min}}{S_{\max} - S_{\min}} + 0.5\right) \cdot \frac{P_b}{2}, & \text{if } r > r_{\text{thresh}}; \\ 1, & \text{otherwise.} \end{cases} \quad (8)$$

Here  $S$  is the scale of skip features, and  $P_b$  is the enhancement factor;  $r$  and  $r_{\text{thresh}}$  are the radius and the radius threshold, respectively. Thus, the high-frequency components are split through masking.

**b) High-Frequency Component Enhancement.** We then multiply the adaptive mask  $\mathcal{M}_{\text{skip}}$  element-wise with the skip features  $\mathbf{x}_{\text{skip}}$  in the frequency domain to amplify and enhance the high-frequency components, represented as,

$$\mathcal{F}(\mathbf{x}_{\text{skip}})' = \mathcal{F}(\mathbf{x}_{\text{skip}}) \odot \mathcal{M}_{\text{skip}}, \quad (9)$$

where  $\odot$  denotes element-wise multiplication. Finally, the inverse Fourier transformation, which is denoted as  $\mathbf{x}'_{\text{skip}} = \mathcal{G}(\mathcal{F}(\mathbf{x}_{\text{skip}})')$ , transfers the enhanced skip features to the spatial feature domain.

**Remark.** In practical applications, the AEM and HEM modules can actually be integrated into any layer of the diffusion U-Net blocks. However, our experimental validation shows that the better setup is to apply the AEM to the backbone features and the HEM to the skip features, as this configuration consistently yields superior performance; we may refer the readers to Appendix D.2 for more discussion. Empirically, both modules can be simultaneously incorporated into diffusion-based ISR models without the need for additional fine-tuning or adjustments. On various ISR benchmarks, FedSR achieves significant performance gains, effectively offering a *free lunch* for ISR.

## 5 EXPERIMENTS

### 5.1 EXPERIMENTAL SETTINGS

**Datasets and Baselines.** We employ the test datasets from StableSR (Wang et al., 2023c) and evaluate our approach on both synthetic and real-world datasets. (1) For the synthetic dataset, we

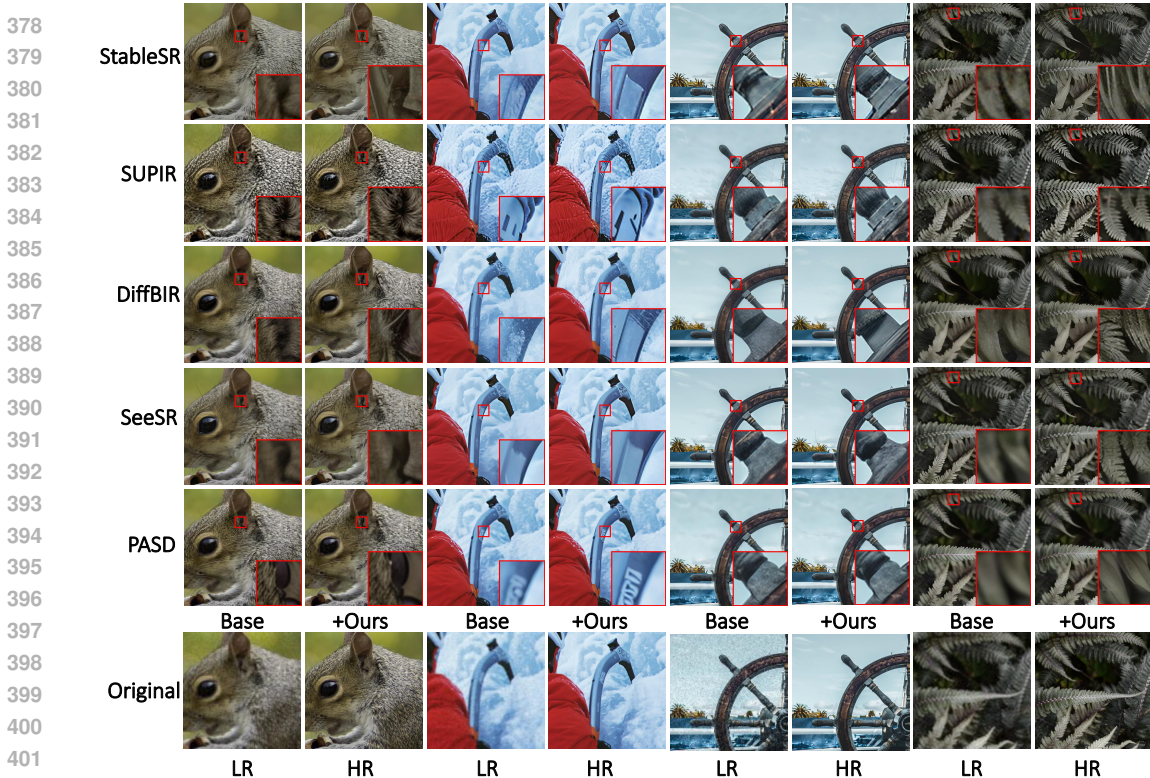


Figure 4: Qualitative comparisons of diffusion model-based ISR methods before and after incorporating our FedSR. It shows that FedSR can reconstruct more realistic HR images.

use 3,000 generated pairs of LR-HR images from the DIV2K validation set (Agustsson & Timofte, 2017), where the LR images have a resolution of  $128 \times 128$ , and the HR images have a resolution of  $512 \times 512$ . (2) For the real-world datasets, we utilize the DRealSR (Wei et al., 2020) and RealSR (Cai et al., 2019) datasets center-cropping the LR images to  $128 \times 128$ .

We select five state-of-the-art (SOTA) diffusion-based ISR models, namely StableSR (Wang et al., 2023c), SUPIR (Yu et al., 2024), SeeSR (Wu et al., 2023), PASD (Yang et al., 2023), and DiffBIR (Lin et al., 2023). And we incorporate our FedSR into these frameworks to evaluate effectiveness.

**Evaluation Metrics.** We adopt a series of full-reference and no-reference metrics to assess the performance of different methods. The full-reference metrics includes PSNR, SSIM (evaluated on the Y channel in the YCbCr color space), and LPIPS (Zhang et al., 2018). For quality evaluation, we employ no-reference image quality assessment (IQA) metrics: CLIP-IQA (Wang et al., 2023b), MUSIQ (Ke et al., 2021), NIQE (Zhang et al., 2015), and MANIQA (Yang et al., 2022b).

**Implementation Details.** To obtain the validation set of LR-HR pairs, we employ the degradation process of BSRGAN (Zhang et al., 2021) on the small random subset of size 100 from the DIV2K training set. Then we adjust the hyper-parameters in FedSR. Based on the default settings (i.e.,  $P_a = 0.5$ ,  $P_{b1} = 1$ , and  $P_{b2} = 1$ ) as default, for further tuning. To determine the selection threshold  $P_s$ , we experiment with  $P_s \in [0, 1]$  (see Figure 6) on a validation set created also by randomly selecting from DIV2K. We find that as the channel selection threshold increases, various metrics gradually stabilize, and set  $P_s$  as 0.3 for better performance, further indicating larger amplitude channels contribute less. Detailed hyper-parameter settings can be found in the Appendix D.

## 5.2 COMPARISON BEFORE AND AFTER FEDSR APPLICATION

**Quantitative Comparisons.** As shown in Table 1, we apply our method to five SOTA diffusion-based ISR frameworks, and the results on DIV2K-Valid indicate that almost all no-reference metrics



Table 2: Quantitative results on the real-world benchmark DRealSR with our FedSR.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
StableSR (IJCV2024)	27.93	0.7442	0.3280	0.6272	58.28	6.475	0.3890
<b>FedSR+StableSR</b>	26.68	0.7206	0.3903	0.6690	67.27	<b>5.373</b>	0.4810
$\Delta$ StableSR	-1.25	-0.0236	+0.0623	<b>+0.0418</b>	<b>+8.99</b>	<b>-1.102</b>	<b>+0.0920</b>
SUPIR (CVPR2024)	24.80	0.6333	0.4323	0.6880	59.73	7.420	0.5040
<b>FedSR+SUPIR</b>	23.18	0.5777	0.4622	<b>0.7232</b>	64.09	6.810	<b>0.5584</b>
$\Delta$ SUPIR	-1.62	-0.0556	+0.0299	<b>+0.0352</b>	<b>+4.36</b>	<b>-0.610</b>	<b>+0.0544</b>
SeeSR (CVPR2024)	28.04	0.7661	0.3188	0.6924	65.08	6.389	0.5134
<b>FedSR+SeeSR</b>	27.33	0.7671	0.3422	0.6944	<b>67.46</b>	6.052	0.5300
$\Delta$ SeeSR	-0.71	+0.0010	+0.0234	<b>+0.0020</b>	<b>+2.38</b>	<b>-0.337</b>	<b>+0.0166</b>
PASD (ECCV2024)	28.96	0.7919	0.3142	0.5122	52.29	6.929	0.3672
<b>FedSR+PASD</b>	28.28	0.7860	0.3203	0.5790	62.16	6.420	0.4232
$\Delta$ PASD	-0.68	-0.0059	+0.0061	<b>+0.0668</b>	<b>+9.87</b>	<b>-0.509</b>	<b>+0.0560</b>
DiffBIR (Arxiv2023)	25.90	0.6220	0.4715	0.7076	66.22	6.309	0.5568
<b>FedSR+DiffBIR</b>	24.53	0.6014	0.5024	<b>0.7167</b>	<b>71.90</b>	<b>5.833</b>	<b>0.5902</b>
$\Delta$ DiffBIR	-1.37	-0.0206	+0.0309	<b>+0.0091</b>	<b>+5.68</b>	<b>-0.476</b>	<b>+0.0334</b>

improved. It suggests that our FedSR can further enhance image quality within these existing frameworks. Table 2 and Table 3 present the results on real-world datasets. For example, on the DIV2K-Val dataset, FedSR improves the original StableSR by **+7.44%** on CLIP-IQA metric. Additionally, on the real-world datasets DRealSR and RealSR, our method improves PASD by **+18.88%** and **+6.88%** on MUSIQ metric, respectively, thus demonstrating the effectiveness of FedSR. Although our method does not show significant improvements in full-reference metrics (PSNR, SSIM, and LPIPS), these metrics only capture certain aspects of performance (Blau & Michaeli, 2018; Ledig et al., 2017). Moreover, Figure 5 shows that solely pursuing improvements in these traditional metrics does not necessarily lead to better visual effects. **Our FedSR, while maintaining reasonable PSNR/SSIM, significantly enhances no-reference metrics (largely improved MUSIQ +10.53%).**



Figure 5: Ours (FedSR+DiffBIR) generates images with better image quality but obtains lower metrics in PSNR, SSIM, and LPIPS, which shows the bias between metric evaluation and image quality.

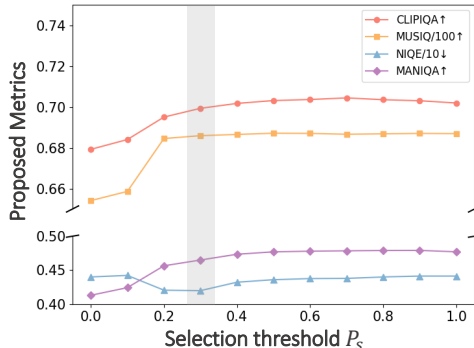


Figure 6: Selection thresholds and performance on our validation set (a random subset DIV2K training data). The gray column represents the best selection threshold  $P_s$ .

**Qualitative Comparisons.** To demonstrate the effectiveness of FedSR, Figure 4 presents a comparison before and after incorporating FedSR. It can be observed that our method significantly enhances the quality of the image generated by diffusion-based ISR methods, particularly in detailed textures and general visual effects. An interesting observation is that there occurs pseudo-textures in some images (e.g., squirrels) when applying FedSR to baselines like SeeSR. However, a closer inspection shows that the original baselines already demonstrate pseudo textures (though

Table 3: Quantitative results on the real-world benchmark RealSR with our FedSR.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
StableSR (IJCV2024)	24.66	0.7003	0.3101	0.6169	65.24	5.924	0.4302
FedSR+StableSR	23.77	0.6832	0.3502	0.6683	70.27	<u>5.094</u>	0.5186
$\Delta$ StableSR	-0.89	-0.0171	+0.0401	<b>+0.0514</b>	<b>+5.03</b>	<b>-0.830</b>	<b>+0.0884</b>
SUPIR (CVPR2024)	23.64	0.6603	0.3511	0.6316	61.34	6.299	0.4952
FedSR+SUPIR	22.25	0.6173	0.3727	<u>0.6807</u>	65.57	5.685	<u>0.5625</u>
$\Delta$ SUPIR	-1.39	-0.043	+0.0216	<b>+0.0491</b>	<b>+4.23</b>	<b>-0.614</b>	<b>+0.0673</b>
SeeSR (CVPR2024)	25.14	0.7182	0.2996	0.6697	69.86	5.419	0.5437
FedSR+SeeSR	24.73	0.7258	0.2982	0.6612	<u>70.93</u>	5.280	0.5591
$\Delta$ SeeSR	-0.41	+0.0076	-0.0014	-0.0085	<b>+1.07</b>	<b>-0.139</b>	<b>+0.0154</b>
PASD (ECCV2024)	26.53	0.7597	0.2783	0.5030	60.61	6.018	0.3894
FedSR+PASD	26.15	0.7596	0.2751	0.5191	64.78	5.744	0.4199
$\Delta$ PASD	-0.38	-0.0001	-0.0032	<b>+0.0161</b>	<b>+4.17</b>	<b>-0.274</b>	<b>+0.0305</b>
DiffBIR (Arxiv2023)	24.83	0.6473	0.3678	0.7017	69.22	5.812	0.5584
FedSR+DiffBIR	23.97	0.6405	0.3667	<b>0.7090</b>	<b>72.83</b>	<b>5.068</b>	<b>0.5812</b>
$\Delta$ DiffBIR	-0.86	-0.0068	-0.0011	<b>+0.0073</b>	<b>+3.61</b>	<b>-0.744</b>	<b>+0.0228</b>

being blurry). Since FedSR does not modify the original parameters of the models, these erroneous textures are inadvertently amplified. However, for models such as PASD and SUPIR, we successfully preserve the natural fur texture while simultaneously enhancing the quality of other fine details. In summary, with better baselines, FedSR is able to output much more realistic details. And one may also regard our FedSR as a detector to verify the true ISR ability of baseline models.

### 5.3 ABLATION STUDY

In this section, we present our ablation results on StableSR to show the effectiveness of FedSR. First, we validate the effectiveness of the AEM in ISR tasks. Compared to the default settings, removing the AEM results in poorer no-reference metrics (see Row 3 of Table 4), while adding it leads to a noticeable improvement in no-reference metrics. For more visual results, please refer to the Appendix E. Next, to validate the effectiveness of the HEM, we removed this module and it results in worse no-reference metrics compared to the default settings (see Row 2 of Table 4). In contrast, simply adding the HEM leads to a noticeable improvement in no-reference metrics.

Table 4: Ablation studies of FedSR on DRealSR and RealSR benchmarks.

Variants		DRealSR/RealSR						
AEM	HEM	PSNR $\uparrow$	SSIM $\downarrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
		27.93 / 24.66	0.7442 / 0.7003	0.3280 / 0.3101	0.6272 / 0.6169	58.28 / 65.24	6.475 / 5.924	0.3890 / 0.4302
✓		27.36 / 24.27	0.7322 / 0.6890	0.3635 / 0.3389	<b>0.6760 / 0.6760</b>	65.23 / 69.59	5.707 / 5.241	0.4681 / 0.5087
	✓	26.91 / 23.85	0.7230 / 0.6859	0.3655 / 0.3351	0.6749 / 0.6522	64.56 / 68.55	5.768 / 5.386	0.4447 / 0.4738
✓	✓	26.68 / 23.77	0.7206 / 0.6832	0.3903 / 0.3502	0.6690 / 0.6683	<b>67.27 / 70.27</b>	<b>5.373 / 5.094</b>	<b>0.4810 / 0.5186</b>

## 6 CONCLUSION

In this work, we propose a *generic and training-free* framework FedSR for enhancing diffusion-based ISR models from a frequency perspective. To achieve this, we first propose a novel channel selection mechanism for enhancing the amplitude information (AEM). Also, we develop a new semantic-aware high-pass filtering algorithm that adaptively determines the thresholds by feature inputs (HEM). As shown in the extensive experimental evaluation, we demonstrate the effectiveness of the FedSR as a plug-in for most diffusion-based ISR models. Additionally, our analysis of the degradation on the frequency domain may also inspire other ISR models, e.g. GAN-based ISR models (see Appendix C.2). We also hope our work will draw more attention from the community toward a broader view of addressing low-level vision tasks like ISR from a frequency perspective.

## 540 ETHICAL STATEMENT

541

542 Although our proposed method does not strictly fall under generative AI, it can serve as a plug-and-  
 543 play framework integrated into diffusion-based ISR algorithms developed by Wang et al. (2023c).  
 544 As diffusion models evolve toward aligning with human preferences, concerns regarding their po-  
 545 tential misuse and malicious purposes (such as generative discrimination or inappropriate content)  
 546 become increasingly prominent. Regarding other potential societal consequences of our work, none  
 547 of which we feel must be specifically highlighted here.

## 548 REPRODUCIBILITY STATEMENT

549

550 We provide our implementation details, including the main algorithm and parameters, which can be  
 551 found in Section 5 and Appendix D. Additionally, our source code is available in the supplementary  
 552 materials. This information provides the necessary resources for reproducing our results.

553

## 554 REFERENCES

555

- 556 Eirikur Agustsson and Radu Timofte. NTIRE 2017 challenge on single image super-resolution:  
 557 Dataset and study. In *CVPR*, pp. 1122–1131. IEEE Computer Society, 2017.
- 558 Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *CVPR*, pp. 6228–6237.  
 559 Computer Vision Foundation / IEEE Computer Society, 2018.
- 560 Jianrui Cai, Hui Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image  
 561 super-resolution: A new benchmark and a new model. In *ICCV*, pp. 3086–3095. IEEE, 2019.
- 562 Mu Cai, Hong Zhang, Huijuan Huang, Qichuan Geng, Yixuan Li, and Gao Huang. Frequency  
 563 domain image translation: More photo-realistic, better identity-preserving. In *ICCV*, pp. 13910–  
 564 13920. IEEE, 2021.
- 565 Chaofeng Chen, Xinyu Shi, Yipeng Qin, Xiaoming Li, Xiaoguang Han, Tao Yang, and Shihui Guo.  
 566 Real-world blind super-resolution via feature matching with implicit high-resolution priors. In  
 567 *ACM MM*, pp. 1329–1338. ACM, 2022.
- 568 Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chun-  
 569 jing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *CVPR*, pp.  
 570 12299–12310. Computer Vision Foundation / IEEE, 2021.
- 571 Julien Cornebise, Ivan Orsolice, and Freddie Kalaitzis. Open high-resolution satellite imagery: The  
 572 worldstrat dataset - with application to super-resolution. In *NeurIPS*, 2022.
- 573 Qinpeng Cui, Yixuan Liu, Xinyi Zhang, Qiqi Bao, Zhongdao Wang, Qingmin Liao, Li Wang, Tian  
 574 Lu, and Emad Barsoum. Taming diffusion prior for image super-resolution with domain shift  
 575 sdes. *arXiv preprint arXiv:2409.17778*, 2024.
- 576 Tao Dai, Jianping Wang, Hang Guo, Jinmin Li, Jinbao Wang, and Zexuan Zhu. Freqformer:  
 577 Frequency-aware transformer for lightweight image super-resolution. In *IJCAI*. ijcai.org, 2024.
- 578 Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep  
 579 convolutional networks. *IEEE TPAMI*, 38(2):295–307, 2016a.
- 580 Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional  
 581 neural network. In *ECCV*, volume 9906 of *Lecture Notes in Computer Science*, pp. 391–407.  
 582 Springer, 2016b.
- 583 Shuting Dong, Feng Lu, Zhe Wu, and Chun Yuan. DFVSR: directional frequency video super-  
 584 resolution via asymmetric and enhancement alignment network. In *IJCAI*, pp. 681–689. ijcai.org,  
 585 2023.
- 586 Dario Fuoli, Luc Van Gool, and Radu Timofte. Fourier space losses for efficient perceptual image  
 587 super-resolution. In *ICCV*, pp. 2340–2349. IEEE, 2021.

593

- 594 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
595 Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pp. 2672–2680,  
596 2014.
- 597 Wenxue Guan, Haobo Li, Dawei Xu, Jiabin Liu, Shenghua Gong, and Jun Liu. Frequency generation  
598 for real-world image super-resolution. *IEEE TCSVT*, 34(8):7029–7040, 2024.
- 600 Xin Guo, Xueyang Fu, Man Zhou, Zhen Huang, Jialun Peng, and Zheng-Jun Zha. Exploring fourier  
601 prior for single image rain removal. In *IJCAI*, pp. 935–941. ijcai.org, 2022.
- 602 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
603 nition. In *CVPR*, pp. 770–778. IEEE Computer Society, 2016.
- 604 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*,  
605 2020.
- 607 Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, pp. 7132–7141.  
608 Computer Vision Foundation / IEEE Computer Society, 2018.
- 609 Linjiang Huang, Rongyao Fang, Aiping Zhang, Guanglu Song, Si Liu, Yu Liu, and Hongsheng  
610 Li. Fouriscale: A frequency perspective on training-free high-resolution image synthesis. *CoRR*,  
611 abs/2403.12963, 2024.
- 613 Xiaozhong Ji, Guangpin Tao, Yun Cao, Ying Tai, Tong Lu, Chengjie Wang, Jilin Li, and Feiyue  
614 Huang. Frequency consistent adaptation for real world super resolution. In *AAAI*, pp. 1664–1672.  
615 AAAI Press, 2021.
- 616 Xinrui Jiang, Nannan Wang, Jingwei Xin, Keyu Li, Xi Yang, Jie Li, Xiaoyu Wang, and Xinbo Gao.  
617 Fabnet: Frequency-aware binarized network for single image super-resolution. *IEEE TIP*, 32:  
618 6234–6247, 2023.
- 619 Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. MUSIQ: multi-scale image  
620 quality transformer. In *ICCV*, pp. 5128–5137. IEEE, 2021.
- 622 Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for  
623 image super-resolution. In *CVPR*, pp. 1637–1645. IEEE Computer Society, 2016.
- 624 Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro  
625 Acosta, Andrew P. Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, and Wenzhe Shi. Photo-  
626 realistic single image super-resolution using a generative adversarial network. In *CVPR*, pp.  
627 105–114. IEEE Computer Society, 2017.
- 628 Ao Li, Le Zhang, Yun Liu, and Ce Zhu. Feature modulation transformer: Cross-refinement of global  
629 representation via high-frequency prior for image super-resolution. In *ICCV*, pp. 12480–12490.  
630 IEEE, 2023a.
- 632 Guangyuan Li, Chen Rao, Juncheng Mo, Zhanjie Zhang, Wei Xing, and Lei Zhao. Rethinking  
633 diffusion model for multi-contrast MRI super-resolution. *CoRR*, abs/2404.04785, 2024.
- 634 Jinmin Li, Tao Dai, Mingyan Zhu, Bin Chen, Zhi Wang, and Shu-Tao Xia. FSR: A general  
635 frequency-oriented framework to accelerate image super-resolution networks. In *AAAI*, pp. 1343–  
636 1350. AAAI Press, 2023b.
- 637 Ziqiang Li, Pengfei Xia, Xue Rui, and Bin Li. Exploring the effect of high-frequency components  
638 in gans training. *ACM TOMCCAP*, 19(5):153:1–153:22, 2023c.
- 639 Jie Liang, Hui Zeng, and Lei Zhang. Efficient and degradation-adaptive network for real-world  
640 image super-resolution. In *ECCV*, volume 13678 of *Lecture Notes in Computer Science*, pp.  
641 574–591. Springer, 2022.
- 642 Jingyun Liang, Jie Zhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir:  
643 Image restoration using swin transformer. In *ICCVW*, pp. 1833–1844. IEEE, 2021.
- 644 Yaoyuan Liang. Unsupervised super resolution reconstruction of traffic surveillance vehicle images.  
645 In *ICMLC*, pp. 336–341. ACM, 2021.
- 647

- 648 Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep resid-  
649 ual networks for single image super-resolution. In *CVPR*, pp. 1132–1140. IEEE Computer Soci-  
650 ety, 2017.
- 651 Xinqi Lin, Jingwen He, Ziyang Chen, Zhaoyang Lyu, Ben Fei, Bo Dai, Wanli Ouyang, Yu Qiao, and  
652 Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior. *CoRR*,  
653 abs/2308.15070, 2023.
- 654 Wu Liu, Xinchun Liu, Huadong Ma, and Peng Cheng. Beyond human-level license plate super-  
655 resolution with progressive vehicle search and domain prior GAN. In *ACM MM*, pp. 1618–1626.  
656 ACM, 2017.
- 657  
658 Feng Luo, Jinxi Xiang, Jun Zhang, Xiao Han, and Wei Yang. Image super-resolution via latent dif-  
659 fusion: A sampling-space mixture of experts and frequency-augmented decoder approach. *CoRR*,  
660 abs/2310.12004, 2023.
- 661  
662 Xintian Mao, Yiming Liu, Fengze Liu, Qingli Li, Wei Shen, and Yan Wang. Intriguing findings of  
663 frequency selection for image deblurring. In *AAAI*, pp. 1905–1913. AAAI Press, 2023a.
- 664  
665 Ye Mao, Lan Jiang, Xi Chen, and Chao Li. Disc-diff: Disentangled conditional diffusion model for  
666 multi-contrast MRI super-resolution. In *MICCAI*, volume 14229 of *Lecture Notes in Computer  
667 Science*, pp. 387–397. Springer, 2023b.
- 668  
669 Brian B. Moser, Stanislav Frolov, Federico Raue, Sebastian Palacio, and Andreas Dengel. Waving  
670 goodbye to low-res: A diffusion-wavelet approach for image super-resolution. In *IJCNN*, pp. 1–8.  
IEEE, 2024.
- 671  
672 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
673 resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10674–10685. IEEE, 2022.
- 674  
675 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-  
676 ical image segmentation. In *MICCAI*, volume 9351 of *Lecture Notes in Computer Science*, pp.  
234–241. Springer, 2015.
- 677  
678 Jacob Shermeyer and Adam Van Etten. The effects of super-resolution on object detection per-  
679 formance in satellite imagery. In *CVPR*, pp. 1432–1441. Computer Vision Foundation / IEEE,  
2019.
- 680  
681 Wenzhe Shi, Jose Caballero, Ferenc Huszar, Johannes Totz, Andrew P. Aitken, Rob Bishop, Daniel  
682 Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient  
683 sub-pixel convolutional neural network. In *CVPR*, pp. 1874–1883. IEEE Computer Society, 2016.
- 684  
685 Chenyang Si, Ziqi Huang, Yuming Jiang, and Ziwei Liu. Freeu: Free lunch in diffusion u-net. *CoRR*,  
abs/2309.11497, 2023.
- 686  
687 Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for  
688 image restoration. In *ICCV*, pp. 4549–4557. IEEE Computer Society, 2017.
- 689  
690 Boyang Wang, Fengyu Yang, Xihang Yu, Chao Zhang, and Hanbin Zhao. APISR: anime production  
691 inspired real-world anime super-resolution. *CoRR*, abs/2403.01598, 2024a.
- 692  
693 Chenyang Wang, Junjun Jiang, Zhiwei Zhong, and Xianming Liu. Spatial-frequency mutual learning  
694 for face super-resolution. In *CVPR*, pp. 22356–22366. IEEE, 2023a.
- 695  
696 Jianyi Wang, Kelvin C. K. Chan, and Chen Change Loy. Exploring CLIP for assessing the look and  
697 feel of images. In *AAAI*, pp. 2555–2563. AAAI Press, 2023b.
- 698  
699 Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin C. K. Chan, and Chen Change Loy. Exploit-  
700 ing diffusion prior for real-world image super-resolution. *CoRR*, abs/2305.07015, 2023c.
- 701  
702 Xingjian Wang, Li Chai, and Jiming Chen. Frequency-domain refinement with multiscale diffusion  
for super resolution. *CoRR*, abs/2405.10014, 2024b.
- 703  
704 Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind  
super-resolution with pure synthetic data. In *ICCVW*, pp. 1905–1914. IEEE, 2021.

- 702 Yufei Wang, Wenhan Yang, Xinyuan Chen, Yaohui Wang, Lanqing Guo, Lap-Pui Chau, Ziwei Liu,  
703 Yu Qiao, Alex C. Kot, and Bihan Wen. Sinsr: Diffusion-based image super-resolution in a single  
704 step. *CoRR*, abs/2311.14760, 2023d.
- 705 Zhengxue Wang, Zhiqiang Yan, and Jian Yang. Sgnet: Structure guided network via gradient-  
706 frequency awareness for depth map super-resolution. In *AAAI*, pp. 5823–5831. AAAI Press,  
707 2024c.
- 708 Pengxu Wei, Ziwei Xie, Hannan Lu, Zongyuan Zhan, Qixiang Ye, Wangmeng Zuo, and Liang Lin.  
709 Component divide-and-conquer for real-world image super-resolution. In *ECCV*, volume 12353  
710 of *Lecture Notes in Computer Science*, pp. 101–117. Springer, 2020.
- 711 Rongyuan Wu, Tao Yang, Lingchen Sun, Zhengqiang Zhang, Shuai Li, and Lei Zhang. Seers:  
712 Towards semantics-aware real-world image super-resolution. *CoRR*, abs/2311.16518, 2023.
- 713 Wenbin Xie, Dehua Song, Chang Xu, Chunjing Xu, Hui Zhang, and Yunhe Wang. Learning  
714 frequency-aware dynamic network for efficient super-resolution. In *ICCV*, pp. 4288–4297. IEEE,  
715 2021.
- 716 Yiran Xu, Taesung Park, Richard Zhang, Yang Zhou, Eli Shechtman, Feng Liu, Jia-Bin Huang, and  
717 Difan Liu. Videogigagan: Towards detail-rich video super-resolution. *CoRR*, abs/2404.12388,  
718 2024.
- 719 Mengping Yang, Zhe Wang, Ziqiu Chi, and Yanbing Zhang. Fregan: Exploiting frequency compo-  
720 nents for training gans under limited data. In *NeurIPS*, 2022a.
- 721 Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and  
722 Yujia Yang. MANIQA: multi-dimension attention network for no-reference image quality assess-  
723 ment. In *CVPR*, pp. 1190–1199. IEEE, 2022b.
- 724 Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Pixel-aware stable diffusion for realistic  
725 image super-resolution and personalized stylization. *CoRR*, abs/2308.14469, 2023.
- 726 Yanchao Yang and Stefano Soatto. FDA: fourier domain adaptation for semantic segmentation. In  
727 *CVPR*, pp. 4084–4094. Computer Vision Foundation / IEEE, 2020.
- 728 Fanghua Yu, Jinjin Gu, Zheyuan Li, Jinfan Hu, Xiangtao Kong, Xintao Wang, Jingwen He, Yu Qiao,  
729 and Chao Dong. Scaling up to excellence: Practicing model scaling for photo-realistic image  
730 restoration in the wild. *CoRR*, abs/2401.13627, 2024.
- 731 Hu Yu, Naishan Zheng, Man Zhou, Jie Huang, Zeyu Xiao, and Feng Zhao. Frequency and spa-  
732 tial dual guidance for image dehazing. In *ECCV*, volume 13679 of *Lecture Notes in Computer  
733 Science*, pp. 181–198. Springer, 2022.
- 734 Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation  
735 model for deep blind image super-resolution. In *ICCV*, pp. 4771–4780. IEEE, 2021.
- 736 Lin Zhang, Lei Zhang, and Alan C. Bovik. A feature-enriched completely blind image quality  
737 evaluator. *IEEE TIP*, 24(8):2579–2591, 2015.
- 738 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
739 diffusion models. In *ICCV*, pp. 3813–3824. IEEE, 2023.
- 740 Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
741 effectiveness of deep features as a perceptual metric. In *CVPR*, pp. 586–595. Computer Vision  
742 Foundation / IEEE Computer Society, 2018.
- 743 Chen Zhao, Weiling Cai, Chenyu Dong, and Chengwei Hu. Wavelet-based fourier information  
744 interaction with frequency diffusion adjustment for underwater image restoration. In *CVPR*, pp.  
745 8281–8291. IEEE, 2024.
- 746 Xiaole Zhao, Yulun Zhang, Tao Zhang, and Xueming Zou. Channel splitting network for single MR  
747 image super-resolution. *IEEE TIP*, 28(11):5649–5662, 2019.
- 748 Qiang Zhu, Pengfei Li, and Qianhui Li. Attention retractable frequency fusion transformer for image  
749 super resolution. In *CVPR*, pp. 1756–1763. IEEE, 2023.

## A IMAGE PROCESS IN FREQUENCY DOMAIN

### A.1 STUDIES ON IMAGE FREQUENCY DOMAIN ANALYSIS

**Applications of frequency components.** In recent years, frequency-domain information has widely applied in various computer vision tasks, with many studies exploring the impact of different components on image quality through various frequency decomposition methods (Yang & Soatto, 2020; Yu et al., 2022; Wang et al., 2023a; Si et al., 2023; Huang et al., 2024; Cai et al., 2021; Mao et al., 2023a). (1) Decomposition based on composition: Frequency-domain information can be divided into amplitude and phase spectra. FDA (Yang & Soatto, 2020) reduces the distribution discrepancy between the source and target domains by swapping their amplitude spectra. FSDGN (Yu et al., 2022) addresses the dehazing problem by investigating the correlation between amplitude and phase spectra in the frequency domain under foggy degradation. (2) Decomposition based on distance with the frequency center: The frequency domain can also be divided into high-frequency and low-frequency components. FreeU (Si et al., 2023) suppresses low-frequency features in the frequency domain to prevent Stable Diffusion from generating overly smooth images. FouriScale (Huang et al., 2024) applies low-pass filtering in the frequency domain to alleviate repetitive patterns and structural distortions in the generation of high-resolution images by pre-trained diffusion models. (3) Decomposition based on properties: The frequency domain can be separated into real and imaginary components. DeepRFT (Mao et al., 2023a) applies ReLU networks to the real and imaginary parts of the frequency domain separately to achieve effective image deblurring.

Table 5: Classification and Comparison of Frequency-Domain-Based Super-Resolution Methods.

Domain	Method	Amplitude and Phase Separate	High- and Low-Frequency Separate	Frequency Loss	Training-Free
ISR	FSN	×	✓	×	×
	FDC	×	✓	✓	×
	ARFFT	×	×	✓	×
	FADN	×	✓	✓	×
	CRAFT	×	✓	×	×
VSR	DFVSR	×	✓	✓	×
	FTVSR	×	✓	×	×
	VideoGigaGAN	×	✓	×	×
	MFPI	×	×	×	×
FSR	SFMNet	✓	×	✓	×
ISR	Ours	✓	✓	×	✓

**Other methods of frequency-based super-resolution.** The main paper summarizes existing methods that apply frequency transform to super-resolution tasks. Table 5 integrates and categorizes frequency-based super-resolution methods from multiple perspectives. It can be observed that other methods fail to consider degradation systematically. In contrast, our training-free method addresses degradation by modulating degraded amplitude and high-frequency components.

### A.2 IMAGE MODELING IN THE FREQUENCY DOMAIN

To better understand the semantic information represented by various frequency-domain components, we perform a visual modeling of them, shown in Figure 7. First, the image is transformed into the frequency domain, and then three types of decomposition are applied: (1) based on composition: amplitude and phase components; (2) based on distance: high- and low-frequency components; and (3) based on properties: real and imaginary components. Afterward, these components are transformed back into the spatial domain directly. The experimental results demonstrate that the amplitude and phase components, as well as the high- and low-frequency components, can convey the semantic information of the image. Specifically, the amplitude component primarily reflects the style characteristics of the image, such as color and contrast, while the phase component reveals the contour information. The low-frequency component captures the overall structure of the image, whereas the high-frequency component highlights the edges and texture details.

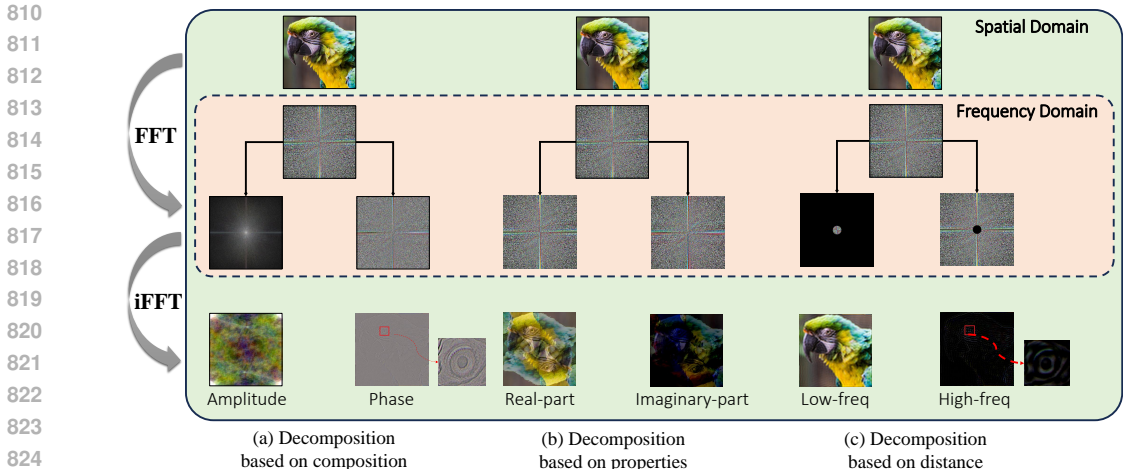


Figure 7: Image modeling methods in the frequency domain. We present the results of three decomposition approaches: (a) decomposition based on composition, (b) decomposition based on properties, and (c) decomposition based on distance. Compared to (b), the method (a) and (c) offer better separation of the intrinsic properties of the image.

## B QUANTITATIVE RESULTS OF THE IMPACT OF ISR DEGRADATION

In the main paper, the impact of the ISR degradation process on various components is visualized. Detailed quantitative results from testing on RealSR (Cai et al., 2019) are presented in Table 6. The first two rows of Table 6 show the results of replacing the amplitude and phase components of HR images with their corresponding LR counterparts. Following the replacement of the amplitude component, the resulting image quality metrics are poor, indicating that the information loss is primarily concentrated in the amplitude component. The last two rows of Table 6 display the outcomes of replacing the high- and low-frequency components of HR images with their LR counterparts. After the replacement of the high-frequency component, the image quality metrics also remained low, further confirming that the information loss is primarily concentrated in the high-frequency component.

Table 6: Quantitative results of the impact of ISR degradation on amplitude and phase components, high- and low-frequency components.

Method	CLIPQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
$\mathcal{A}_{LR} + \mathcal{P}_{HR}$	0.2320	25.44	6.849	0.1906
$\mathcal{A}_{HR} + \mathcal{P}_{LR}$	0.3060	28.78	6.410	0.2373
$\mathcal{H}_{LR} + \mathcal{L}_{HR}$	0.2731	25.54	9.99	0.2447
$\mathcal{H}_{HR} + \mathcal{L}_{LR}$	0.4447	56.98	6.035	0.3227

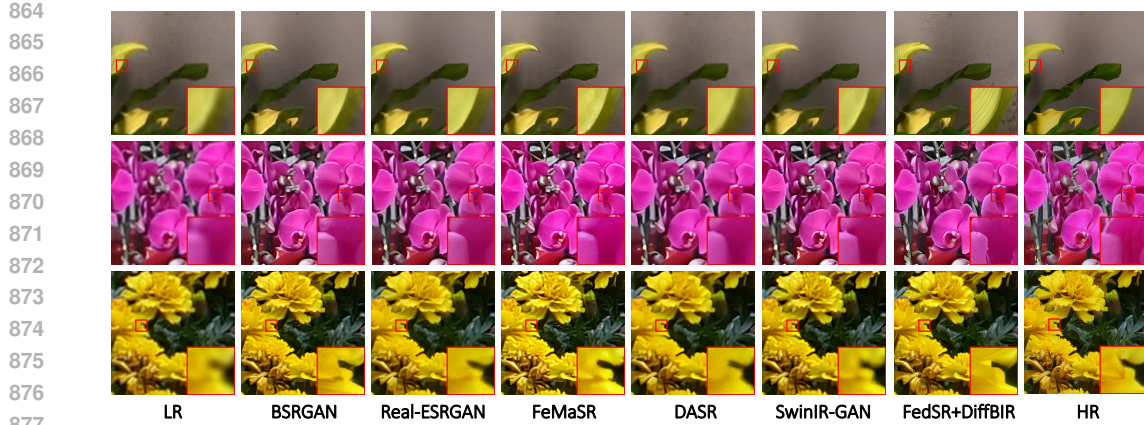
To further validate our argument regarding the impact of the ISR degradation, we performed the degradation process on phase and low-frequency components that have a minimal influence on the ISR task. Figure 9 demonstrates that the degradation of the phase component directly results in the loss of image structural information, and the degradation of the low-frequency leads to the disappearance of color information, rendering the issue no longer within the scope of ISR research.

## C COMPARE WITH GAN-BASED METHODS

### C.1 QUANTITATIVE AND QUALITATIVE COMPARISONS ON GAN-BASED MODELS

**GAN-based ISR Methods.** Based on the results in Table 1, 2, and 3, we further demonstrate the effectiveness of applying our method to DiffBIR and compare its superiority with GAN-based approaches, including BSRGAN (Zhang et al., 2021), Real-ESRGAN (Wang et al., 2021), FeMaSR





878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899

Figure 8: Qualitative comparisons of GAN-based methods and our FedSR applied to DiffBIR (Lin et al., 2023) on real-world examples.

(Chen et al., 2022), DASR (Liang et al., 2022), and SwinIR-GAN (Liang et al., 2021). We conduct the tests using publicly available codes and models from the comparison methods.

**Quantitative Comparisons.** Compared to GAN-based methods, our approach demonstrates superior performance in no-reference metrics across three datasets, as shown in Table 7. We also find that GAN-based methods generally perform better in PSNR and SSIM scores. This is because diffusion models generate more realistic details that may not perfectly match the ground truth (GT), thus leading to lower full-reference metrics compared to GAN-based methods.

**Qualitative Comparisons.** To validate the effectiveness of our method, we present a comparison between our approach and GAN-based methods in Figure 8. Our method has a significant advantage in detail generation. Specifically, the DiffBIR (Lin et al., 2023) model combined with our FedSR can produce sharp contours and realistic details, as shown in the second-to-last column of Figure 8, whereas other methods tend to generate blurred results.

## 900 C.2 DISCUSSION ON THE FEDSR APPLICATION EFFECTIVENESS OF GAN

901  
902  
903  
904  
905  
906  
907  
908  
909

According to Li et al. (2023c), in the training process of most GAN models, the discriminator tends to overemphasize high-frequency components, which weakens the generator’s ability to fit low-frequency components. As a result, while GANs can generate sharper images compared to Diffusion models, these images often exhibit unnatural details or artifacts. To explore the model’s generalizability, we applied the proposed method to the classic GAN-based ISR model BSRGAN (Zhang et al., 2021), by modulating its high-frequency and amplitude components to enhance the generated results and achieve fine control over each component. Specifically, we introduce the AEM module into the RRDB backbone network to adjust the amplitude components. We also incorporate the HEM module into the residual connections to reduce the impact of high-frequency components. Figure 11 presents the visual results, and the quantitative comparisons are detailed in Table 8.

## 910 D IMPLEMENTATION DETAILS

### 911 D.1 DISCUSSIONS OF OUR AMPLITUDE MODULATION

912  
913  
914  
915  
916  
917

**The relationship between low amplitude and high-frequency components.** Although physics indicates that, under the same energy (power), the amplitude of high-frequency waves is usually smaller than that of low-frequency waves, there is currently no evidence in frequency domain analysis of image processing to suggest a one-to-one correspondence between low-amplitude components and high-frequency details, especially in the feature maps of black-box deep learning models. For

Table 7: Quantitative comparison with other GAN-based methods on both synthetic and real-world benchmarks. The **bold** and underline represent the best and second-best performance, respectively.

DIV2K-Val							
Methods	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
BSRGAN	<b>24.57</b>	0.6232	0.3354	0.5255	<u>61.23</u>	4.751	0.3561
Real-ESRGAN	24.29	<b>0.6328</b>	<b>0.3115</b>	0.5283	61.11	<b>4.674</b>	<u>0.3823</u>
FeMaSR	23.05	0.5816	<u>0.3125</u>	<u>0.5997</u>	60.82	4.746	0.3457
DASR	<u>24.46</u>	<u>0.6267</u>	0.3542	0.5036	55.20	5.033	0.3186
SwinIR-GAN	23.92	0.6235	0.3159	0.5340	60.22	<u>4.706</u>	0.3656
DiffBIR+Ours	22.38	0.5222	0.4236	<b>0.7382</b>	<b>73.04</b>	4.729	<b>0.5838</b>
DrealSR							
Methods	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
BSRGAN	<u>28.68</u>	0.8021	0.2885	0.5104	<u>57.25</u>	6.518	0.3407
Real-ESRGAN	28.61	<u>0.8044</u>	<u>0.2848</u>	0.4525	54.26	6.701	<u>0.3422</u>
FeMaSR	26.87	0.7557	0.3179	<u>0.5534</u>	53.32	<b>5.775</b>	0.3121
DASR	<b>29.74</b>	<b>0.8257</b>	0.3143	0.3807	42.43	7.522	0.2822
SwinIR-GAN	28.46	0.8036	<b>0.2801</b>	0.4389	52.65	6.388	0.3265
DiffBIR+Ours	24.53	0.6014	0.5024	<b>0.7167</b>	<b>71.90</b>	<u>5.833</u>	<b>0.5902</b>
RealSR							
Methods	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
BSRGAN	<u>26.37</u>	0.7643	<u>0.2652</u>	0.5105	<u>63.19</u>	<u>5.690</u>	<u>0.3800</u>
Real-ESRGAN	25.65	0.7592	0.2720	0.4491	60.49	5.910	0.3769
FeMaSR	25.06	0.7342	0.2896	<u>0.5450</u>	59.20	5.807	0.3648
DASR	<b>27.01</b>	<u>0.7702</u>	0.3047	0.3135	40.95	6.682	0.2459
SwinIR-GAN	26.30	<b>0.7719</b>	<b>0.2479</b>	0.4367	58.83	5.800	0.3455
DiffBIR+Ours	23.97	0.6405	0.3667	<b>0.7090</b>	<b>72.83</b>	<b>5.068</b>	<b>0.5812</b>

Table 8: Quantitative results of BSRGAN (Zhang et al., 2021) method on the RealSR with FedSR.

Method	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
BSRGAN	26.37	0.7643	0.2652	0.5105	63.19	5.690	0.3800
BSRGAN+ours	25.33	0.7541	0.2711	0.5327	64.66	5.522	0.4141
$\Delta$ BSRGAN	-1.04	-0.0102	+0.0059	<b>+0.0222</b>	<b>+1.47</b>	<b>-0.168</b>	<b>+0.0341</b>

instance, in images containing abundant details, textures, and sharp edges, the amplitude of high-frequency components may be large, such as in a dense forest with lush leaves; likewise, the amplitude of high-frequency components increases in the presence of high-frequency noise. In contrast, low-frequency components may have relatively small amplitudes in patterns primarily composed of high-frequency information, such as fine lines or repetitive textures. Therefore, our research attempts to enhance both amplitude and high-frequency components. Our supplementary experiments in Appendix D.2 also indicate that enhancing low-amplitude channels in the backbone is not equivalent to enhancing high-frequency components. This further illustrates the differences between low-amplitude components and high-frequency details.

**The illustrations of reweighting of the AEM** Due to the obscurity of DNN’s features, our motivation for optimizing ISR primarily stems from the visible aspects of image space. However, insights derived from image space do not fully apply to the feature space due to the differences between the two. Although the effect of applying amplitude reweighting in AEM is not obvious in the image space (see Figure 10), this operation can effectively enhance the quality of ISR reconstruction in the feature space. This method of reducing amplitude values seems contrary to the conventional approach of increasing values to enhance results, and it is even somewhat counterintuitive. However, Row 2 and 3 in Table 9 indicate that increasing amplitude values actually leads to a deterioration in performance metrics. To explain this phenomenon, we first investigate the denoising process in diffusion models. In the previous paragraph, we detail the differences between low-amplitude

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

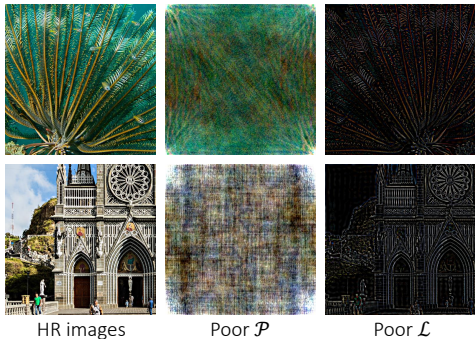


Figure 9: Visualization of degradation in phase and low-frequency components. Phase and low-frequency component degradation is denoted as poor  $\mathcal{P}$ , and poor  $\mathcal{L}$ .

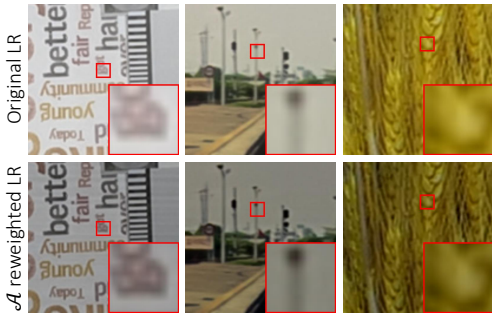


Figure 10: Visual comparison between original images and after amplitude  $\mathcal{A}$  reweighting on real-world LR images. It shows that the effect is not obvious in image space.

and high-frequency components. Since noise typically belongs to high-frequency components, we cannot establish a direct connection between optimizing denoising and reducing amplitude values. Regrettably, to our knowledge, there is currently no relevant literature evidence suggesting that modulating feature amplitudes to smaller values can improve image quality. Considering the internal structure of U-Net, we hypothesize that the attention modules within U-Net may prefer lower amplitude channel information. Therefore, further reducing the already information-rich low-amplitude channels may assist U-Net in better understanding and representing features, thereby enhancing super-resolution quality. Please note that the preference of deep networks for feature space exceeds the scope of this study, and we will explore this in greater depth in future work. Meanwhile, we encourage researchers in the community to provide more reasonable explanations.

## D.2 DISCUSSION ON THE MODULE CONFIGURATION

To demonstrate the effectiveness of our module configuration, we conduct supplementary experiments on AEM and HEM by replacing the positions of their effects. Specifically, we apply the AEM to the skip features and observe that its performance metrics are inferior to those obtained with the default settings for skip features (see Row 1, 2 of Table 9). Similarly, when applying HEM to the backbone features, the generated results were very poor (see Row 3, 5 of Table 9). This indicates that enhancing high-frequency components in the backbone does not equate to enhancing lower-amplitude components, further confirming the differences between high-frequency and low-amplitude components. Furthermore, we note that FreeU (Si et al., 2023) also includes frequency-related operations. To validate the effectiveness of our adaptive masking for high-frequency components, we replace HEM with the operations of skip features in FreeU, resulting in a significant decrease in no-reference metrics, specifically MUSIQ and MANIQA (see Row 5, 6 of Table 9).

Table 9: Supplementary experiments of the AEM and HEM on DRealSR and RealSR benchmarks.

Strategy		DRealSR/RealSR						
Module	Place	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	CLIP-IQA $\uparrow$	MUSIQ $\uparrow$	NIQE $\downarrow$	MANIQA $\uparrow$
AEM	Skip	27.32 / 24.09	0.7236 / 0.6805	0.3508 / 0.3275	0.6683 / 0.6554	62.51 / 67.88	6.275 / 5.642	0.4154 / 0.4504
AEM	Backbone (A $\uparrow$ )	28.12 / 24.80	0.7456 / 0.6886	0.3709 / 0.3404	0.4806 / 0.4703	45.91 / 55.26	6.561 / 6.284	0.2944 / 0.3218
AEM	Backbone (A $\downarrow$ )	27.36 / 24.27	0.7322 / 0.6890	0.3635 / 0.3389	<b>0.6760 / 0.6760</b>	<b>65.23 / 69.59</b>	<b>5.707 / 5.241</b>	<b>0.4681 / 0.5087</b>
HEM	Backbone	18.67 / 15.98	0.4093 / 0.2961	0.7266 / 0.7455	0.2011 / 0.1907	33.92 / 35.64	9.410 / 9.612	0.2840 / 0.2691
HEM	Skip (FreeU)	27.46 / 24.44	0.7241 / 0.6867	0.3520 / 0.3196	0.6732 / <b>0.6563</b>	61.10 / 66.29	<b>5.613 / 5.294</b>	0.3972 / 0.4245
HEM	Skip	26.91 / 23.85	0.7230 / 0.6859	0.3655 / 0.3351	<b>0.6749 / 0.6522</b>	<b>64.56 / 68.55</b>	5.768 / 5.386	<b>0.4447 / 0.4738</b>

## D.3 THE ALGORITHM AND PARAMTERS

As stated in the main paper, the AEM and the HEM are two key modules embedded in the skip connections of the U-Net within the Diffusion model. Each module contains its respective enhance-

**Algorithm 1** FedSR Algorithm

---

```

1026 for each  $t \in [1, \text{Sampling Steps}]$  do
1027
1028   2: Initialize the backbone features  $x_{bone}$  and the skip features  $x_{skip}$  in the skip connection;
1029   3:  $\mathbf{f}_{bone} = \mathcal{F}(x_{bone}), \mathbf{f}_{skip} = \mathcal{F}(x_{skip})$ 
1030   4: // (1) Amplitude Enhancement Module
1031   5:  $\mathcal{A}(x_{bone}), \mathcal{P}(x_{bone}) = \text{FFTSplit}(\mathbf{f}_{bone})$ ;
1032   6: // a) Channel Split;
1033   7: Split  $\mathcal{A}(x_{bone})$  by channel, then obtain  $\mathcal{S}_A = \{\mathcal{A}(x_{bone})_i\}_{i=1}^C$ ;
1034   8: // b) Average & Order;
1035   9: for each  $i \in [1, n]$  do
1036     10: Average amplitude value  $a_i = \frac{1}{HW} \sum_{h=1}^H \sum_{w=1}^W \mathcal{A}(x_{bone})_i^{(h,w)}$ ;
1037     11: end for
1038     12: Order  $(\mathcal{S}_A)$  by  $a_i$ 
1039     13: // c) Channel Selection;
1040     14:  $\mathcal{S}_S = \{\mathcal{A}(x_{bone})_i | a_i \leq a_{\min} + P_s \times (a_{\max} - a_{\min})\}$ ;
1041     15: // d) Amplitude Modulation;
1042     16: if  $\mathcal{A}(x_{bone})_i \in \mathcal{S}_S$  then
1043       17:  $\mathcal{M}_{bone} = 1 - P_a \cdot (\bar{\mathcal{A}} - \bar{\mathcal{A}}_{\min}) / (\bar{\mathcal{A}}_{\max} - \bar{\mathcal{A}}_{\min}), \mathcal{A}(x_{bone})'_i = \mathcal{A}(x_{bone})_i \odot \mathcal{M}_{bone}$ ;
1044     18: end if
1045     19:  $\mathbf{f}'_{bone} = \text{FFTCombine}(\mathcal{A}(x_{bone})', \mathcal{P}(x_{bone}))$ ,  $\mathbf{x}'_{bone} = \mathcal{G}(\mathbf{f}'_{bone})$ 
1046     20: // (2) High-frequency Enhancement Module
1047     21: for  $r \in [0, S/2]$  do
1048       22:  $\mathcal{M}_{skip}(r) = 1 + (r > r_{\text{thresh}}) \cdot [(S - S_{\min}) / (S_{\max} - S_{\min}) + 0.5] \cdot P_b / 2$ 
1049     23: end for
1050     24:  $\mathbf{f}'_{skip} = \mathbf{f}_{skip} \odot \mathcal{M}_{skip}$ ,  $\mathbf{x}'_{skip} = \mathcal{G}(\mathbf{f}'_{skip})$ ;
1051     25: end for

```

---

ment parameters, as detailed in Table 10. The detailed algorithmic process for these two modules is shown in Algorithm D.2.

Table 10: The parameters and their definitions for the AEM and HEM, which are set within five state-of-the-art diffusion-based ISR models.

Module	Parameter	Definition	StableSR	DiffBIR	SUPIR	SeeSR	PASD
AEM	$P_a$	The linearization param in Eq. (6)	0.3	0.3	0.05	0.3	0.3
	$P_s$	The selection threshold of Figure 6	0.3	0.3	0.3	0.3	0.3
HEM	$P_{b_1}$	The scaling factor in Eq. (8)	0.9	0.9	0.3	0.1	0.1
	$P_{b_2}$	The scaling factor in Eq. (8)	0.2	0.2	0.2	0.4	0.2

#### D.4 DISCUSSION ON THE METRICS

We show the DISTS metrics on the RealSR dataset (see Table 11). In the literature, the trade-off between fidelity and visual quality remains a long-standing challenge in the field of SR, and there is currently no definitive optimal evaluation metric. As noted by (Blau & Michaeli, 2018), this trade-off implies that solely optimizing distortion metrics may not only be ineffective but could also degrade visual quality. Meanwhile, we find recent Diffusion-based SR methods tends to emphasizing more on perceptual metrics such as MUSIQ and CLIP-IQA. (Wang et al., 2023c; Yu et al., 2024). Notably, our fluctuations on metrics like PNSR/SSIM are deems acceptable, much lower than the gap between SOTA diffusion-based methods themselves (e.g. SUPIR and StableSR differ by 0.1109 in SSIM, while DiffBIR and StableSR differ by 2.03 points in PSNR).

#### D.5 DISCUSSION ON THE TIMESTEP

Our FedSR is a highly flexible framework which can be adapted to specific timesteps. We conduct a preliminary experiment and observe that our FedSR demonstrates a greater impact during the early

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

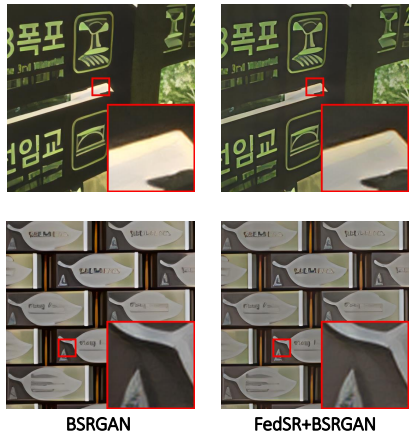


Figure 11: Visual comparisons of BSRGAN. After applying our method, BSRGAN (Zhang et al., 2021) presents the results with more natural image details and contrast.

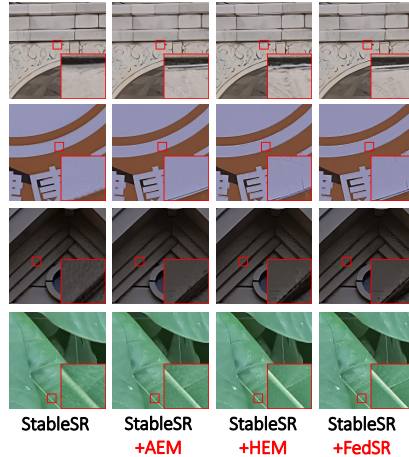


Figure 12: Visual effects of AEM and HEM. The results show that AEM primarily enhances the overall appearance, while HEM improves the clarity of details.

Table 11: Quantitative results of DISTS metrics on RealSR dataset.

Matrics	StableSR	DiffBIR	SeeSR	PASD	SUPIR
Baseline	0.2202	0.2401	0.2227	0.1989	0.2494
FedSR+	0.2422	0.2562	0.2294	0.2029	0.2632

denoising stages (see last row of Table 12). Additionally, incorporating FedSR in segments proves less effective than applying it as a whole.

Table 12: Quantitative results on specific timesteps on the RealSR dataset.

Matrics	PSNR	SSIM	LIPIPS	MUSIQ	CLIP-IQA	NIQE	MANIQA
StableSR	24.66	0.7003	0.3101	65.24	0.6169	5.924	0.4302
StableSR+FedSR(totally 1-200)	23.77	0.6832	0.3502	70.27	0.6683	5.094	0.5186
StableSR+FedSR (1-100)	24.04	0.6860	0.3302	68.62	0.6489	5.343	0.4681
StableSR+FedSR (101-200)	24.35	0.6966	0.3272	68.65	0.6628	5.511	0.5031

## D.6 COMPLEXITY ANALYSIS

In this section, we evaluate the complexity of the FedSR method using StableSR and DiffBIR as examples. We list the parameters and FLOPs of the denoising models in each framework below, which demonstrate almost the same statistics after integrating FedSR.

Table 13: Parameters and FLOPs of denoising models before and after integrating FedSR.

Matrics	StableSR	FedSR+StableSR	DiffBIR	FedSR+DiffBIR
Param (M)	918.93	918.93	1666.75	1666.75
FLOPs (G)	375.55	375.59	61.45	61.49

1134 E ADDITIONAL VISUAL RESULT  
1135

1136 In this section, we present additional experimental results. Figure 12 illustrates the visual effects  
1137 of AEM and HEM when applied individually and in combination. The results show that AEM  
1138 primarily enhances the overall image appearance, such as contrast, while HEM mainly improves the  
1139 clarity of high-frequency details.  
1140

1141 F LIMITATIONS AND FUTURE WORK  
1142

1143 Although our proposed FedSR achieves significant results, there are still some limitations. Similar  
1144 to other ISR studies on natural scenes, this work focuses only on existing natural image datasets and  
1145 synthetic datasets for ISR tasks. Applying ISR on a larger scale to AI-generated datasets remains  
1146 an interesting avenue for further exploration. Additionally, we only employ a training-free imple-  
1147 mentation, without delving into model training and fine-tuning. In future work, we will explore  
1148 how to leverage the network’s preference for frequency domain components to fine-tune the model  
1149 architecture, thereby further enhancing ISR quality.  
1150

1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187