

# DIFF-SSR: DIFFUSION MODEL WITH STRUCTURE-MODULATED FOR IMAGE SUPER-RESOLUTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Diffusion-based super-resolution (SR) models have recently garnered significant attention due to their potent restoration capabilities. But conventional diffusion models perform noise sampling from a single distribution, constraining their ability to handle real-world scenes and complex textures across semantic regions. With the success of segment anything model (SAM), generating sufficiently fine-grained region masks can enhance the detail recovery of diffusion-based SR model. However, directly integrating SAM into SR models will result in much higher computational cost. In this paper, we propose the Diff-SSR model, which can utilize the fine-grained structure information from SAM in the process of sampling noise to improve the image quality without additional computational cost during inference. In the process of training, we encode structural position information into the segmentation mask from SAM. Then the encoded mask is integrated into the forward diffusion process by modulating it to the sampled noise. This adjustment allows us to independently adapt the noise mean within each corresponding segmentation area. The diffusion model is trained to estimate this modulated noise. Crucially, our proposed framework does NOT change the reverse diffusion process and does NOT require SAM at inference. Experimental results demonstrate the effectiveness of our proposed method, showcasing superior performance in suppressing artifacts.

## 1 INTRODUCTION

Single-image super-resolution (SR) has remained a longstanding research focus in computer vision, aiming to restore a high-resolution (HR) image based on a low-resolution (LR) reference image. The applications of SR span various domains, including mobile phone photography (Ignatov et al., 2022), medical imaging (Huang et al., 2017; Isaac & Kulkarni, 2015), and remote sensing (Wang et al., 2022a; Haut et al., 2018). Considering the inherently ill-posed nature of the SR problem, deep learning models (Dong et al., 2014; Kim et al., 2016; Chen et al., 2021) have been employed. These models leverage deep neural networks to learn informative hierarchical representations, allowing them to effectively approximate HR images.

Conventional deep learning-based SR models typically process an LR image progressively through CNN blocks (Zhang et al., 2018a) or transformer blocks (Liang et al., 2021; Chen et al., 2021; 2023). The final output is then compared with the corresponding HR image using distance measurement (Dong et al., 2014; Zhang et al., 2018a) or adversarial loss (Ledig et al., 2017; Wang et al., 2018b). Despite the significant progress achieved by these methods, there remains a challenge in generating satisfactory textures (Li et al., 2023). The introduction of diffusion models (Ho et al., 2020a; Rombach et al., 2022a) marked a new paradigm for image generation, exhibiting remarkable performance. Motivated by this success, several methods have incorporated diffusion models into the image SR task (Saharia et al., 2022b; Li et al., 2022; Shang et al., 2023; Xia et al., 2023). Saharia *et al.* (Saharia et al., 2022b) introduced diffusion models to predict residuals, enhancing convergence speed. Building upon this framework, Li *et al.* (Li et al., 2022) further integrated a frequency domain-based loss function to improve the prediction of high-frequency details.

In comparison with traditional CNN-based methods, diffusion-based image SR has shown significant performance improvements in texture-level prediction. However, existing approaches in this domain often employ independent and identically distributed noise during the diffusing process, ig-

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

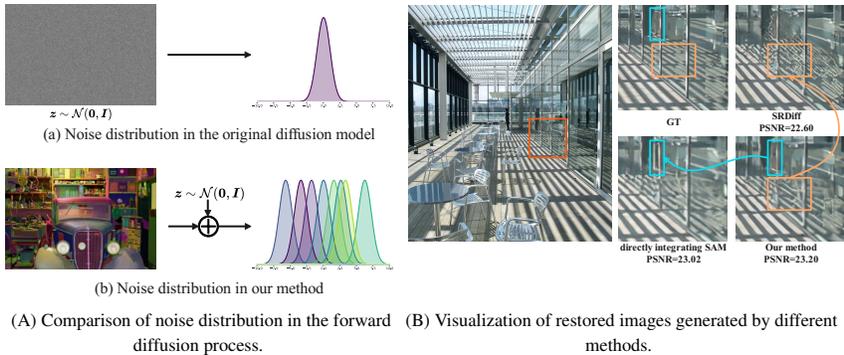


Figure 1: (A) is comparison of noise distribution in the forward diffusion process between existing diffusion-based image SR methods and our Diff-SSR. Our approach enhances the restoration of different image areas by modulating the corresponding noise with guidance from segmentation masks generated by SAM. (B) is Visualization of restored images generated by different methods. Our method can achieve similar reconstruction performance to directly integrating SAM into diffusion model.

noring the fact that different local areas of an image may exhibit distinct data distributions. This oversight can lead to inferior structure-level restoration and chaotic texture distribution in generated images due to confusion of information across different regions. In the visualization of SR images, this manifests as distorted structures and bothersome artifacts.

Recently, the segment anything model (SAM) has emerged as a novel approach capable of extracting exceptionally detailed segmentation masks from given images (Kirillov et al., 2023). For instance, SAM can discern between a feather and beak of a bird in a photograph, assigning them to distinct areas in the mask, which provides a sufficiently fine-grained representation of the original image at the structural level. This structure-level ability is exactly what diffusion model lacks. But directly integrating SAM into diffusion model may result in significant computational costs at inference stage. Motivated by these problems, we are intrigued by the question: *Can we introduce structure-level ability to distinguish different regions in the diffusion model, ensuring the generation of correct texture distribution and structure in each region without incurring additional inference time?*

In this paper, we verified the feasibility of controlling the generated images by modulating the distribution of noise during training stage, and the theory is illustrated in Figure 1(A). Based on this theory, we proposed the structure-modulated diffusion framework named Diff-SSR for image SR task. This framework utilizes the fine-grained structure segmentation ability to guide image restoration. By enabling the denoise model (U-Net) to approximate the SAM ability, it can modulate the structure information into the noise during the diffusion process.

The training and inference process are illustrated in Figure 3(b). Our method does not change the inference process, and the training process is as follows: (1)For each HR image in the training set, SAM is employed to generate a fine-grained segmentation masks. (2) Subsequently, the Structural Position Encoding (SPE) module is introduced to incorporate masks by position information and generate SPE mask. (3) Finally, the SPE mask is utilized to modulate the mean of the diffusing noise in each fine-grained area separately, thereby enhancing accuracy of structure and texture distribution during the forward diffusion process.

To achieve the goal of reducing the cost of training and inference, our method have with the following advantages:

- During the training, our method *have negligible extra training cost*. We use SAM to pre-generated mask of training samples, and reused them in all epochs. And the cost of modulate noise process is negligible.
- During the inference, our method *have no additional inference cost*. The diffusion model has already acquired structure-level knowledge during training, it can restore SR images without requiring access to the oracle SAM.

We conduct extensive experiments on several commonly used image SR benchmarks, and our method showcases superior performance over existing diffusion-based methods. Furthermore, our

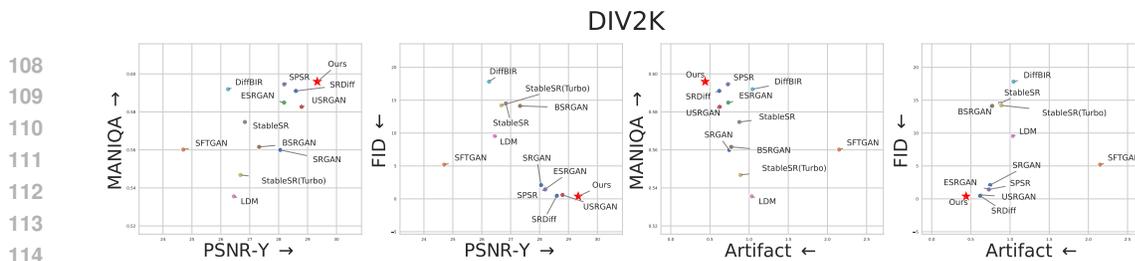


Figure 2: We compared the metrics MANIQA, FID, PSNR, and Artifact(5.3) on the DIV2K dataset. In this context, higher values of MANIQA and PSNR are better, while lower values of FID and Artifact are preferred. The red arrow indicates the direction of the best performance based on the combined horizontal and vertical metrics.

method has the fewest artifacts in generated models such as GAN and diffusion models. Our model achieved a balanced advantage across various metric combinations, as shown in Figure 2.

## 2 RELATED WORKS

### 2.1 DISTANCE-BASED SUPER-RESOLUTION

Neural network-based methods have become the dominant approach in image super-resolution (SR). The introduction of convolutional neural networks (CNN) to the image SR task, as exemplified by SRCNN (Dong et al., 2015), marked a significant breakthrough, showcasing superior performance over conventional methods. Subsequently, numerous CNN-based networks has been proposed to further enhance the reconstruction quality. This is achieved through the design of new residual blocks (Ledig et al., 2017) and dense blocks (Wang et al., 2018b; Zhang et al., 2018b). Moreover, the incorporation of attention mechanisms in several studies (Dai et al., 2019; Mei et al., 2021) has led to notable performance improvements.

Recently, the Transformer architecture (Vaswani et al., 2017) has achieved significant success in the computer vision field. Leveraging its impressive performance, Transformer has been introduced for low-level vision tasks (Tu et al., 2022; Wang et al., 2022b; Zamir et al., 2022). In particular, IPT (Chen et al., 2021) develops a Vision Transformer (ViT)-style network and introduces multi-task pre-training for image processing. SwinIR (Liang et al., 2021) proposes an image restoration Transformer based on the architecture introduced in (Liu et al., 2021). VRT (Liang et al., 2022b) introduces Transformer-based networks to video restoration. EDT (Li et al., 2021) validates the effectiveness of the self-attention mechanism and a multi-related-task pre-training strategy. These Transformer-based approaches consistently push the boundaries of the image SR task.

### 2.2 GENERATIVE SUPER-RESOLUTION

To enhance the perceptual quality of SR results, Generative Adversarial Network (GAN)-based methods have been proposed, introducing adversarial learning to the SR task. SRGAN (Ledig et al., 2017) introduces an SRResNet generator and employs perceptual loss (Johnson et al., 2016) to train the network. ESRGAN (Wang et al., 2018b) further enhances visual quality by adopting a residual-in-residual dense block as the backbone for generator.

In recent times, diffusion models (Ho et al., 2020a) have emerged as influential in the field of image SR. SR3 (Saharia et al., 2022b) and SRdiff (Li et al., 2022) have successfully integrated diffusion models into image SR, surpassing the performance of GAN-based methods. Additionally, Palette (Saharia et al., 2022a) draws inspiration from conditional generation models (Mirza & Osindero, 2014) and introduces a conditional diffusion model for image restoration. Despite their success, generated models often suffer from severe perceptually unpleasant artifacts. SPSR (Ma et al., 2020) addresses the issue of structural distortion by introducing a gradient guidance branch. LDL (Liang et al., 2022a) models the probability of each pixel being an artifact and introduces an additional loss during training to inhibit artifacts.

## 2.3 SEMANTIC GUIDED SUPER-RESOLUTION

As image SR is a low-level vision task with a pixel-level optimization objective, SR models inherently lack the ability to distinguish between different semantic structures. To address this limitation, some works introduce segmentation masks generated by semantic segmentation models as conditional inputs for generated models. For instance, (Gatys et al., 2017) utilizes semantic maps to control perceptual factors in neural style transfer, while (Ren et al., 2017) employs semantic segmentation for video deblurring. SFTGAN (Wang et al., 2018a) demonstrates the possibility of recovering textures faithful to semantic classes. SSG-RWSR (Aakerberg et al., 2022) utilizes an auxiliary semantic segmentation network to guide the super-resolution learning process.

Image segmentation tasks have undergone significant evolution in recent years, wherein the most recent development is the SAM (Kirillov et al., 2023), showcasing superior improvements in segmentation capability and granularity. The powerful segmentation ability of SAM has opened up new ideas and tools for addressing challenges in various domains. For instance, (Xiao et al., 2023) leverages semantic priors generated by SAM to enhance the performance of image restoration models. Similarly, (Lu et al., 2023) improves both alignment and fusion procedures by incorporating semantic information from SAM. However, these approaches necessitate segmentation models to provide semantic information during inference, resulting in much higher latency. In contrast, our method endows SR models with the ability to distinguish different semantic distributions in images without incurring additional costs at inference.

## 3 PRELIMINARY

### 3.1 DIFFUSION MODEL

The diffusion model is an emerging generative model that has demonstrated competitive performance in various computer vision fields (Ho et al., 2020a; Rombach et al., 2022a). The basic idea of diffusion model is to learn the reverse of a forward diffusion process. Sampling in the original distribution can then be achieved by putting a data point from a simpler distribution through the reverse diffusion process. Typically, the forward diffusion process is realized by adding standard Gaussian noise to a data sample  $\mathbf{x}_0 \in \mathbb{R}^{c \times h \times w}$  from the original data distribution step by step:

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

where  $\mathbf{x}_t$  represents the latent variable at diffusion step  $t$ . The hyperparameters  $\beta_1, \dots, \beta_T \in (0, 1)$  determine the scale of added noise for  $T$  steps. With a proper configuration of  $\beta_t$  and a sufficiently large number of diffusing steps  $T$ , a data sample from the original distribution transforms into a noise variable following the standard Gaussian distribution. During training, a model is trained to learn the reverse diffusion process, *i.e.*, predicting  $\mathbf{x}_{t-1}$  given  $\mathbf{x}_t$ . At inference time, new samples are generated by using the trained model to transform a data point sampled from the Gaussian distribution back into the original distribution.

As illustrated in Equation 1, identical Gaussian noise is added to each pixel of the sample during the forward diffusion process, indicating that all spatial positions are treated equally. Existing approaches (Saharia et al., 2022b; Li et al., 2022; Shang et al., 2023; Xia et al., 2023) introduce the diffusion model into the image SR task following this default setting of noise. However, image SR is a low-level vision task aiming at learning a mapping from the LR space to the HR space. This implies that data distributions in corresponding areas of an LR image and an HR image are highly correlated, while other areas are nearly independent of each other. The adoption of identical noise in diffusion-based SR overlooks this local correlation property and may result in an inferior restoration of structural details due to the confusion of information across different areas in an image. Therefore, injecting spatial priors into diffusion models to help them learn local projections is a promising approach to improve diffusion-based image SR.

### 3.2 SEGMENT ANYTHING MODEL

Segment Anything Model (SAM) is proposed as a foundational model for segmentation tasks, comprising a prompt encoder, an image encoder, and a lightweight mask decoder. The mask decoder generates a segmentation mask by incorporating both the encoded prompt and image as input.

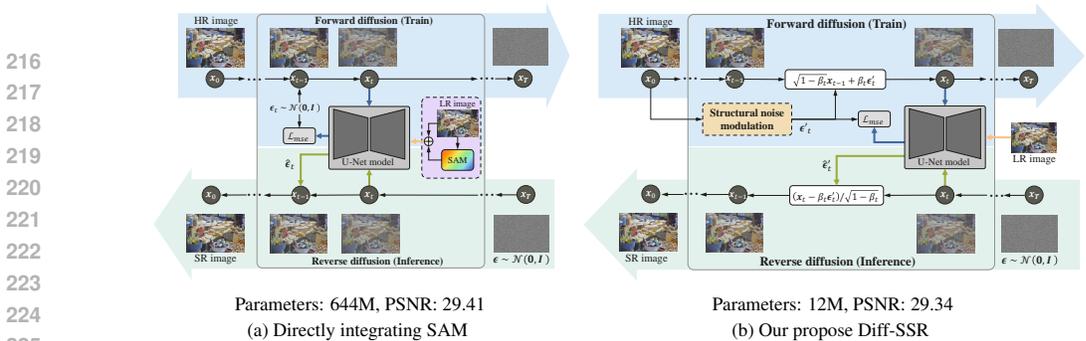


Figure 3: Comparison between (a) directly integrating SAM into the diffusion model and (b) our proposed Diff-SSR reveals distinct approaches, and the PSNR evaluate on DIV2K dataset. In (a), mask information predicted by SAM is utilized during both the training and inference stages. In contrast, (b) only employs modulated noise generated by the structural noise modulation model during training. The details of structural noise modulation can be found in Figure 4(a), and our method achieves comparable reconstruction performance to (b) as demonstrated in Figure 1(B).

In comparison to conventional cluster-based models and image segmentation models, SAM is preferable for generating segmentation masks in image SR tasks. Cluster-based models lack the ability to extract high-level information from images, resulting in the generation of low-quality masks. Deep-learning image segmentation models, while capable of differentiating between different objects, produce coarse masks that struggle to segment areas within an object. In contrast, SAM excels in generating extraordinarily fine-grained segmentation masks for given images, owing to its advanced model architecture and high-quality training data. It can generate mask for each different texture region. This ability to distinguish different texture distribution is we aspire to incorporate into diffusion model.

Table 1: Comparison of the effectiveness and efficiency of various diffusion-based image super-resolution methods.

	SRDiff	SAM+SRDiff	Diff-SSR
Parameter	12M	632M+12M	12M
Train time	10h16min/100k step	48h52min/100k step	10h21min/100k step
Inference time	37.64s/per img	65.72s/per img	37.62s/per img
PSNR	28.6	29.41	29.34
FID	0.4649	0.3938	0.3809

### 3.3 DIRECTLY INTEGRATING SAM INTO DIFFUSION MODEL

To validate the enhancing effect of structure level information on the diffusion process, we devised a straightforward diffusion model (SAM+SRDiff) to utilize the mask information predicted by SAM. Specifically, we concatenated the LR image with the embedding mask information to guide the denoising model in predicting noise. The model structure is detailed in Figure 3(a). Results indicate that the images generated by this simple model exhibit more accurate texture and fewer artifacts.

However, this approach introduces additional inference time as SAM predicts the mask, as shown in Table 1. Can we enable the diffusion model to learn the capability of distinguishing different texture distributions without relying on an auxiliary model? Furthermore, is it possible to train the denoising model to acquire this capability?

## 4 METHOD

### 4.1 OVERVIEW

In this paper, we present Diff-SSR, a structure-modulated diffusion framework designed to improve the performance of diffusion-based image SR models by leveraging fine-grained segmentation masks. As illustrated in Figure 3(b), these masks play a crucial role in a structural noise modulation

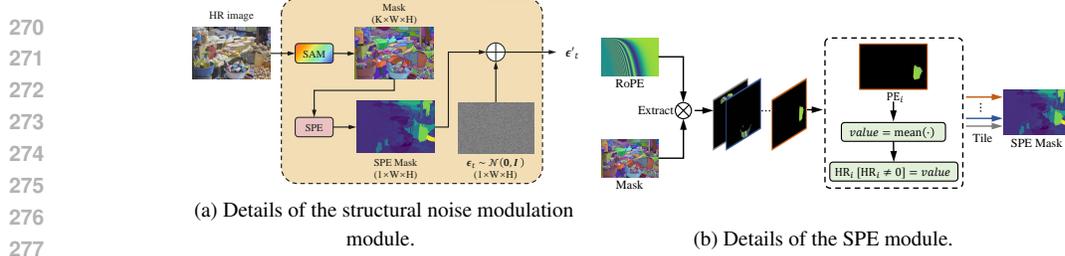


Figure 4: (a) During training, a SAM generates a segmentation mask for an HR image, and a structural position encoding (SPE) module encodes structure-level position information in the mask. The encoded mask is then added to the noise to modulate its mean in each segmentation area separately. At inference time, the framework utilizes only the trained diffusion model for image restoration, eliminating the inference cost of SAM. (b) This module encodes structural position information in the mask generated by SAM.

module, modulating the mean of added noise in different segmentation areas during the forward process. Additionally, a structural position encoding (SPE) module is integrated to enrich the masks with structure-level position information.

We elaborate on the forward process in the proposed framework.\* As discussed in Section 3.1, the added noise at each spatial point is independent and follows the same distribution, treating different areas in sample  $\mathbf{x}_0$  equally during the forward process, even though they may possess different structural information and distributions. To address this limitation, we utilize a SAM to generate segmentation masks for modulating the added noise. The corresponding segmentation mask of  $\mathbf{x}_0$  generated by SAM is denoted as  $\mathbf{M}_{\text{SAM}}$ . We then encode structural information into the mask using the SPE module, and the resulting encoded embedding mask is denoted as  $\mathbf{E}_{\text{SAM}}$ . Details of the SPE module will be provided in Section 4.2. At each step of the forward process,  $\mathbf{E}_{\text{SAM}}$  is added to the standard Gaussian noise to inject structure-level information into the diffusion model. This modified process can be formulated as:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \mathbf{E}_{\text{SAM}}, \beta_t \mathbf{I}). \quad (2)$$

Compared with the original forward diffusion process defined in Equation 1, the modified process adds noise with different means to different segmentation areas. This makes local areas in an image distinguishable during forward diffusion, further aiding the diffusion model in learning a reverse process that makes more use of local information when generating an SR restoration for each area. Since the added Gaussian noise is independently sampled at each step, we can obtain the conditional distribution of  $\mathbf{x}_t$  given  $\mathbf{x}_0$  by iteratively applying the modified forward process:

$$q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \varphi_t \mathbf{E}_{\text{SAM}}, (1 - \bar{\alpha}_t) \mathbf{I}), \quad (3)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ , and  $\varphi_t = \sum_{i=1}^t \sqrt{\frac{\alpha_t \beta_i}{\alpha_i}}$ . With this formula, we can directly derive the latent variable  $\mathbf{x}_t$  from  $\mathbf{x}_0$  in one step.

To achieve the SR image from restoration of an LR image, learning the reverse of the forward diffusion process is essential, characterized by the posterior distribution  $p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{E}_{\text{SAM}})$ . However, the intractability arises due to the known marginal distributions  $p(\mathbf{x}_{t-1})$  and  $p(\mathbf{x}_t)$ . This challenge is addressed by incorporating  $\mathbf{x}_0$  into the condition. Employing Bayes' theorem, the posterior distribution  $p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}})$  can be formulated as:

$$\begin{aligned} \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) &= \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \left( \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon \right) \right), \\ \tilde{\beta}_t &= \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t, \\ p(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) &= \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}), \tilde{\beta}_t \mathbf{I}), \end{aligned} \quad (4)$$

\*For additional details regarding the derivation, please refer to the supplementary material.

where  $\epsilon \sim \mathcal{N}(0, 1)$ . To generate an SR image of an unseen LR image, we need to estimate the weighted summation of  $\mathbf{E}_{\text{SAM}}$  and  $\epsilon$ , as these variables are only defined in the forward process and cannot be accessed during inference. We adopt a denoising network  $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{x}_{LR}, t)$  for approximation. The associated loss function is formulated as:

$$\mathcal{L}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[ \left\| \frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon - \epsilon_{\theta}(\mathbf{x}_t, \mathbf{x}_{LR}, t) \right\|_2^2 \right]. \quad (5)$$

The denoising network  $\epsilon_{\theta}(\mathbf{x}_t, \mathbf{x}_{LR}, t)$  predicts the weighted summation based on latent variable  $\mathbf{x}_t$ , LR image  $\mathbf{x}_{LR}$ , and step  $t$ . During training,  $\mathbf{x}_t$  is derived by sampling from the distribution defined in Equation 3. At inference time, the restored sample at step  $t$  is used as  $\mathbf{x}_t$ .

**Discussion.** The structure-level information encoded by the mask can be injected into the diffusion model through two distinct approaches. One method involves using the mask to modulate the input of the diffusion model, while the other method entails modulating the noise in the forward process, which is the approach adopted in our proposed method. In comparison to directly modulating the input, our method only requires the oracle SAM during training. Subsequently, the trained diffusion model can independently restore the SR image of an unseen LR image by iteratively applying the posterior distribution defined in Equation 4. This highlights that our Diff-SSR method incurs *no additional inference cost* during inference.

## 4.2 STRUCTURAL POSITION ENCODING

After obtaining the original segmentation mask using SAM, we employ an SPE module to encode structural position information in the mask. Details of this module are illustrated in Figure 4(b).

The fundamental concept behind the SPE module is to assign a unique value to each segmentation area. The segmentation mask generated by SAM comprises a series of 0-1 masks, where each mask corresponds to an area in the original image sharing the same semantic information. Consequently, for HR image  $\mathbf{x}_{HR}^{3 \times h \times w}$ , we can represent the  $K$  segmentation masks as  $\mathbf{M}_{\text{SAM}, i}$ , where  $i = 1, 2, \dots, K$  is the index of different areas in the original image. Specifically, the value of a point in  $\mathbf{M}_{\text{SAM}, i} \in \{0, 1\}^{h \times w}$  equals 1 when its position is within the  $i$ -th area in the original image and 0 otherwise. To encode position information, we generate a rotary position embedding (RoPE) (Su et al., 2021)  $\mathbf{x}_{\text{RoPE}} \in \mathbb{R}^{1 \times h \times w}$ , where the width is considered the length of the sequence and the height is considered the embedding dimension in RoPE. We initialize  $\mathbf{x}_{\text{RoPE}}$  with a 1-filled tensor. Similarly,  $\mathbf{x}_{\text{RoPE}}$  can be decomposed as:  $\mathbf{x}_{\text{RoPE}} = \sum_i \mathbf{x}_{\text{RoPE}, i} = \sum_i \mathbf{x}_{\text{RoPE}} \cdot \mathbf{M}_{\text{SAM}, i}$ . Subsequently, we obtain the structurally positioned embedded mask by:

$$\mathbf{E}_{\text{SAM}} = \sum_i \mathbf{M}_{\text{SAM}, i} \cdot \text{mean}(\mathbf{x}_{\text{RoPE}, i}), \quad (6)$$

which means to assign the average value of  $\mathbf{x}_{\text{RoPE}, i}$  to  $i$ -th segmentation area.

## 4.3 TRAINING AND INFERENCE

The training of the diffusion model necessitates segmentation masks for all HR images in the training set. We employ SAM to generate these masks. This process is executed once before training, and the generated masks are reused in all epochs. Therefore, our method incurs only a negligible additional training cost from the integration of SAM. Subsequently, a model is trained to estimate the modulated noise in the forward diffusion process using the loss function outlined in Equation 5.

During inference, we follow the practice of SRDiff(Li et al., 2022), the restoration of SR images can be accomplished by applying the reverse diffusion process to LR images. By iteratively applying the posterior distribution in Equation 4 and utilizing the trained model to estimate the mean, the restoration of the corresponding SR image is achieved. It is noteworthy that we opted for the  $\mathbf{x}_T$  sample from  $\mathcal{N}(0, \mathbf{I})$  instead of  $\mathcal{N}(\varphi_T \mathbf{E}_{\text{SAM}}, \mathbf{I})$ . Because the denoising model can generate the correct noise distribution, the initial distribution is not expected to exert a significant impact on the ultimately reconstructed image during the iterative denoising process. Simultaneously, such choice also ensures that our Diff-SSR method without additional inference cost.

Table 2: Results on test sets of several public benchmarks and the validation set of DIV2K. We report the results achieved by GAN-based and diffusion-based methods. ( $\uparrow$ ) and ( $\downarrow$ ) indicate that a larger or smaller corresponding score is better, respectively. Best and second best performance are marked with **best** and **second**.

Method	Urban100					BSDS100					Manga109					General100					DIV2K				
	PSNR $\uparrow$	SSIM $\uparrow$	MANIQAT $\uparrow$	FID $\downarrow$	Artifact $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	MANIQAT $\uparrow$	FID $\downarrow$	Artifact $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	MANIQAT $\uparrow$	FID $\downarrow$	Artifact $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	MANIQAT $\uparrow$	FID $\downarrow$	Artifact $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	MANIQAT $\uparrow$	FID $\downarrow$	Artifact $\downarrow$
SRGAN	22.86	0.6846	0.6162	10.50	2.7320	24.76	0.6400	0.6058	54.50	1.2158	28.08	0.8616	0.5822	4.18	0.4736	25.99	0.7470	0.6172	36.23	1.4216	28.05	0.7738	0.5600	2.0889	0.7456
SFTGAN	21.95	0.6457	0.6153	9.14	4.4362	24.69	0.6365	0.6173	49.29	1.2137	20.72	0.7008	0.5687	9.65	5.7064	22.20	0.6432	0.6253	37.06	3.6220	24.70	0.6929	0.5602	5.1979	2.1495
ESRGAN	22.99	0.6940	0.6678	7.38	2.7006	24.66	0.6374	0.6449	<b>45.88</b>	1.2331	28.60	0.8553	0.6026	3.12	<b>0.4042</b>	26.04	0.7449	0.6452	<b>30.93</b>	1.4331	28.19	0.7709	0.5849	1.4586	0.7335
USRGAN	25.24	0.7060	0.6785	6.44	2.5297	25.13	0.6604	0.6517	<b>48.58</b>	1.0947	20.70	0.7092	0.6226	8.61	5.7367	26.35	0.7631	0.6411	35.22	1.3029	<b>28.70</b>	<b>0.7945</b>	0.5827	0.5938	0.6239
SPSR	23.05	0.6973	<b>0.6823</b>	7.84	2.7069	24.61	0.6375	<b>0.6648</b>	48.81	1.2467	23.26	0.7740	0.6211	6.64	2.6665	25.96	0.7435	0.6571	<b>30.94</b>	1.4701	28.19	0.7727	<b>0.5945</b>	1.4315	0.7295
BSRGAN	22.37	0.6628	0.6334	33.74	2.9030	24.95	0.6365	0.5993	114.08	1.1467	26.10	0.8272	0.6105	33.51	1.0150	25.23	0.7309	0.6337	86.14	1.5147	27.33	0.7577	0.5616	14.1312	0.7718
LDM	22.23	0.6546	0.6239	23.07	3.4932	23.56	0.5812	0.6194	109.77	1.8173	24.26	0.7941	0.5870	20.75	1.9994	25.32	0.6779	0.5683	265.82	1.5201	26.45	0.7340	0.5356	9.5518	1.0334
StableSR	21.16	0.6529	<b>0.7025</b>	28.94	4.0859	24.64	0.6523	0.6606	68.77	1.2014	21.22	0.7456	<b>0.6576</b>	31.41	4.5033	18.39	0.5324	<b>0.6749</b>	73.51	8.4946	26.83	0.7653	0.5747	14.5232	0.8749
StableSR(Turbo)	21.22	0.6658	0.6633	30.61	3.9131	24.61	0.6691	0.6347	74.04	1.2598	22.68	0.7819	0.5875	29.16	2.9956	18.63	0.5421	0.6446	68.68	8.3225	26.68	0.7776	0.5468	14.2137	0.8883
DiffBIR	22.40	0.6417	0.6336	30.64	2.9967	25.09	0.6254	0.6626	69.18	1.0446	21.81	0.7197	<b>0.6281</b>	30.64	4.0051	24.37	0.6678	<b>0.6762</b>	66.35	1.8129	26.25	0.7051	0.5919	17.8206	1.0440
PSA-SR	24.71	0.6444	0.4483	61.86	2.7320	13.72	0.4802	0.3449	58.53	1.0160	24.36	0.7115	0.5492	59.39	4.0033	25.24	0.7094	0.4556	23.56	1.7406	-	-	-	-	-
StructSR	16.50	0.4153	0.4231	138.71	15.2055	12.15	0.4452	0.3813	87.62	6.9455	-	-	-	-	-	14.45	0.4654	0.3146	100.26	15.9960	13.91	0.3572	0.4408	83.8778	7.2400
SRDiff	<b>25.98</b>	<b>0.7692</b>	0.6604	<b>5.22</b>	<b>1.4163</b>	<b>25.86</b>	<b>0.6895</b>	0.6478	56.27	1.2226	<b>28.78</b>	<b>0.8761</b>	0.5967	<b>2.49</b>	0.4047	<b>29.82</b>	<b>0.8223</b>	0.6590	36.35	<b>0.5370</b>	28.60	0.7908	0.5910	<b>0.4649</b>	<b>0.6185</b>
Diff-SSR (Ours)	<b>25.54</b>	<b>0.7721</b>	0.6709	<b>4.53</b>	<b>1.1453</b>	<b>26.47</b>	<b>0.7003</b>	<b>0.6667</b>	60.81	<b>0.9236</b>	<b>29.43</b>	<b>0.8899</b>	0.6046	<b>2.40</b>	<b>0.3081</b>	<b>30.29</b>	<b>0.8353</b>	0.6346	39.15	<b>0.3145</b>	<b>29.34</b>	<b>0.8108</b>	<b>0.5959</b>	<b>0.3809</b>	<b>0.4391</b>

## 5 EXPERIMENT

### 5.1 EXPERIMENTAL SETUP

**Dataset.** We evaluate the proposed method on the general SR ( $4\times$ ) task. The training data in DIV2K (Agustsson & Timofte, 2017) and all data in Flickr2K (Wang et al., 2019) are adopted as the training set. For images in the training set, we adopt a SAM to obtain their corresponding segmentation masks. After that, structural position information is encoded into the mask by the SPE module in our proposed framework. We adopt a patch size settings of  $160\times 160$  to crop each image and its corresponding mask. For evaluation, several commonly-used SR testing dataset are used, including Set14 (Zeyde et al., 2012), Urban100 (Huang et al., 2015), BSDS100 (Martin et al., 2001), Manga109 (Fujimoto et al., 2016), General100 (Dong et al., 2016). Besides, the validation set of DIV2K (Agustsson & Timofte, 2017) is also used for evaluation.

**Baseline.** We choose a wide range of methods for comparison. Among them, SRGAN (Ledig et al., 2017), SFTGAN (Wang et al., 2018a), ESRGAN (Wang et al., 2018b), BSRGAN (Zhang et al., 2021), USRGAN (Zhang et al., 2020), and SPSR (Ma et al., 2020) are GAN-based generative methods. Besides, we also comparison with diffusion-base generative methods such as LDM (Rombach et al., 2022b), StableSR (Wang et al., 2023), StructSR (Li et al., 2025), PiSA-SR (Sun et al., 2025), DiffBIR (Lin et al., 2023), and SRDiff (Li et al., 2022), .

**Model architecture.** Architecture of the used denoising model in our experiments follows Li *et al.* (Li et al., 2022). As for configuration of the forward diffusion process, we set the number of diffusing steps  $T$  to 100 and scheduling hyperparameters  $\beta_1, \dots, \beta_T$  following Nichol *et al.* (Nichol & Dhariwal, 2021)

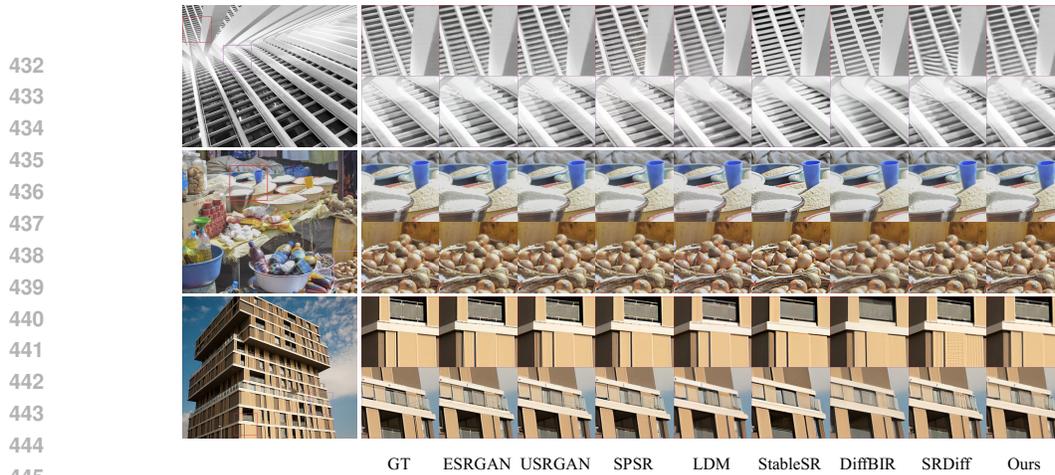
**Optimization.** We train the diffusion model for 400K iterations with a batch size of 16, and adopt Adam (Kingma & Ba, 2014) as the optimizer. The initial learning rate is  $2 \times 10^{-4}$  and the cosine learning rate decay is adopted. The training process requires approximately 75 hours and 30GB of GPU memory on a single GPU card.

**Metric.** Both objective and subjective metrics are used in our experiment. PSNR and SSIM (Wang et al., 2004) serve as objective metrics for quantitative measurements, which are computed over the Y-channel after converting SR images from the RGB space to the YUV space. To evaluate the perceptual quality, we also adopt Fréchet inception distance (FID) (Heusel et al., 2017) and MANIQA (Yang et al., 2022) as the subjective metric, which measures the fidelity and diversity of generated images.

### 5.2 PERFORMANCE OF IMAGE SR

We compare the performance of the proposed Diff-SSR method with baselines on several commonly used benchmarks for image SR. The quantitative results are presented in Table 2. In the results, our method outperforms the diffusion-based baseline SRDiff in terms of all three metrics, except a slightly higher FID score on BSDS100 and General100. Moreover, Diff-SSR can even achieves better performance when compared to conventional approaches.

Figure 5 presents several images by generated different methods. Compared with the baselines, our methods is able to generate more realistic details of the given image. Moreover, the reconstructed

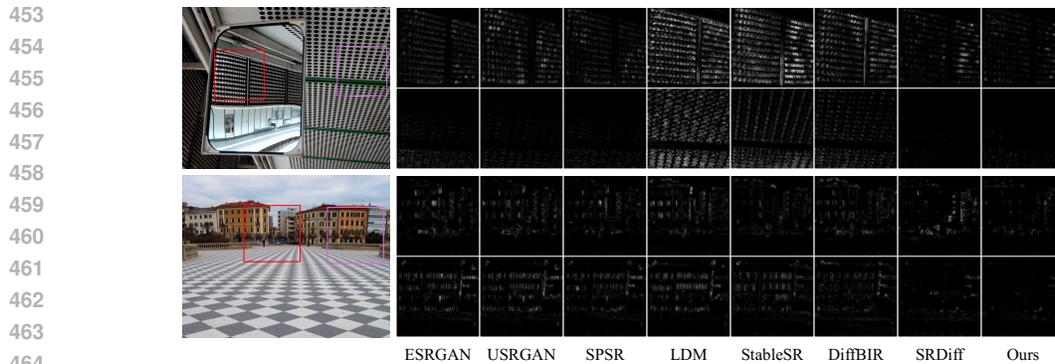


445  
446  
447  
448  
449

Figure 5: Visualization of restored images generated by different methods. Our Diff-SSR surpasses other approaches in terms of both higher reconstruction quality and fewer artifacts. Additional visualization results can be found in our supplementary material.

450  
451  
452

images contain less artifact, which refers to the unintended distortion or anomalies in the SR image. We further evaluate the proposed method in terms of inhibiting artifact in Section 5.3.



465  
466  
467

Figure 6: Visualization of artifact maps. Bright regions indicate artifacts in the restored images. Our proposed method generates images with fewer artifacts compared to other methods.

### 468 5.3 PERFORMANCE OF INHIBITING ARTIFACT

469  
470  
471  
472  
473

Generative image SR models excel at recovering sharp images with rich details. However, they are prone to unintended distortions or anomalies in the restored images (Liang et al., 2022a), commonly referred to as artifacts. In our experiments, we closely examine the performance of our method in inhibiting artifacts.

474  
475  
476  
477  
478

Following the approach outlined in (Liang et al., 2022a), we calculate the artifact map for each SR image. Table 2 presents the averaged values of artifact maps on four datasets, and Figure 6 visually showcases the artifact maps. When compared with other methods, our Diff-SSR demonstrates the ability to generate SR images with fewer artifacts, as supported by both quantitative and qualitative assessments.

## 480 6 CONCLUSION

481  
482  
483  
484  
485

This paper introduces Diff-SSR, a framework that incorporates structural position information from SAM-generated masks into the diffusion process to enhance structure-level detail restoration in image super-resolution. By modulating noise within segmentation regions, our method effectively improves structural fidelity and reduces artifacts without extra inference cost, as demonstrated by extensive experiments on standard benchmarks.

## REFERENCES

- 486  
487  
488 Andreas Aakerberg, Anders S Johansen, Kamal Nasrollahi, and Thomas B Moeslund. Semantic  
489 segmentation guided real-world super-resolution. In *Proceedings of the IEEE/CVF Winter Con-*  
490 *ference on Applications of Computer Vision*, pp. 449–458, 2022.
- 491 Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution:  
492 Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recog-*  
493 *niton workshops*, pp. 126–135, 2017.
- 494 Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chun-  
495 jing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *CVPR*, pp.  
496 12299–12310, 2021.
- 497 Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in  
498 image super-resolution transformer. In *Proceedings of the IEEE/CVF Conference on Computer*  
499 *Vision and Pattern Recognition*, pp. 22367–22377, 2023.
- 500 Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network  
501 for single image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision*  
502 *and pattern recognition*, pp. 11065–11074, 2019.
- 503 Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional  
504 network for image super-resolution. In *ECCV*, pp. 184–199. Springer, 2014.
- 505 Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep  
506 convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):  
507 295–307, 2015.
- 508 Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional  
509 neural network. In *European Conference on Computer Vision*, pp. 391–407. Springer, 2016.
- 510 Azuma Fujimoto, Toru Ogawa, Kazuyoshi Yamamoto, Yusuke Matsui, Toshihiko Yamasaki, and  
511 Kiyoharu Aizawa. Manga109 dataset and creation of metadata. In *Proceedings of the 1st inter-*  
512 *national workshop on comics analysis, processing and understanding*, pp. 1–5, 2016.
- 513 Leon A Gatys, Alexander S Ecker, Matthias Bethge, Aaron Hertzmann, and Eli Shechtman. Con-  
514 trolling perceptual factors in neural style transfer. In *Proceedings of the IEEE conference on*  
515 *computer vision and pattern recognition*, pp. 3985–3993, 2017.
- 516 Juan Mario Haut, Ruben Fernandez-Beltran, Mercedes E Paoletti, Javier Plaza, Antonio Plaza, and  
517 Filiberto Pla. A new deep generative network for unsupervised remote sensing single-image  
518 super-resolution. *IEEE Transactions on Geoscience and Remote sensing*, 56(11):6792–6810,  
519 2018.
- 520 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.  
521 Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in*  
522 *neural information processing systems*, 30, 2017.
- 523 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
524 *neural information processing systems*, 33:6840–6851, 2020a.
- 525 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
526 *neural information processing systems*, 33:6840–6851, 2020b.
- 527 Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from trans-  
528 formed self-exemplars. In *Proceedings of the IEEE conference on computer vision and pattern*  
529 *recognition*, pp. 5197–5206, 2015.
- 530 Yawen Huang, Ling Shao, and Alejandro F Frangi. Simultaneous super-resolution and cross-  
531 modality synthesis of 3d medical images using weakly-supervised joint convolutional sparse cod-  
532 ing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
533 6070–6079, 2017.

- 540 Andrey Ignatov, Radu Timofte, Maurizio Denna, Abdel Younes, Ganzorig Gankhuyag, Jingang  
541 Huh, Myeong Kyun Kim, Kihwan Yoon, Hyeon-Cheol Moon, SeungHo Lee, et al. Efficient  
542 and accurate quantized image super-resolution on mobile npus, mobile ai & aim 2022 challenge:  
543 report. In *ECCV*, pp. 92–129. Springer, 2022.
- 544 Jithin Saji Isaac and Ramesh Kulkarni. Super resolution techniques for medical image processing. In  
545 *2015 International Conference on Technologies for Sustainable Development (ICTSD)*, pp. 1–6.  
546 IEEE, 2015.
- 547 Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and  
548 super-resolution. In *European Conference on Computer Vision*, pp. 694–711. Springer, 2016.
- 549 Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep  
550 convolutional networks. In *CVPR*, pp. 1646–1654, 2016.
- 551 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*  
552 *arXiv:1412.6980*, 2014.
- 553 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
554 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv*  
555 *preprint arXiv:2304.02643*, 2023.
- 556 Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro  
557 Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single  
558 image super-resolution using a generative adversarial network. In *Proceedings of the IEEE*  
559 *conference on computer vision and pattern recognition*, pp. 4681–4690, 2017.
- 560 Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting  
561 Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*,  
562 479:47–59, 2022.
- 563 Wenbo Li, Xin Lu, Jiangbo Lu, Xiangyu Zhang, and Jiaya Jia. On efficient transformer and image  
564 pre-training for low-level vision. *arXiv preprint arXiv:2112.10175*, 3(7):8, 2021.
- 565 Xin Li, Yulin Ren, Xin Jin, Cuiling Lan, Xingrui Wang, Wenjun Zeng, Xinchao Wang, and Zhibo  
566 Chen. Diffusion models for image restoration and enhancement—a comprehensive survey. *arXiv*  
567 *preprint arXiv:2308.09388*, 2023.
- 568 Yachao Li, Dong Liang, Tianyu Ding, and Sheng-Jun Huang. Structsr: Refuse spurious details in  
569 real-world image super-resolution, 2025. URL <https://arxiv.org/abs/2501.05777>.
- 570 Jie Liang, Hui Zeng, and Lei Zhang. Details or artifacts: A locally discriminative learning approach  
571 to realistic image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer*  
572 *Vision and Pattern Recognition*, pp. 5657–5666, 2022a.
- 573 Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir:  
574 Image restoration using swin transformer. In *Proceedings of the IEEE/CVF international confer-*  
575 *ence on computer vision*, pp. 1833–1844, 2021.
- 576 Jingyun Liang, Jiezhong Cao, Yuchen Fan, Kai Zhang, Rakesh Ranjan, Yawei Li, Radu Timofte, and  
577 Luc Van Gool. Vrt: A video restoration transformer. *arXiv preprint arXiv:2201.12288*, 2022b.
- 578 Xinqi Lin, Jingwen He, Ziyang Chen, Zhaoyang Lyu, Ben Fei, Bo Dai, Wanli Ouyang, Yu Qiao,  
579 and Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior. *arXiv*  
580 *preprint arXiv:2308.15070*, 2023.
- 581 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.  
582 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the*  
583 *IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021.
- 584 Zhihe Lu, Zeyu Xiao, Jiawang Bai, Zhiwei Xiong, and Xinchao Wang. Can sam boost video super-  
585 resolution? *arXiv preprint arXiv:2305.06524*, 2023.

- 594 Cheng Ma, Yongming Rao, Yean Cheng, Ce Chen, Jiwen Lu, and Jie Zhou. Structure-preserving  
595 super resolution with gradient guidance. In *Proceedings of the IEEE/CVF conference on computer  
596 vision and pattern recognition*, pp. 7769–7778, 2020.
- 597
- 598 David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented  
599 natural images and its application to evaluating segmentation algorithms and measuring ecological  
600 statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*,  
601 volume 2, pp. 416–423. IEEE, 2001.
- 602 Yiqun Mei, Yuchen Fan, and Yuqian Zhou. Image super-resolution with non-local sparse attention.  
603 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
604 3517–3526, 2021.
- 605
- 606 Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint  
607 arXiv:1411.1784*, 2014.
- 608
- 609 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
610 In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 611 Wenqi Ren, Jinshan Pan, Xiaochun Cao, and Ming-Hsuan Yang. Video deblurring via semantic seg-  
612 mentation and pixel-wise non-linear kernel. In *Proceedings of the IEEE International Conference  
613 on Computer Vision*, pp. 1077–1085, 2017.
- 614
- 615 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
616 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
617 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022a.
- 618
- 619 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
620 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
621 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022b.
- 622
- 623 Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David  
624 Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH*,  
2022a.
- 625
- 626 Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad  
627 Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis  
628 and Machine Intelligence*, 45(4):4713–4726, 2022b.
- 629
- 630 Shuyao Shang, Zhengyang Shan, Guangxing Liu, and Jinglin Zhang. Resdiff: Combining cnn and  
diffusion model for image super-resolution. *arXiv preprint arXiv:2303.08714*, 2023.
- 631
- 632 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: En-  
633 hanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021.
- 634
- 635 Lingchen Sun, Rongyuan Wu, Zhiyuan Ma, Shuaizheng Liu, Qiaosi Yi, and Lei Zhang. Pixel-  
636 level and semantic-level adjustable super-resolution: A dual-lora approach, 2025. URL <https://arxiv.org/abs/2412.03017>.
- 637
- 638 Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao  
639 Li. Maxim: Multi-axis mlp for image processing. In *Proceedings of the IEEE/CVF Conference  
640 on Computer Vision and Pattern Recognition*, pp. 5769–5780, 2022.
- 641
- 642 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-  
643 tion processing systems*, 30, 2017.
- 644
- 645 Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting  
646 diffusion prior for real-world image super-resolution. *arXiv preprint arXiv:2305.07015*, 2023.
- 647
- Peijuan Wang, Bulent Bayram, and Elif Sertel. A comprehensive review on deep learning based  
remote sensing image super-resolution methods. *Earth-Science Reviews*, pp. 104110, 2022a.

- 648 Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image  
649 super-resolution by deep spatial feature transform. In *Proceedings of the IEEE conference on*  
650 *computer vision and pattern recognition*, pp. 606–615, 2018a.
- 651 Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen  
652 Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *European*  
653 *Conference on Computer Vision*, pp. 0–0, 2018b.
- 654 Yingqian Wang, Longguang Wang, Jungang Yang, Wei An, and Yulan Guo. Flickr1024: A large-  
655 scale dataset for stereo image super-resolution. In *Proceedings of the IEEE/CVF International*  
656 *Conference on Computer Vision Workshops*, pp. 0–0, 2019.
- 657 Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li.  
658 Uformer: A general u-shaped transformer for image restoration. In *Proceedings of the IEEE/CVF*  
659 *conference on computer vision and pattern recognition*, pp. 17683–17693, 2022b.
- 660 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment:  
661 from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–  
662 612, 2004.
- 663 Bin Xia, Yulun Zhang, Shiyin Wang, Yitong Wang, Xinglong Wu, Yapeng Tian, Wenming Yang,  
664 and Luc Van Gool. Diffir: Efficient diffusion model for image restoration. *arXiv preprint*  
665 *arXiv:2303.09472*, 2023.
- 666 Zeyu Xiao, Jiawang Bai, Zhihe Lu, and Zhiwei Xiong. A dive into sam prior in image restoration.  
667 *arXiv preprint arXiv:2305.13620*, 2023.
- 668 Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and  
669 Yujie Yang. Maniqa: Multi-dimension attention network for no-reference image quality assess-  
670 ment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
671 pp. 1191–1200, 2022.
- 672 Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-  
673 Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *Proceed-*  
674 *ings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 5728–5739,  
675 2022.
- 676 Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-  
677 representations. In *Curves and Surfaces: 7th International Conference, Avignon, France, June*  
678 *24-30, 2010, Revised Selected Papers 7*, pp. 711–730. Springer, 2012.
- 679 Chaoning Zhang, Dongshen Han, Yu Qiao, Jung Uk Kim, Sung-Ho Bae, Seungkyu Lee, and  
680 Choong Seon Hong. Faster segment anything: Towards lightweight sam for mobile applications.  
681 *arXiv preprint arXiv:2306.14289*, 2023.
- 682 Kai Zhang, Luc Van Gool, and Radu Timofte. Deep unfolding network for image super-resolution.  
683 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.  
684 3217–3226, 2020.
- 685 Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation  
686 model for deep blind image super-resolution. In *IEEE International Conference on Computer*  
687 *Vision*, pp. 4791–4800, 2021.
- 688 Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-  
689 resolution using very deep residual channel attention networks. In *European Conference on*  
690 *Computer Vision*, pp. 286–301, 2018a.
- 691 Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for  
692 image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern*  
693 *recognition*, pp. 2472–2481, 2018b.

700  
701

## A ALGORITHM DETAILS

Here we provide algorithm details of our SAM-DiffSR framework. We adopt the original notations in denoising diffusion probabilistic model (DDPM) (Ho et al., 2020b). Given a data sample  $\mathbf{x}_0 \in p_{\text{data}}$ , the proposed framework in DDPM is defined as:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \quad (7)$$

where  $\mathbf{x}_t$  is the noise latent variable at step  $t$ .  $\beta_1, \dots, \beta_T \in (0, 1)$  are hyperparameters scheduling the scale of added noise for  $T$  steps. Given  $\mathbf{x}_{t-1}$  We can sample  $\mathbf{x}_t$  from this distribution by:

$$\mathbf{x}_t = \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\boldsymbol{\epsilon}_t, \quad (8)$$

where  $\boldsymbol{\epsilon}_t \sim \mathcal{N}(0, \mathbf{I})$ .

In our SAM-DiffSR framework, we use a structural position encoded segmentation mask  $\mathbf{E}_{\text{SAM}}$  to modulate the standard Gaussian noise used in the original DDPM by adding  $\mathbf{E}_{\text{SAM}}$  to  $\boldsymbol{\epsilon}_t$ . Then the sampling of  $\mathbf{x}_t$  becomes:

$$\mathbf{x}_t = \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}(\boldsymbol{\epsilon}_t + \mathbf{E}_{\text{SAM}}), \quad (9)$$

and its corresponding conditional distribution is:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\mathbf{E}_{\text{SAM}}, \beta_t\mathbf{I}). \quad (10)$$

Let  $\alpha_t = 1 - \beta_t$  and iteratively apply Equation 9, we have:

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}(\dots) + \sqrt{\beta_{t-1}}\mathbf{E}_{\text{SAM}} + \sqrt{\beta_{t-1}}\boldsymbol{\epsilon}_{t-1}) + \sqrt{\beta_t}\mathbf{E}_{\text{SAM}} + \sqrt{\beta_t}\boldsymbol{\epsilon}_t \\ &= \sqrt{\alpha_t \dots \alpha_1}\mathbf{x}_0 + (\sqrt{\alpha_t \dots \alpha_2\beta_1} + \dots + \sqrt{\beta_t})\mathbf{E}_{\text{SAM}} + (\sqrt{\alpha_t \dots \alpha_2\beta_1}\boldsymbol{\epsilon}_1 + \dots + \sqrt{\beta_t}\boldsymbol{\epsilon}_t) \\ &= \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \varphi_t\mathbf{E}_{\text{SAM}} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, \end{aligned} \quad (11)$$

where  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ ,  $\varphi_t = \sqrt{\alpha_t \dots \alpha_2\beta_1} + \dots + \sqrt{\beta_t} = \sum_{i=1}^t \sqrt{\bar{\alpha}_t \frac{\beta_i}{\alpha_i}}$ , and  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ .

The corresponding conditional distribution is:

$$q(\mathbf{x}_t|\mathbf{x}_0, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \varphi_t\mathbf{E}_{\text{SAM}}, (1 - \bar{\alpha}_t)\mathbf{I}). \quad (12)$$

Then similar to the original DDPM, we are interested in the posterior distribution that defines the reverse diffusion process. With Bayes' theorem, it can be formulated as:

$$\begin{aligned} p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) &= \frac{p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{x}_0, \mathbf{E}_{\text{SAM}})}{p(\mathbf{x}_t|\mathbf{x}_0, \mathbf{E}_{\text{SAM}})} \\ &\propto \exp\left(-\frac{1}{2}\left(\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}\right)\mathbf{x}_{t-1}^2 - 2\left(\frac{\sqrt{\alpha_t}(\mathbf{x}_t - \sqrt{\beta_t}\mathbf{E}_{\text{SAM}})}{\beta_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \varphi_{t-1}\mathbf{E}_{\text{SAM}}}{1 - \bar{\alpha}_{t-1}}\right)\mathbf{x}_{t-1}\right)\right) \\ &\quad + C(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}), \end{aligned} \quad (13)$$

where  $C(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}})$  not involves  $\mathbf{x}_{t-1}$ . The posterior is also a Gaussian distribution. By using the following notations:

$$\tilde{\beta}_t = 1/\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}\right) = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t, \quad (14)$$

$$\begin{aligned} \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) &= \left(\frac{\sqrt{\alpha_t}(\mathbf{x}_t - \sqrt{\beta_t}\mathbf{E}_{\text{SAM}})}{\beta_t} + \frac{\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \varphi_{t-1}\mathbf{E}_{\text{SAM}}}{1 - \bar{\alpha}_{t-1}}\right) \cdot \tilde{\beta}_t \\ &= \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\left(\frac{\sqrt{1 - \bar{\alpha}_t}}{\sqrt{\beta_t}}\mathbf{E}_{\text{SAM}} + \boldsymbol{\epsilon}\right)\right), \end{aligned} \quad (15)$$

the posterior distribution can be formulated as:

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}}), \tilde{\beta}_t\mathbf{I}). \quad (16)$$

Given latent variable  $\mathbf{x}_t$ , we want to sample from the posterior distribution to obtain the denoised latent variable  $\mathbf{x}_{t-1}$ . This requires the estimation of  $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}})$ , *i.e.*, the estimation of  $\frac{\sqrt{1-\bar{\alpha}_t}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon$ . This is achieved by a parameterized denoising network  $\epsilon_\theta(\mathbf{x}_t, t)$ . The loss function is:

$$\begin{aligned} \mathcal{L}(\theta) &= \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\| \frac{\sqrt{1-\bar{\alpha}_t}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon - \epsilon_\theta(\mathbf{x}_t, t) \|_2^2] \\ &= \mathbb{E}_{t, \mathbf{x}_0, \epsilon} [\| \frac{\sqrt{1-\bar{\alpha}_t}}{\sqrt{\beta_t}} \mathbf{E}_{\text{SAM}} + \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \varphi_t \mathbf{E}_{\text{SAM}} + \sqrt{1-\bar{\alpha}_t} \epsilon, t) \|_2^2]. \end{aligned} \quad (17)$$

This is the loss function in our main paper. Note that the form of  $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0, \mathbf{E}_{\text{SAM}})$  is same to that in the original DDPM. Therefore, our framework requires no change of the generating process and brings no additional inference cost.

## B ABLATION STUDY

**Quality of mask.** Segmentation masks provide the diffusion model structure-level information during training. We conduct experiments to study the impact of using masks with different qualities. Specifically, masks with three qualities are considered: those that are generated by MobileSAM (Zhang et al., 2023) using LR images, those that are generated by MobileSAM using HR images, and those that are generated the original SAM (Kirillov et al., 2023) using HR images. These three kinds of masks are referred to as “Low”, “Medium”, and “High”, respectively. The results of comparing masks with varying qualities are presented in Table 3, indicating that the final performance of the trained model improves as the mask quality increases on both the Urban100 and DIV2k datasets. These findings demonstrate the critical role of high-quality masks in achieving exceptional performance.

**Structural position embedding.** In our SPE module, the RoPE is adopted to generate a 2D position embedding map for obtaining the value assigned to each segmentation area. Here we consider two other approaches: one is using a cosine function to generate a 2D grid as the position embedding map, and the other one is using a linear function whose output value ranges from 0 to 1 to generate the 2D grid. Table 4 shows the corresponding results. Compared with using 2D grids generated with cosine or linear functions, utilizing that generated by RoPE to calculate the value assigned to each segmentation area results in superior performance, thereby showcasing the effectiveness of our SPE module design.

Table 3: Comparison of masks with different qualities.

Mask quality	Urban100			DIV2K		
	PSNR (↑)	SSIM (↑)	FID (↓)	PSNR (↑)	SSIM (↑)	FID (↓)
Low	25.33	0.7702	4.7100	29.09	0.8062	0.4480
Medium	25.40	0.7700	4.7576	29.30	0.8103	0.4176
High	<b>25.54</b>	<b>0.7721</b>	<b>4.5276</b>	<b>29.34</b>	<b>0.8109</b>	<b>0.3809</b>

Table 4: Comparison of different schemes for position embedding.

Position embedding	Urban100			DIV2K		
	PSNR (↑)	SSIM (↑)	FID (↓)	PSNR (↑)	SSIM (↑)	FID (↓)
Cosine	25.28	0.7670	4.7790	28.98	0.8033	0.4689
Linear	25.31	0.7693	4.6197	29.09	0.8073	0.4731
RoPE	<b>25.54</b>	<b>0.7721</b>	<b>4.5276</b>	<b>29.34</b>	<b>0.8109</b>	<b>0.3809</b>

**Non-informative segmentation mask.** There are cases where all pixels in a training sample belongs to the same segmentation area because of the patch-splitting scheme used during training. Two schemes are considered to cope with such non-informative segmentation mask: directly using the original mask, or adopting a special mask filled with fixed values, *i.e.*, zeros. Table 5 presents the comparison results of the above two schemes. Based on the results, it is advantageous to convert non-informative segmentation masks into an all-zero matrix. Our speculation is that the model may be confused by various values in non-informative segmentation masks across different samples, if no reduction is applied to unify such scenarios.

Table 5: Comparison of two schemes for handling non-informative masks. "Reduce" indicates that the mask is replaced with a zero-filled matrix when all pixels belong to the same segmentation area. Otherwise, the original mask is used.

Reduce	Urban100			DIV2K		
	PSNR (↑)	SSIM (↑)	FID (↓)	PSNR (↑)	SSIM (↑)	FID (↓)
$\times$	25.40	0.7687	4.7149	29.18	0.8064	0.4673
$\checkmark$	<b>25.54</b>	<b>0.7721</b>	<b>4.5276</b>	<b>29.34</b>	<b>0.8109</b>	<b>0.3809</b>

**Model performance at different super-resolution scales.** We conducted experiments on the X2 setting, and the results show that our method has a significant performance improvement over the baseline on the reference metric, while maintaining the same level on the unreferenced metric.

Table 6: X2 scale results on test sets of several public benchmarks. (↑) and (↓) indicate that a larger or smaller corresponding score is better, respectively.

Method	Urban100				BSDS100				Manga109				General100				DIV2K			
	PSNR (↑)	SSIM (↑)	MANIQA (↑)	FID (↓)																
SRDiff (X2)	30.84	0.9080	0.5265	0.2067	36.87	0.9667	0.4176	0.0679	30.05	0.8541	0.4545	10.2967	36.43	0.9431	0.4852	6.2866	34.05	0.9178	0.3853	0.0292
SAM-DiffSR (X2)	30.88	0.9095	0.5246	0.2145	37.08	0.9679	0.4192	0.0692	30.36	0.8628	0.4346	10.4271	36.69	0.9458	0.4824	6.4630	34.33	0.9230	0.3832	0.0287

## C DISCUSSION

### C.1 EFFECT OF SEGMENTATION CUES AND COMPUTATIONAL COST

Table 7 is framed as an ablation over three configurations: **SRDiff**, **SRDiff with SAM** (SAM+SRDiff; segmentation cues are directly integrated during training and inference), and **SAN-Diff** (ours; an improvement built upon SRDiff). This ablation targets two questions: (i) whether injecting additional segmentation information benefits diffusion-based SR; and (ii) whether comparable gains can be achieved without increasing training or inference cost.

Comparing **SRDiff** and **SAM+SRDiff** shows that adding segmentation information improves fidelity and realism (e.g., PSNR ↑, FID ↓), but at substantial computational overhead (notably larger parameter count and longer training/inference time). In contrast, **SAN-Diff** achieves improvements similar to **SAM+SRDiff** while keeping the parameter scale and time cost on par with **SRDiff**, thereby preserving efficiency.

For completeness, we include additional diffusion baselines (LDM, StableSR, DiffBIR) as reference. Training times for these models were recorded on heterogeneous hardware and should be interpreted qualitatively. Unless otherwise stated, all other evaluations were conducted on DIV2K using a single V100 GPU.

Table 7: Ablation on segmentation cues and efficiency. The first block (SRDiff, SAM+SRDiff, SAN-Diff) constitutes the main ablation; the second block lists other diffusion models for reference. PSNR (↑) and FID (↓).

	Ablation (ours vs. variants)			Other diffusion models (reference)		
	SRDiff	SAM+SRDiff	SAN-Diff (ours)	LDM	StableSR	DiffBIR
#Params (M)	12	644	12	169	960	1101.8
Train time	2 days	10 days	2 days	6 days	10 days + SD pretrain	6 days + SD pretrain
Inference (s/img)	37.64	65.72	37.62	26.3	238.6	112.4
PSNR ↑	28.6	29.41	29.34	26.45	26.83	26.25
FID ↓	0.4649	0.3938	<b>0.3809</b>	9.5518	14.5232	17.8206

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

## C.2 EXTENSION TO OTHER DIFFUSION TASKS

Our framework has the flexibility to accommodate such tasks seamlessly, as the SAM information functions like a plugin without necessitating alterations to the original diffusion framework. Previous works [1] have demonstrated the efficacy of diffusion-based frameworks across various low-level tasks such as inpainting and deblurring. We are confident that our framework can similarly excel in these areas. However, it's worth noting that our method modifies the diffusion process, which means that simple fine-tuning of pretrained models using parameter efficient approaches like LoRA is not suitable. Instead, retraining the model becomes necessary, which poses computational challenges due to resource constraints. Given these limitations, our paper primarily focuses on the image SR task. Nonetheless, we are committed to expanding our method to encompass a broader range of tasks in the future. We eagerly anticipate collaboration with the computer vision community to further explore these possibilities.

## C.3 REALISTIC FINE-GRAINED TEXTURES

In the field of Image SR, models sometimes generate images with seemingly fine-grained textures, even though the LR images do not contain recognizable textures to the human eye. Defining the correctness of generated texture in such cases presents a challenge. In addressing this issue, we believe that exploring how to generate realistic fine-grained textures within our framework by integrating other kinds of prior information into the model would be a valuable research direction.

## C.4 LIMITATION FROM THE ABILITY OF THE SEGMENTATION MODEL

Compared to the original diffusion model without structural guidance, masks generated by existing SAM models can improve performance, as demonstrated in our experimental results.

However, the performance of our model does depend on the quality of the segmentation masks, as they capture the structural information of the corresponding image. Our model benefits from SAM's fine-grained segmentation capability and its strong generalization ability across diverse objects and textures in the real world. Nevertheless, the performance of our model is also limited by the capabilities of the segmentation model itself. For instance, SAM may struggle to identify structures with low resolution in certain scenes. While the model can partially mitigate this issue by learning from a large amount of data during training, it is undeniable that higher segmentation precision (e.g., SAM2) and finer segmentation granularity would significantly enhance the performance of our approach.

## C.5 SOCIETAL IMPACT

Although our work focuses on improving the performance of diffusion models in super-resolution tasks, the proposed framework can be applied to any task based on diffusion models. This may result in generative models producing higher-quality and more difficult-to-detect deepfakes.

## C.6 SAM INFERENCE RESULT VISUALIZATION

# D MULTI-METRIC COMPARISON

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

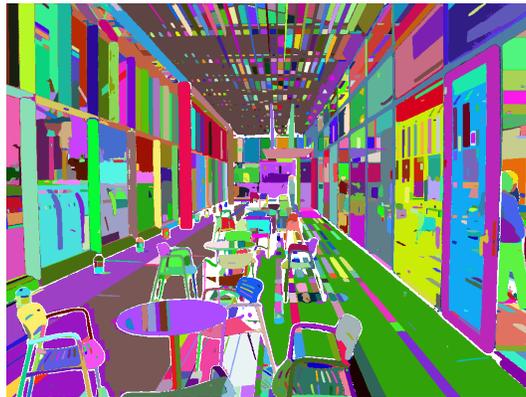


Figure 7: We visualized the results obtained by applying SAM inference to the original images in Figure 1(B). These results are not involved in the inference process. It is only used as a reference for analyzing the super-resolution result.

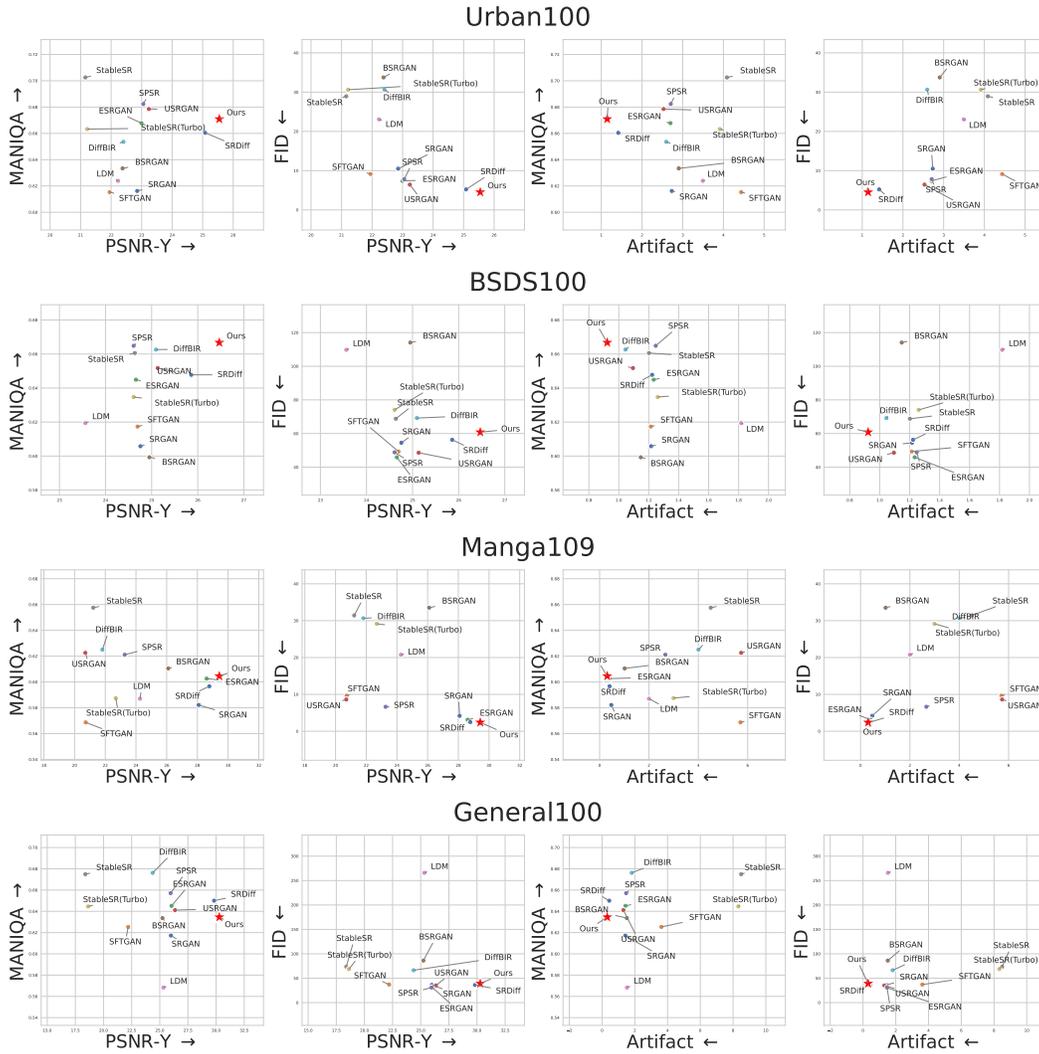


Figure 8: We compared the metrics MANIQA, FID, PSNR, and Artifact across different datasets. In this context, higher values of MANIQA and PSNR are better, while lower values of FID and Artifact are preferred.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

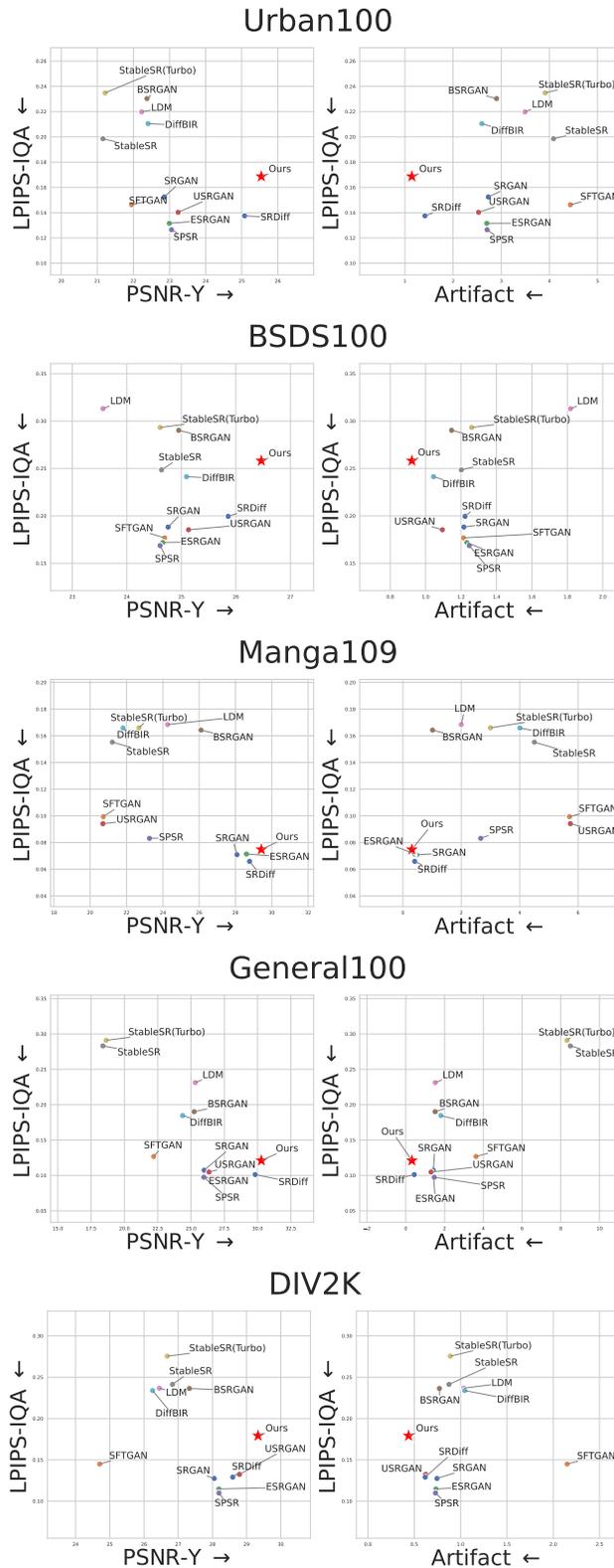


Figure 9: We compared the metrics LPIPS, FID, PSNR, and Artifact across different datasets. In this context, higher values of PSNR is better, while lower values of LPIPS, FID and Artifact are preferred.