# SeeSR: Towards Semantics-Aware Real-World Image Super-Resolution

Rongyuan Wu[1,2], Tao Yang[3], Lingchen Sun[1,2], Zhengqiang Zhang[1,2],
Shuai Li[1,2], Lei Zhang[1,2*]

[1]The Hong Kong Polytechnic University    [2]OPPO Research Institute    [3]ByteDance Inc

## Abstract

*Owe to the powerful generative priors, the pre-trained text-to-image (T2I) diffusion models have become increasingly popular in solving the real-world image super-resolution problem. However, as a consequence of the heavy quality degradation of input low-resolution (LR) images, the destruction of local structures can lead to ambiguous image semantics. As a result, the content of reproduced high-resolution image may have semantic errors, deteriorating the super-resolution performance. To address this issue, we present a semantics-aware approach to better preserve the semantic fidelity of generative real-world image super-resolution. First, we train a degradation-aware prompt extractor, which can generate accurate soft and hard semantic prompts even under strong degradation. The hard semantic prompts refer to the image tags, aiming to enhance the local perception ability of the T2I model, while the soft semantic prompts compensate for the hard ones to provide additional representation information. These semantic prompts can encourage the T2I model to generate detailed and semantically accurate results. Furthermore, during the inference process, we integrate the LR images into the initial sampling noise to mitigate the diffusion model's tendency to generate excessive random details. The experiments show that our method can reproduce more realistic image details and hold better the semantics. Code: https://github.com/cswry/SeeSR*

## 1. Introduction

Images inevitably undergo degradation due to factors such as subpar imaging devices, unfavorable capturing environments, transmission losses. This degradation manifests in various forms, including low-resolution, blurriness and noise. Image super-resolution (ISR) aims to reconstruct a high-resolution (HR) image from the given low-resolution (LR) input. Traditionally, researchers investigate the ISR problem by assuming simple and known image degrada-

tions (*e.g.*, bicubic downsampling), and developed many successful models [7, 8, 13, 28, 68, 69]. However, these methods often yield over-smoothed outcomes due to their fidelity-focused learning objectives. To enhance visual perception, generative adversarial networks (GANs) [14] have been adopted to solve the ISR problem [50]. By using the adversarial loss in training, the ISR models can be supervised to generate perceptually realistic details, enhancing the visual quality but in the price of sacrificing fidelity.

Despite the remarkable advancements, when applying the above mentioned models to real-world LR images, whose degradations are much more complex and even unknown, the output HR images can have low visual quality with many artifacts. This is mainly because the domain gap between the synthetic training data and the real-world test data. The goal of real-world ISR (Real-ISR) is to reproduce a perceptually realistic HR image from its LR observation with complex and unknown degradation. To this end, some researchers proposed to collect real-world LR-HR image pairs using long-short camera focal lens [3, 55]. Another more cost-effective way is to simulate the complex real-world image degradation process using random combinations of basic degradation operations. The representative work along this line include BSRGAN [64], Real-ESRGAN [51] and their variants [6, 29, 30, 57]. With the abundant amount of more realistic synthetic training pairs, the GAN-based Real-ISR methods can generate more authentic details. However, they still tend to introduce many unpleasant visual artifacts due to the unstable adversarial training. The LDL [29] can suppress much the visual artifacts by detecting the problematic pixels using local image statistics. Unfortunately, it is not able to generate additional details.

Recently, denoising diffusion probabilistic models (DDPMs) [21] have exhibited remarkable performance in the realm of image generation, gradually emerging as successors to GANs in various downstream tasks [40, 42]. Some researchers [25, 52] have leveraged pretrained DDPMs to effectively tackle the inverse image restoration problems. However, their application to the challenging Real-ISR scenarios is hindered by the assumptions of known linear degradation model. Considering that the

---

*Corresponding author.

1

large-scale pretrained text-to-image (T2I) models [40, 42], which are trained on a dataset exceeding 5 billion image-text pairs, encompass more potent natural image priors, some methods have recently emerged to harness their potentials to address the Real-ISR problem, including StableSR [48], PASD [59] and DiffBIR [35]. These diffusion prior based Real-ISR methods have demonstrated highly promising capability to generate realistic image details; however, they still have some limitations. StableSR [48] and Diff-BIR [35] solely rely on input LR images as control signals, overlooking the role of semantic text information in the pretrained T2I models. PASD [59] attempts to utilize off-the-shelf high-level models to extract semantic prompts as additional control conditions for the T2I model. However, it encounters difficulties when dealing with scenes containing a variety of objects or severely degraded images.

In this work, we investigate in-depth the problem that how to extract more effective semantic prompts to harness the generative potential of pretrained T2I models so that better Real-ISR results can be obtained. By analyzing the effects of different types of semantic prompts on the Real-ISR outcomes, we conclude two major criteria. Firstly, the prompt should cover as many objects in the scene as possible, helping the T2I model to understand different local regions of the LR image. Secondly, the prompt should be degradation-aware to avoid erroneous semantic restoration results. (Please refer to Section 3.1 for more discussions.) While the prompt extractor undergoes low-level data augmentation during training [17], there still exists much gap between this augmentation and real-world degradation. Hence, it is not suitable for directly extracting semantic prompts from real-world LR inputs.

Based on the aforementioned criteria, we present a **Se**mantic-awar**e SR** (**SeeSR**) approach, which utilizes high-quality semantic prompts to enhance the generative capacity of pretrained T2I models for Real-ISR. SeeSR consists of two stages. In the first stage, the semantic prompt extractor is fine-tuned to acquire degradation-aware capabilities. This enables it to extract accurate semantic information from LR images as soft and hard prompts. In the second stage, the pristine semantic prompts collaborate with LR images to exert precise control over the T2I model, facilitating the generation of rich and semantically correct details. Moreover, during inference stage, we incorporate the LR image into the initial sampling noise to alleviate the diffusion model's propensity for generating excessive random details. Our extensive experiments demonstrate the superior realistic detail generation performance of SeeSR while preserving well the image semantics of Real-ISR outputs.

## 2. Related Work

**GAN-based Real-ISR.** Starting from SRCNN [13], deep learning based ISR has become prevalent. A variety of methods focusing on model design have been proposed [7–10, 28, 31, 68–70] to improve the accuracy of ISR reconstruction. However, most of these methods assume simple and known degradations such as bicubic downsampling, limiting their effectiveness when dealing with complex and unknown degradations in the real world. Recent advancements in Real-ISR have explored more complex degradation models to approximate the real-world degradations. Specifically, BSRGAN [64] introduces a randomly shuffled degradation modeling strategy, while Real-ESRGAN [51] employs a high-order degradation modeling process. Using the training samples with more realistic degradations, both BSRGAN and Real-ESRGAN utilize GANs [14] to reconstruct desired HR images. While generating more perceptually realistic details, training GANs is unstable and Real-ISR outputs often suffer from unnatural visual artifacts. Many following works such as LDL [29] and DeSRA [57] can suppress much the artifacts, yet they are difficult to generate more natural details.

**Diffusion Probabilistic Models.** Inspired by the non-equilibrium thermodynamics theory [23] and sequential MonteCarlo [37], Sohl-Dickstein *et al*. [45] proposed the diffusion model to model complex datasets. Subsequently, a series of fruitful endeavors [11, 21, 46] have been made to apply diffusion models in the realm of image generation, especially since the development of DDPM [21]. Rombach *et al*. [40] expanded the training of DDPMs to the latent space, greatly facilitating the development of large-scale pretrained text-to-image (T2I) diffusion models such as stable diffusion (SD) [1] and Imagen [41]. It has been demonstrated that T2I diffusion priors are powerful in image editing [36, 66], video generation [44, 56], 3D content generation [32, 54, 54], *etc*.

**Diffusion Prior based Real-ISR.** Early attempts [25, 43, 52] using DDPMs to address the ISR problem are mostly assuming simple downsampling degradation. However, such an assumption of known linear image degradation restricts their practical application in complex scenarios like Real-ISR. Recently, some researchers [35, 48, 59] have employed powerful pretrained T2I models such as SD [1] to tackle the real-ISR problem. Having been trained on billions of image-text pairs, these models can perceive strong image priors for tackling Real-ISR challenges. StableSR [48] achieves this goal by training a time-aware encoder to fine-tune the SD model and employing feature warping to balance between fidelity and perceptual quality. DiffBIR [35] adopts a two-stage strategy to tackle the Real-ISR problem. It first reconstructs the image as an initial estimation, and then utilizes the SD prior to enhance image details.

The aforementioned methods solely rely on images as conditions to activate the generative capability of the T2I model. In contrast, PASD [59] goes further by utilizing off-the-shelf high-level models (*i.e*., ResNet [16], Yolo [39] and
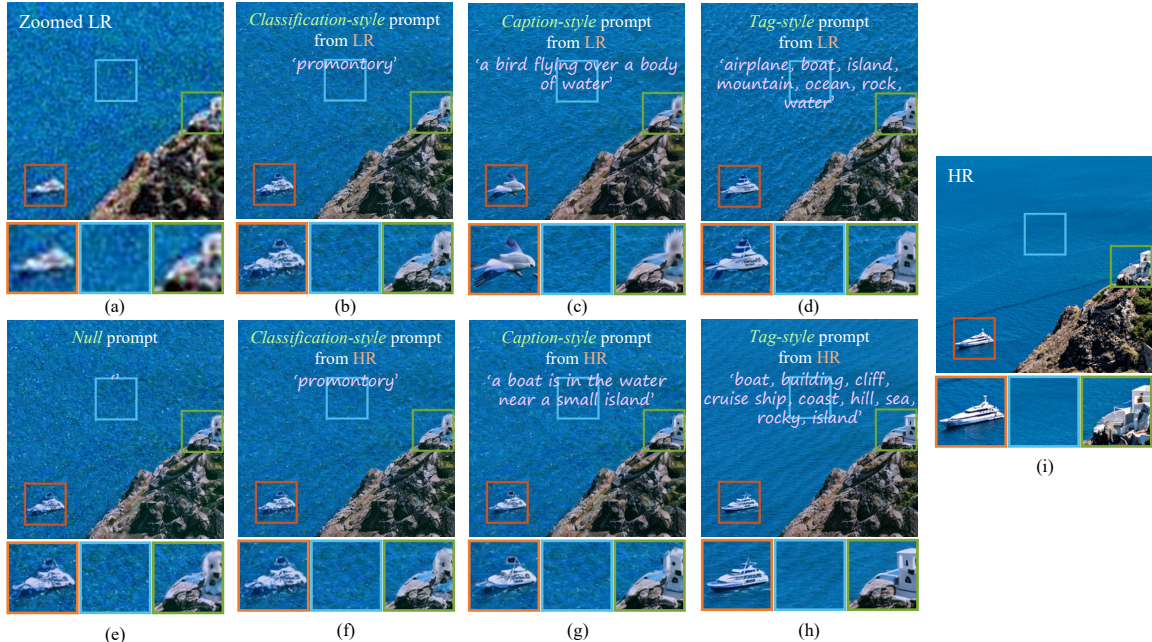
Figure 1. The comparison of different styles of prompts and their corresponding Real-ISR results with PASD [59]. (a) Input LR image. (b)-(d) show the extracted classification-style, caption-style and tag-style prompts from LR image and the corresponding Real-ISR results. (e) Null prompt and its corresponding Real-ISR result. (f)-(h) show the extracted classification-style, caption-style and tag-style prompts from HR image and their corresponding Real-ISR results. (i) HR image.

BLIP [27]) to extract semantic information to guide the diffusion process, stimulating more generative capacity of the T2I model. However, ResNet and Yolo have limited object recognition ability, leading to a diminished recall rate. The captions generated by BLIP struggle to comprehensively describe the semantic information in images, particularly in scenes with a rich diversity of objects. Therefore, how to introduce prompts to more effectively elicit the potential of pretrained T2I models in assisting Real-ISR needs deep investigation, which is the goal of this work.

## 3. Methodology

### 3.1. Motivation and Framework Overview

**Motivation.** To unleash the generative potential of pretrained T2I model while avoiding semantic distortion in Real-ISR outputs, we investigate the use of three representative styles of semantic prompts, including *classification-style*, *caption-style* and *tag-style*. In specific, we use the methods in [16] , [27] and [71] to extract classification-style, caption-style and tag-style prompts, respectively.

The classification-style prompt provides only one category label for the entire image, which is robust to image degradation due to its global view. However, such kind of prompts lack the ability to provide semantic support of local objects, particularly in scenes containing multiple entities. As shown in Figs. 1(b) and 1(f), by using the classification-

style prompts extracted from the LR and the HR images, the Real-ISR results are almost indistinguishable from that obtained by using the null prompt (see Fig. 1(e)).

The caption-style prompt provides a sentence to describe the corresponding image, offering richer information compared to the classification-style prompt. However, it still has two shortcomings. Firstly, the redundant prepositions and adverbs in this type of prompt may scatter the attention of T2I models towards degraded objects [18]. Secondly, it is prone to semantic errors due to the influence of degradation in LR images. As shown in Fig. 1(c), the T2I model mistakenly reconstructs a bird instead of a ship due to the incorrect caption extracted from the LR image.

The tag-style prompt provides category information for all objects in the image, offering a more detailed description of the entities compared to caption-style prompt. Even without providing object location information, it is found that the T2I model can align the semantic prompts with the corresponding regions in the image due to its underlying semantic segmentation capability [62]. Unfortunately, similar to the captioning models, the tagging models are also susceptible to image degradations, resulting in erroneous semantic cues and semantic distortion in the reconstructed results. As shown in Fig. 1(d), the wrong semantic prompt "airplane" leads to distorted reconstruction of the ship.

We summarize the characteristics of different styles of prompts in Table 1. This motivates us that if we can adapt

3

Table 1. Comparison of different prompt styles.

| | Rich Objects | Concise Description | Degradation Aware |
|---|---|---|---|
| Classification-style | ✗ | ✓ | ✓ |
| Caption-style | ✓ | ✗ | ✗ |
| Tag-style | ✓✓ | ✓ | ✗ |
| Our DAPE | ✓✓ | ✓ | ✓ |

the tag-style prompt to be degradation-ware, then it may help the T2I models generate high-quality Real-ISR outputs while preserving correct image semantics.

**Framework Overview.** Based on the above discussions, we propose to extract high-quality tag-style prompts from the LR image to guide the pretrained T2I model, such as stable diffusion (SD) [40], for producing semantics-preserved Real-ISR results. The framework of our proposed method, namely **Se**mantic-awar**e SR** (**SeeSR**), is shown in Fig. 2. The training of SeeSR goes through two stages. In the first stage (Fig. 2(a)), we learn a degradation-aware prompt extractor (DAPE), which consists of an image encoder and a tagging head. It is expected that both the feature representations and tagging outputs of the LR image can be as close as possible to that of the corresponding HR image by using the original tag model. The learned DAPE is copied to the second stage (Fig. 2(b)) to extract the feature representations and tags (as text prompts) from the input LR image, which serve as control signals over the pretrained T2I model to generate visually pleasing and semantically correct Real-ISR results. During inference, only the second stage is needed to process the input image. Fig. 2(c) illustrates the collaborative interplay between the image branch, feature representation branch, and text prompt branch in governing the pretrained T2I model.

## 3.2. Degradation-Aware Prompt Extractor

The DAPE is fine-tuned from a pretrained tag model, *i.e.*, RAM [71]. As depicted in Fig. 2(a), the HR image $x$ goes through a frozen tag model to output representation embedding $f_x^{rep}$ and logits embedding $f_x^{logits}$ as anchor points to supervise the training of DAPE. LR images $y$ are obtained by applying random degradations to $x$, and they are fed into the trainable image encoder and tagging head. To make DAPE robust to image degradation, we force the representation embedding and logits embedding from the LR branch to be close to that of the HR branch. The training objective is as follows:

$$\mathcal{L}_{DAPE} = \mathcal{L}_r(f_y^{rep}, f_x^{rep}) + \lambda \mathcal{L}_l(f_y^{logits}, f_x^{logits}), \quad (1)$$

where $\lambda$ is a balance parameter, $f_y^{rep}$ and $f_y^{logits}$ are the representation embedding and logits embedding from LR

branch. $\mathcal{L}_r$ is the mean squared error (MSE) loss, while $\mathcal{L}_l$ is the cross-entropy loss [16]. By aligning the outputs from LR and HR branches, DAPE is learned to predict high-quality semantic prompts from corrupted image inputs.

Once trained, DAPE undertakes the crucial role of extracting reliable semantic prompts from the LR images. The prompts can be classified into two categories: hard prompts (*i.e.*, tag texts from the tagging head) and soft prompts (*i.e.*, representation embeddings from the image encoder). As shown in Figs. 2(b) and 2(c), hard prompts are directly passed to the frozen text encoder built into the T2I model to enhance its local understanding capability. The abundance of text prompts is controlled by a preset threshold. If the threshold is too high, the accuracy of predicted categories will improve but the recall rate can be affected, and vice versa. Therefore, the soft label prompts are used to compensate for the limitations of hard prompt, which are free of the impact threshold and avoid the low information entropy issue caused by one-hot categories [20].

## 3.3. Training of SeeSR Model

Fig. 2(c) illustrates the detailed structure of the controlled T2I diffusion model. Given the successful application of ControlNet [66] in conditional image generation, we utilize it as the controller of the T2I model for Real-ISR purpose. In specific, we clone the encoder of the Unet in pre-trained SD model as a trainable copy to initialize the ControlNet. To incorporate soft prompts into the diffusion process, we adopt the cross-attention mechanism proposed in PASD [59] to learn semantic guidance. The representation cross-attention (RCA) modules are added to the Unet and placed after the text cross-attention (TCA) modules. Note that the randomly initialized RCA modules are cloned simultaneously with the encoder. In addition to the text branch and representation branch, the image branch also plays a role in reconstructing the desired HR image. We pass the LR images through a trainable image encoder to obtain the LR latent, which is input to ControlNet. The structure of trainable image encoder is kept the same as that in [66].

The training process of the SeeSR model is as follows. The latent representation of an HR image is obtained by the encoder of pretrained VAE [40], denoted as $z_0$. The diffusion process progressively introduces noise to $z_0$, resulting in a noisy latent $z_t$, where $t$ represents the randomly sampled diffusion step. With the diffusion step $t$, LR latent $z_{lr}$, hard prompts $p_h$ and soft prompts $p_s$, we train our SeeSR network, denoted as $\epsilon_\theta$, to estimate the noise added to the noisy latent $z_t$. The optimization objective is:

$$\mathcal{L} = \mathbb{E}_{z_0, z_{lr}, t, p_h, p_s, \epsilon \sim \mathcal{N}} \left[ \| \epsilon - \epsilon_\theta (\mathbf{z}_t, z_{lr}, t, p_h, p_s) \|_2^2 \right]. \quad (2)$$

For saving the training cost, we freeze the parameters of the SD model while training solely on the newly added mod-

(a) Degradation-aware prompt extractor



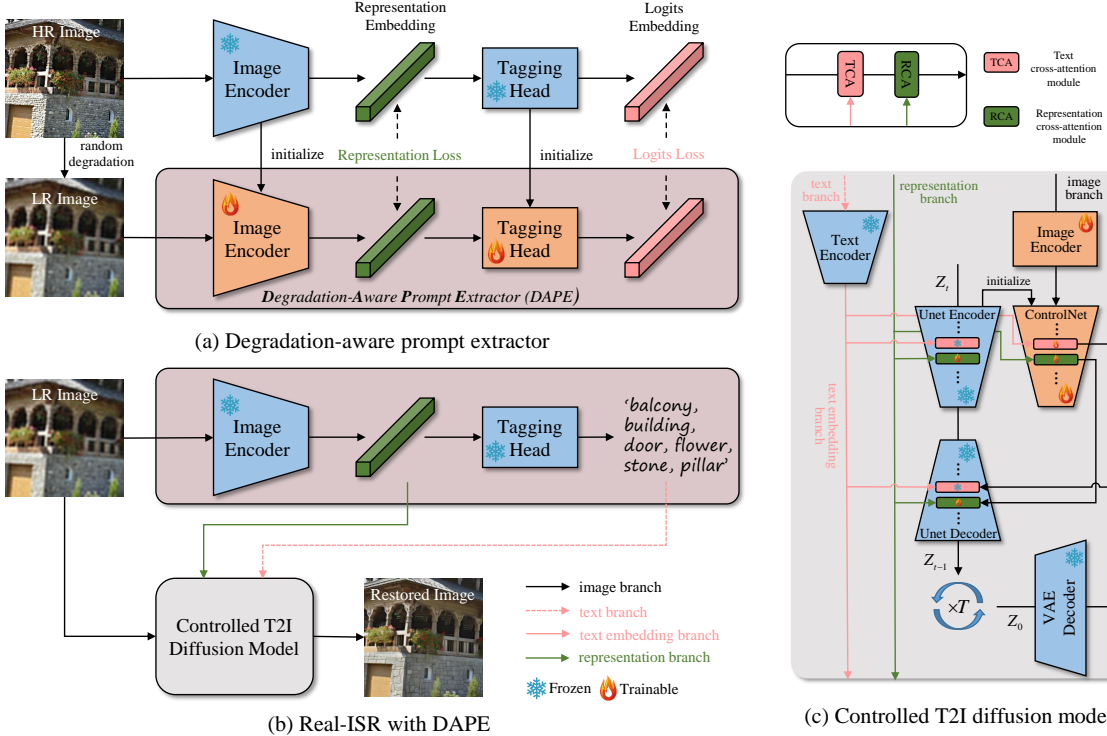(b) Real-ISR with DAPE



(c) Controlled T2I diffusion model

Figure 2. Overview of SeeSR. (a) In the first stage, we train a degradation-aware prompt extractor (DAPE), which is initialized from a tag model. DAPE is trained to align the encoding of the degraded LR image to the encoding of the corresponding HR image by a tag model (*e.g.*, RAM [71] in our work), enabling DAPE the degradation-awareness. (b) In the second stage, the well-trained DAPE provides both soft prompts (representation embedding) and hard prompts (tagging text), which are combined with the LR image to control a pretrained T2I model (*e.g.*, SD [40] in our work). The detailed structure of the controlled T2I diffusion model is shown in (c).

ules, including the image encoder, the ControlNet and the RCA modules within the Unet.

## 3.4. LR Embedding in Inference

The pretrained T2I models such as SD, during their training phase, do not completely convert the images into random Gaussian noises. However, during the inference process, most of existing SD-based Real-ISR methods [35, 48, 59] take a random Gaussian noise as their start point, leading to a discrepancy on the noise handling procedure between training and inference [33]. In the Real-ISR task, we observe that this discrepancy can confuse the model to perceive degradation as content to be enhanced, particularly in smooth regions such as the sky, as shown in the top row of Fig. 3. To address this issue, we propose to directly embed the LR latent into the initial random Gaussian noise according to the training noise scheduler. This strategy is applicable to most of the SD-based Real-ISR methods [35, 48, 59]. As shown in the bottom row of Fig. 3, the proposed LR embedding (LRE) strategy alleviate much the inconsistency between training and inference, providing a more faithful start point for the diffusion model and consequently suppressing much the artifacts in the sky region. Note that all

experiments of SeeSR in the subsequent sections utilize the LRE strategy by default.



Figure 3. Effectiveness of the LR embedding (LRE) strategy in alleviating the discrepancy between training and inference of SD-based Real-ISR methods [35, 48, 59]. Top row: results without using LRE. Bottom row: results with LRE. We see that many falsely generated details in the sky area are removed.

## 4. Experiments

Following previous works [51, 64], we focus on the challenging ×4 Real-ISR tasks, while the proposed method can be applied to other scaling factors. Furthermore, to validate the semantic aware capability of SeeSR, we compare

Table 2. Quantitative comparison with state-of-the-art methods on both synthetic and real-world benchmarks. The best and second best results of each metric are highlighted in **red** and blue, respectively. LDM is not tested on *RealLR200* because the related codebase does not provide tiled functionality, which results in the issue of out-of-memory when testing on higher-resolution inputs.

| Datasets | Metrics | BSRGAN [64] | Real-ESRGAN [51] | LDL [29] | DASR [30] | FeMaSR [5] | LDM [40] | StableSR [48] | ResShift [60] | PASD [59] | DiffBIR [35] | SeeSR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *DIV2K-Val* | PSNR ↑ | 21.87 | **21.94** | 21.52 | 21.72 | 20.85 | 21.26 | 20.84 | 21.75 | 20.77 | 20.94 | 21.04 |
| | SSIM ↑ | 0.5539 | **0.5736** | 0.5690 | 0.5536 | 0.5163 | 0.5239 | 0.4887 | 0.5422 | 0.4958 | 0.4938 | 0.5341 |
| | LPIPS ↓ | 0.4136 | **0.3868** | 0.3995 | 0.4266 | 0.3973 | 0.4154 | 0.4055 | 0.4284 | 0.4410 | 0.4270 | 0.3876 |
| | DISTS ↓ | 0.2737 | 0.2601 | 0.2688 | 0.2688 | 0.2428 | 0.2500 | 0.2542 | 0.2606 | 0.2538 | 0.2471 | **0.2298** |
| | FID ↓ | 64.28 | 53.46 | 58.94 | 67.22 | 53.70 | 41.93 | 36.57 | 55.77 | 40.77 | 40.42 | **32.79** |
| | NIQE ↓ | 4.7615 | 4.9209 | 5.0249 | 4.8596 | **4.5726** | 6.4667 | 4.6551 | 6.9731 | 4.8328 | 4.7211 | 5.0120 |
| | MANIQA ↑ | 0.4834 | 0.5251 | 0.5127 | 0.4346 | 0.4869 | 0.5237 | 0.5914 | 0.5232 | 0.6049 | 0.6205 | **0.6236** |
| | MUSIQ ↑ | 59.11 | 58.64 | 57.90 | 54.22 | 58.10 | 56.52 | 62.95 | 58.23 | 66.85 | 65.23 | **68.29** |
| | CLIPIQA ↑ | 0.5183 | 0.5424 | 0.5313 | 0.5241 | 0.5597 | 0.5695 | 0.6486 | 0.5948 | 0.6799 | 0.6664 | **0.6834** |
| *RealSR* | PSNR ↑ | 26.39 | 25.69 | 25.28 | **27.02** | 25.07 | 25.48 | 24.70 | 26.31 | 24.29 | 24.77 | 24.90 |
| | SSIM ↑ | 0.7654 | 0.7616 | 0.7567 | **0.7708** | 0.7358 | 0.7148 | 0.7085 | 0.7421 | 0.6630 | 0.6572 | 0.7126 |
| | LPIPS ↓ | **0.2670** | 0.2727 | 0.2766 | 0.3151 | 0.2942 | 0.3180 | 0.3018 | 0.3460 | 0.3435 | 0.3658 | 0.3133 |
| | DISTS ↓ | 0.2121 | **0.2063** | 0.2121 | 0.2207 | 0.2288 | 0.2213 | 0.2135 | 0.2498 | 0.2259 | 0.2310 | 0.2318 |
| | FID ↓ | 141.28 | 135.18 | 142.71 | 132.63 | 141.05 | 132.72 | **128.51** | 141.71 | 129.76 | 128.99 | 134.90 |
| | NIQE ↓ | 5.6567 | 5.8295 | 6.0024 | 6.5311 | 5.7885 | 6.5200 | 5.9122 | 7.2635 | **5.3628** | 5.5696 | 5.4612 |
| | MANIQA ↑ | 0.5399 | 0.5487 | 0.5485 | 0.3878 | 0.4865 | 0.5423 | 0.6221 | 0.5285 | 0.6493 | 0.6253 | **0.6506** |
| | MUSIQ ↑ | 63.21 | 60.18 | 60.82 | 40.79 | 58.95 | 58.81 | 65.78 | 58.43 | 68.69 | 64.85 | **69.67** |
| | CLIPIQA ↑ | 0.5001 | 0.4449 | 0.4477 | 0.3121 | 0.5270 | 0.5709 | 0.6178 | 0.5444 | 0.6590 | 0.6386 | **0.6630** |
| *DrealSR* | PSNR ↑ | 28.75 | 28.64 | 28.21 | **29.77** | 26.90 | 27.98 | 28.13 | 28.46 | 27.00 | 26.76 | 27.90 |
| | SSIM ↑ | 0.8031 | 0.8053 | 0.8126 | **0.8264** | 0.7572 | 0.7453 | 0.7542 | 0.7673 | 0.7084 | 0.6576 | 0.7628 |
| | LPIPS ↓ | 0.2883 | 0.2847 | **0.2815** | 0.3126 | 0.3169 | 0.3405 | 0.3315 | 0.4006 | 0.3931 | 0.4599 | 0.3299 |
| | DISTS ↓ | 0.2142 | **0.2089** | 0.2132 | 0.2271 | 0.2235 | 0.2259 | 0.2263 | 0.2656 | 0.2515 | 0.2749 | 0.2363 |
| | FID ↓ | 155.63 | **147.62** | 155.53 | 155.58 | 157.78 | 156.01 | 148.98 | 172.26 | 159.24 | 166.79 | 151.88 |
| | NIQE ↓ | 6.5192 | 6.6928 | 7.1298 | 7.6039 | 5.9073 | 7.1677 | 6.5354 | 8.1249 | **5.8595** | 6.2935 | 6.4893 |
| | MANIQA ↑ | 0.4878 | 0.4907 | 0.4914 | 0.3879 | 0.4420 | 0.5043 | 0.5591 | 0.4586 | 0.5850 | 0.5923 | **0.6005** |
| | MUSIQ ↑ | 57.14 | 54.18 | 53.85 | 42.23 | 53.74 | 53.73 | 58.42 | 50.60 | 64.81 | 61.19 | **65.11** |
| | CLIPIQA ↑ | 0.4915 | 0.4422 | 0.4310 | 0.3684 | 0.5464 | 0.5706 | 0.6206 | 0.5342 | **0.6773** | 0.6346 | 0.6708 |
| *RealLR200* | NIQE ↓ | 4.3817 | 4.2048 | 4.3845 | 4.3360 | 4.6357 | - | 4.2516 | 6.2878 | **4.1715** | 4.9330 | 4.3332 |
| | MANIQA ↑ | 0.5462 | 0.5582 | 0.5519 | 0.4877 | 0.5295 | - | 0.5841 | 0.5417 | 0.6066 | 0.5902 | **0.6198** |
| | MUSIQ ↑ | 64.87 | 62.94 | 63.11 | 55.67 | 64.14 | - | 63.30 | 60.18 | 68.20 | 62.06 | **69.52** |
| | CLIPIQA ↑ | 0.5679 | 0.5389 | 0.5326 | 0.4659 | 0.6522 | - | 0.6068 | 0.6486 | 0.6797 | 0.6509 | **0.6814** |

the Real-ISR methods using the well-known COCO dataset [34].

## 4.1. Experimental Settings

**Training Datasets.** We train SeeSR on DIV2K [2], DIV8K [15], Flickr2K [47], OST [49], and the first 10K face images from FFHQ [24]. The degradation pipeline of Real-ESRGAN [51] is used to synthesize LR-HR training pairs.

**Test Datasets.** We employ the following test datasets to comprehensively evaluate SeeSR. (1) First, we randomly crop 3K patches (resolution: $512 \times 512$) from the DIV2K validation set [2] and degrade them using the same pipeline as that in training. We name this dataset as *DIV2K-Val*. (2) We employ the two real-world datasets, *RealSR* [3] and *DRealSR* [55], by using the same configuration as [48] to center-crop the LR image to $128 \times 128$ [1]. (3) We build another real-world dataset, named *RealLR200*, which comprises 38 LR images used in recent literature [29, 51, 63], 47 LR images from DiffBIR [35], 50 LR images from VideoLQ (the last frame of each video sequence) [4], and 65 LR images collected from the internet by ourselves.

**Implementation Details.** We finetune the entire DAPE module from RAM [71] using LORA ($r = 8$) [22] for 20k iterations, where the batch size and the learning rate are set

to 32 and $10^{-4}$, respectively. The SD 2.1-base[2] is used as the pretrained T2I model. The whole controlled T2I model is finetuned for 100K iterations with Adam [26] optimizer, where the batch size and learning rate are respectively set to 32 and $5 \times 10^{-5}$. The training process is conducted on $512 \times 512$ resolution images with 8 NVIDIA Tesla 32G-V100 GPUs. For inference, we adopt spaced DDPM sampling [38] with 50 timesteps. $\lambda$ in Eq. (1) is set to 1.

**Evaluation Metrics.** In order to provide a comprehensive and holistic assessment of the performance of different methods, we employ a range of reference and non-reference metrics. PSNR and SSIM [53] (calculated on the Y channel in YCbCr space) are reference-based fidelity measures, while LPIPS[3] [67], DISTS [12] are reference-based perceptual quality measures. FID [19] evaluates the distance of distributions between original and restored images. NIQE [65], MANIQA [58], MUSIQ [58], and CLIPIQA [58] are non-reference image quality measures.

**Compared Methods.** We compare our SeeSR with several state-of-the-art Real-ISR methods, which can be categorized into two groups. The first group consists of GAN-based methods, including BSRGAN [64], Real-ESRGAN [51], LDL [29], FeMaSR [5] and DASR [30]. The second group consists of recent diffusion-based methods, includ-

---

[1] https://huggingface.co/datasets/Iceclear/StableSR-TestSets

[2] https://huggingface.co/stabilityai/stable-diffusion-2-1-base
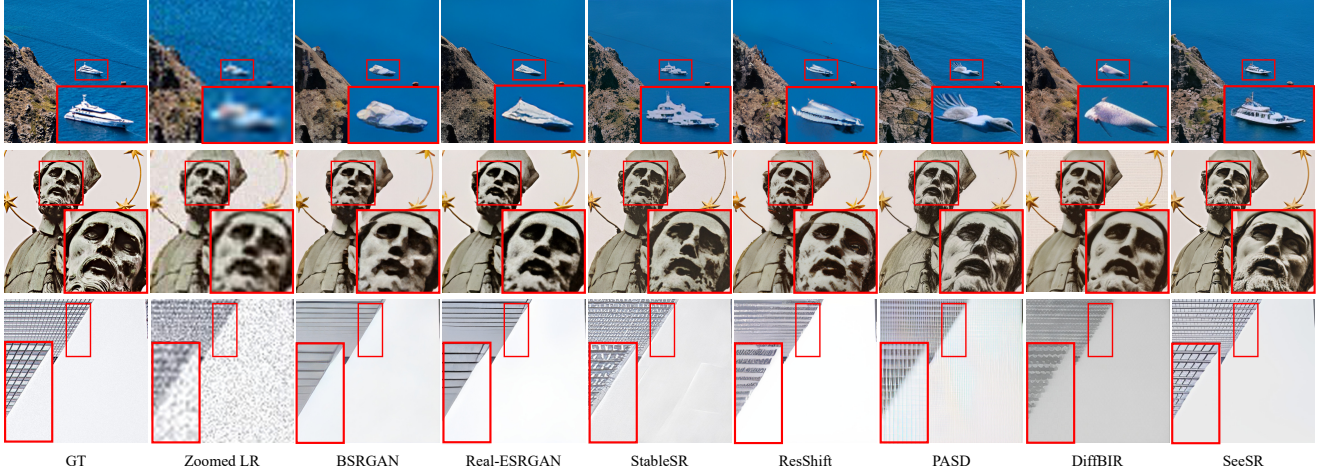[3] We use LPIPS-Alex by default.

Figure 4. Qualitative comparisons of different Real-ISR methods on three images. Please zoom in for a better view.

ing LDM [40], StableSR [48], ResShift [60], PASD [59], and DiffBIR [35]. We use the publicly released codes and models of the competing methods for testing.

## 4.2. Comparisons with State-of-the-Arts

**Quantitative Comparisons.** We first show the quantitative comparison on the four synthetic and real-world datasets in Table 2. We have the following observations. (1) First, our SeeSR consistently achieves the best scores in MANIQA, CLIPIQA and MUSIQ across all the four datasets, except for a 1% lower score in CLIPIQA on *DrealSR*. (2) Second, SeeSR achieves the best FID and DISTS scores on *DIV2K-Val*, surpassing the second-best method by more than 10.3% and 5.3%, respectively. (3) GAN-based methods achieve better PSNR/SSIM scores than DM-based methods. This is mainly because DM-based methods can generate more realistic details, which however sacrifice the fidelity. (4) BSRGAN, Real-ESRGAN and LDL show advantages in terms of reference perceptual metrics LPIPS/DISTS, but they perform poorer in no-reference perceptual metrics such as CLIPIQA, MUSIQ and MANIQA. This is also because DM-based methods will generate some structures and textures that may not match the GTs, making them disadvantageous in full-reference metrics.

Overall, compared with other DM-based methods, our SeeSR achieves better no-reference metric scores, while keep competitive full-reference measures.

**Qualitative Comparisons.** Figs. 4 and 5 present visual comparisons on synthetic and real-world images, respectively. As illustrated in the first row of Fig. 4, BSRGAN and Real-ESRGAN fail to reconstruct the details of the ship, which suffers from severe degradation. Meanwhile, some DM-based methods can reconstruct clear details but exhibit obvious semantic errors. Due to the ambiguous output of its degradation removal stage, DiffBIR mis-generates the

ship into fish. The caption model of PASD provides a text prompt with semantic errors, wrongly generating a bird. In comparison, our well-trained DAPE module in SeeSR can still provide accurate prompt even with strong degradation, aiding SeeSR to generate semantically-accurate and details-rich results. Additionally, StableSR, DiffBIR and PASD all exhibit certain degree of artifacts in smooth regions (*e.g.*, sea, sky). Thanks to our LRE strategy, SeeSR demonstrates better stability in the smooth regions.

Similar conclusions can be drawn from Fig. 5 for real-world LR images. GAN-based methods generate limited and unnatural details. Some DM-based methods can generate more realistic details but with less accurate semantics. DiffBIR generates dense stripes in smooth regions (row 1). Additionally, it adds extra glasses to the person (row 3). StableSR, ResShift and PASD fail to recover edges (row 1). Although PASD does a decent job in restoring the eyes of the person, it generates artifacts in the suit (row 3). In contrast, SeeSR produces sharper and semantically more accurate results, such as the edges on the wall, the fur of camel hump, the neat suit, and the vivid eyes. More visual examples can be found in the Fig. 7.

**User Study.** To further validate the effectiveness of our method, we conduct separate user studies on synthetic data and real data. On synthetic data, inspired by SR3 [43], participants were presented with an LR image placed between two HR images each time: one is the GT and another is the Real-ISR output by one model. They were asked to determine *'Which HR image better corresponds to the LR image?'* When making decision, participants considered two factors: the perceptual quality of the HR image and its semantic similarity to the LR image. Then the confusion rates can be calculated, which indicate the participants' preference to the GT or the Real-ISR output. On real data, participants were presented with LR images alongside two

Table 3. The comparison of semantic restoration performance among different Real-ISR methods.

| Metrics | GT | Zoomed LR | BSRGAN | Real-ESRGAN | LDL | DASR | FeMaSR | LDM | StableSR | ResShift | PASD | DiffBIR | SeeSR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Panoptic Segmentation (PQ) | 52.5 | 9 | 16.2 | 19.4 | 17.8 | 15.5 | 15.6 | 18.7 | 26.8 | 21.4 | 23.7 | 27.2 | **29.6** |
| Object Detection (AP) | 49.1 | 5 | 10.5 | 13.1 | 11.9 | 9.9 | 10.1 | 11.4 | 18.3 | 14.3 | 15.5 | 18.9 | **21.1** |
| Instance Segmentation (AP) | 43.8 | 4 | 9.2 | 11.4 | 10.3 | 8.6 | 8.8 | 9.9 | 16.2 | 12.4 | 13.4 | 16.5 | **18.4** |
| Semantic Segmentation (mIOU) | 62.0 | 12 | 21.7 | 26.0 | 24.2 | 20.5 | 20.4 | 25.3 | 34.5 | 30.4 | 33.3 | 37.7 | **40.8** |



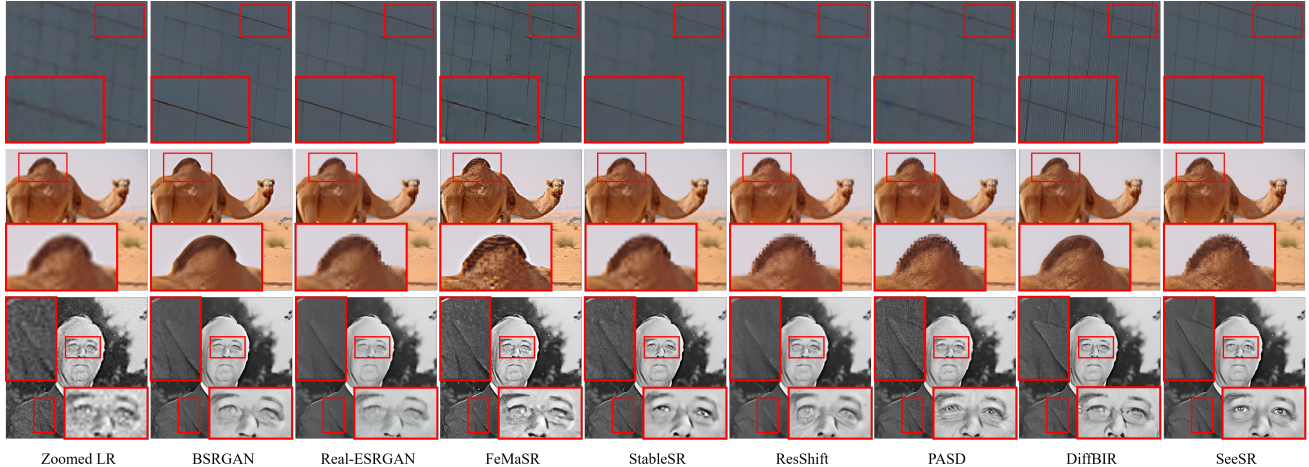Zoomed LR   BSRGAN   Real-ESRGAN   FeMaSR   StableSR   ResShift   PASD   DiffBIR   SeeSR

Figure 5. Qualitative comparisons of different methods on real-world examples. Please zoom in for a better view.

Table 4. Results of user study on synthetic and real-world data.

| Methods | Confusion rates on synthetic data | Best rates on real-world data |
|---|---|---|
| Real-ESRGAN | 5.4% | 0% |
| StableSR | 5.4% | 14.3% |
| ResShift | 3.6% | 0% |
| PASD | 10.7% | 13.4% |
| DiffBIR | 12.5% | 16.1% |
| SeeSR | **36.6%** | **56.2%** |

Table 5. The Real-ISR performance of our SeeSR model with and without LRE on *DIV2K-Val* and *DrealSR* [55] benchmarks.

| Metrics | *DIV2K-Val* | | *DrealSR* | |
|---|---|---|---|---|
| | w/o LRE | w/ LRE | w/o LRE | w/ LRE |
| PSNR ↑ | 20.58 | 21.04 | 26.55 | 27.90 |
| LPIPS ↓ | 0.3942 | 0.3876 | 0.3952 | 0.3299 |
| FID ↓ | 32.53 | 32.79 | 158.04 | 151.88 |
| CLIPIQA ↑ | 0.7314 | 0.6834 | 0.7248 | 0.6708 |

Real-ISR outputs from two models, and they were asked to answer *'Which image is the best SR result of the LR image?'* In this experiment, best rates were calculated, which represent the probability of the model being selected.

We invited 20 participants to test six representative methods (Real-ESRGAN, StableSR, PASD, DiffBIR, ResShift and SeeSR). There are 16 synthetic test sets and 16 real-world test sets. The synthetic data are randomly sampled from *DIV2K-Val*, and the real-world data are randomly sampled from *RealLR200*. Each of the 20 participants was asked to make 112 selections ($16 \times 6 + 16$). As shown in Table 4, our SeeSR significantly outperforms others in terms of selection rate on both synthetic and real data. In the user study on synthetic data, the SR results of all models cannot compete with the GT, while our SeeSR achieve a confusion rate of 36.6, which is three times higher than the second-ranked method. This implies that there is still enough room to improve for the Real-ISR methods. In the user study on real-world data where there is no GT, our method achieves a best selection rate of 56.2%, approximately 3.5 times higher than the second-ranked method.

### 4.3. Semantics Preservation Test

To further validate our model's ability to preserve semantic fidelity, we conduct detection and segmentation tasks on the Real-ISR output images. We resize the original images from COCO-Val (5K images) [34] to $512 \times 512$ as GT, and then degrade them to generate LR images as in training. We employ OpenSeeD [61] trained on COCO as the detector and segmentor since it is a strong transformer-based unified model for segmentation and detection tasks. As shown in Table 3, compared to Zoomed LR, SeeSR achieves a remarkable $3 \sim 4$ times improvement in all four tasks, surpassing all existing Real-ISR methods and showcasing its strong semantics preservation capability.

### 4.4. Ablation Study

We first discuss the effectiveness of the proposed LRE strategy. Then, we discuss the effectiveness of the proposed DAPE module, including its tagging capability and the roles of hard and soft prompts.

**Effectiveness of LRE.** We first show the Real-ISR perfor-

Table 6. Comparison between RAM and DAPE on degraded images of *COCO-val* benchmark [34] for the tagging task.

| | OP ↑ | OR ↑ | AP ↑ |
|---|---|---|---|
| RAM [71] | 0.7929 | 0.3711 | 52.3 |
| DAPE | **0.8940** | **0.3751** | **63.0** |

mance of our SeeSR model on the *DIV2K-Val* and *DrealSR* datasets with and without the LRE strategy. The results are shown in Table 5. One can see that the LRE strategy improves the reference-based metrics, including both fidelity and perception based ones, while it weakens the non-reference metrics such as CLIPIQA. This is because the LRE strategy reduces the model's tendency to generate additional (but maybe unfaithful) textures by narrowing the gap between training and testing (see discussions in Section 3.4). Such an over-generation ability can be favorable by metrics like CLIPIQA, but they will introduce visually unpleasant artifacts, as shown in Fig. 3.

**Tagging Performance of DAPE.** In Table 6, we present the tagging performance of our DAPE module on the degraded images of *COCO-val* benchmark [34] based on three metrics: overall precision (OP), overall recall (OR), and average precision (AP). AP is the averaged precision calculated on different recall rates, which is similar to the detection metric. OP and OR are defined as:

$$\text{OP} = \frac{\sum_i N_i^t}{\sum_i N_i^p}, \quad \text{OR} = \frac{\sum_i N_i^t}{\sum_i N_i^g}, \tag{3}$$

where $C$ is the number of classes, $N_i^p$ is the number of images predicted for label $i$, $N_i^t$ is the number of images correctly predicted for label $i$, and $N_i^g$ is the number of ground truth images for label $i$.

We evaluate RAM [71] and DAPE with the default threshold. DAPE surpasses RAM in terms of OP and AP by 0.1 and 10.7, respectively. It also maintains superiority in OR, indicating that DAPE achieves significant improvements in tagging accuracy for degraded images. This improvement assists the T2I model in generating semantically accurate details when performing the Real-ISR task.

**Effectiveness of DAPE and Hard/Soft Prompts for Real-ISR.** DAPE improves the model's tagging performance on degraded images and consequently enhances the Real-ISR capability. To investigate the effectiveness of DAPE and the roles of its hard/soft prompts, we conducted the following four experiments in Real-ISR tasks.

1. We retrain SeeSR by removing the DAPE and RCA modules, which can be considered as applying Control-Net [66] directly to the Real-ISR task.
2. We replace DAPE with RAM [71] and retrain the model.
3. During the inference of SeeSR, we provide only the hard prompts (*i.e.*, the tag) generated by DAPE to the text encoder of the T2I model.
4. During the inference of SeeSR, we provide only the soft prompts (*i.e.*, the representation embedding features) generated by DAPE to the T2I model.

The results of the four experiments are shown in Table 7. Moreover, the visual comparisons are shown in Fig. 6. From Table 7 and Fig. 6, we can have the following conclusions. First, directly applying ControlNet to the Real-ISR task cannot achieve satisfactory results. Second, replacing DAPE with the original RAM would lead to a decrease in all perceptual metrics (*e.g.*, LPIPS and CLIPIQA). The semantics of the image content may also be changed (see Fig. 6). This is because the original RAM may generate inaccurate prompts (*e.g.*, the tag 'broccoli') from the degraded image. Third, the soft prompts work better in improving the numerical indices than the hard prompts, as well as sharper images. However, without hard prompts, the image semantics can be damaged, as can be seen from the lemons in Exp. (4) of Fig. 6. Finally, with both the hard and soft prompts in DAPE, perceptually realistic and semantically correct Real-ISR outputs can be produced.

### 4.5. Complexity Analysis

Table 8 compares the number of parameters of different Real-ISR models and their inference time to synthesize a $512 \times 512$ image from $128 \times 128$ input. All tests are conducted on one NVIDIA Tesla 32G-V100 GPU. We can have the following observations.

First, the GAN-based methods Real-ESRGAN and Fe-MaSR have much less model parameters and much faster inference speed than DM-based methods. Second, among the DM-based models, LDM and ResShift are much smaller than others because they employ relatively lightweight diffusion models. ResShift runs faster than LDM because it samples only 15 steps while LDM samples 200 steps. Thrid, StableSR, PASD, DiffBIR and our SeeSR are all based on the pre-trained T2I model. SeeSR has more parameters because it has a DAPE module (about 300M) finetuned from the RAM model. In terms of inference speed, PASD, Diff-BIR and SeeSR are comparable, while StableSR is the slowest one because it samples 200 steps.

### 4.6. More Visualization Comparisons

We provide additional qualitative comparisons on real-world images. As shown in Fig. 7, SeeSR can generate sharper edges (case 2) and semantically faithful details (the window railing in case 1, the teeth in case 3, and the vein textures in case 4). Other methods are either blurry or produce unpleasant artifacts.

## 5. Conclusion

We proposed SeeSR, a Real-ISR method that utilizes semantic prompts to enhance the generative capability of pre-

9

Hard prompt from RAM:
'broccoli, garbanzo, green, potato, vegetable'

Hard prompt from DAPE:
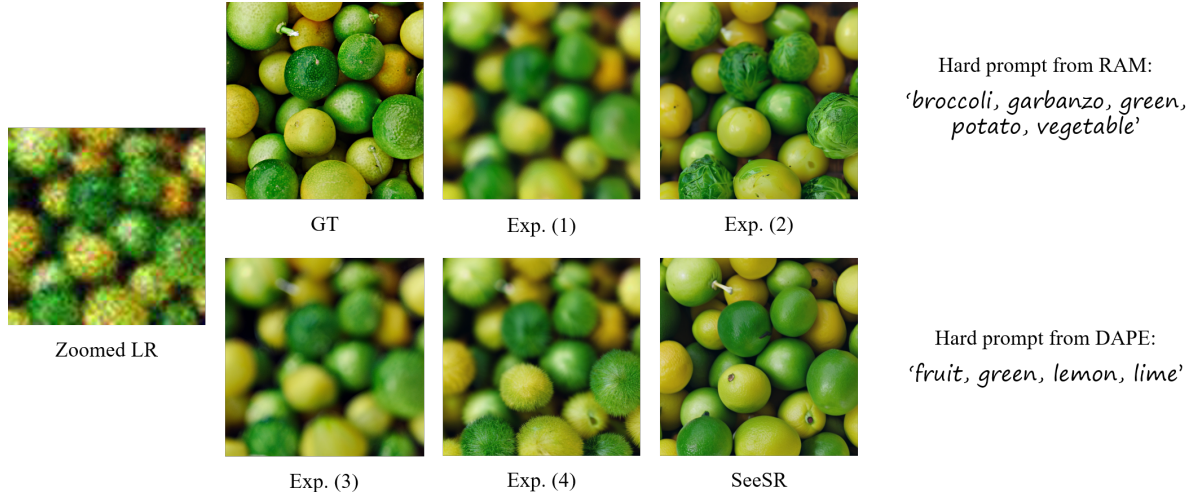'fruit, green, lemon, lime'

Figure 6. Visual comparison for the ablation study on DAPE. Exp. (1) directly applies ControlNet to perform Real-ISR, leading to blurry results. Exp. (2) replaces DAPE with RAM for generating prompts, which can produce sharper but semantically incorrect details. Exp. (3) applies hard prompts only and generates blurry results. Exp. (4) applies soft prompts only and exhibits semantic errors in details generation. With both hard and soft prompts in DAPE, SeeSR produces clear and semantically correct outputs.

Table 7. Ablation studies of DAPE on *DIV2K-Val* and *DrealSR* [55] benchmarks for the Real-ISR task.

| Exp | | (1) | (2) | (3) | (4) | SeeSR |
|---|---|---|---|---|---|---|
| Prompt Extractor | RAM [71] | ✗ | ✓ | ✗ | ✗ | ✗ |
| | DAPE | ✗ | ✗ | ✓ | ✓ | ✓ |
| Prompt Format | Hard Prompt | ✗ | ✓ | ✓ | ✗ | ✓ |
| | Soft Prompt | ✗ | ✓ | ✗ | ✓ | ✓ |
| *DIV2K-Val* | PSNR ↑ | 20.96 | 21.15 | 20.91 | 21.19 | 21.04 |
| | LPIPS ↓ | 0.4236 | 0.4156 | 0.4289 | 0.3859 | 0.3876 |
| | FID ↓ | 37.35 | 46.34 | 38.92 | 38.77 | 32.79 |
| | CLIPIQA ↑ | 0.6343 | 0.6097 | 0.6471 | 0.6751 | 0.6834 |
| *DrealSR* | PSNR ↑ | 27.64 | 27.31 | 27.45 | 28.14 | 27.90 |
| | LPIPS ↓ | 0.3130 | 0.3272 | 0.3285 | 0.3174 | 0.3299 |
| | FID ↓ | 176.26 | 161.69 | 164.57 | 157.63 | 151.88 |
| | CLIPIQA ↑ | 0.5693 | 0.6436 | 0.6410 | 0.6431 | 0.6708 |

Table 8. Complexity comparison between different methods. All the tests are conducted on one NVIDIA Tesla 32G-V100 GPU to synthesize $512 \times 512$ images from $128 \times 128$ inputs.

| Methods | Params | Inference Time-steps | Inference Time |
|---|---|---|---|
| Real-ESRGAN [51] | 16.7M | 1 | 0.09s |
| FeMaSR [5] | 28.3M | 1 | 0.12s |
| LDM [40] | 169.0M | 200 | 5.21s |
| StableSR [48] | 1409.1M | 200 | 18.70s |
| ResShift [60] | 173.9M | 15 | 1.12s |
| PASD [59] | 1900.4M | 20 | 6.07s |
| DiffBIR [35] | 1716.7M | 50 | 5.85s |
| SeeSR | 2283.7M | 50 | 7.24s |

trained T2I diffusion models. Through exploring the im-

pact of different styles of text prompts on the generated results, we found that the image tags can greatly enhance the local perception ability of the T2I model. However, the tags are susceptible to complex image degradation, and they are influenced by manually set thresholds. Therefore, we proposed DAPE, which minimizes the influence of image degradation on semantic prompts and simultaneously outputs soft and hard semantic prompts to guide the diffusion process in image super-resolution. Furthermore, to address the adverse effects of training-test inconsistency in diffusion models, we proposed a simple yet effective LRE strategy, which embeds LR latent at the starting point of diffusion process, avoiding the generation of artifacts in smooth areas. Our work made a step towards better leveraging generative priors to synthesize semantically correct Real-ISR images, as demonstrated in our extensive experiments.
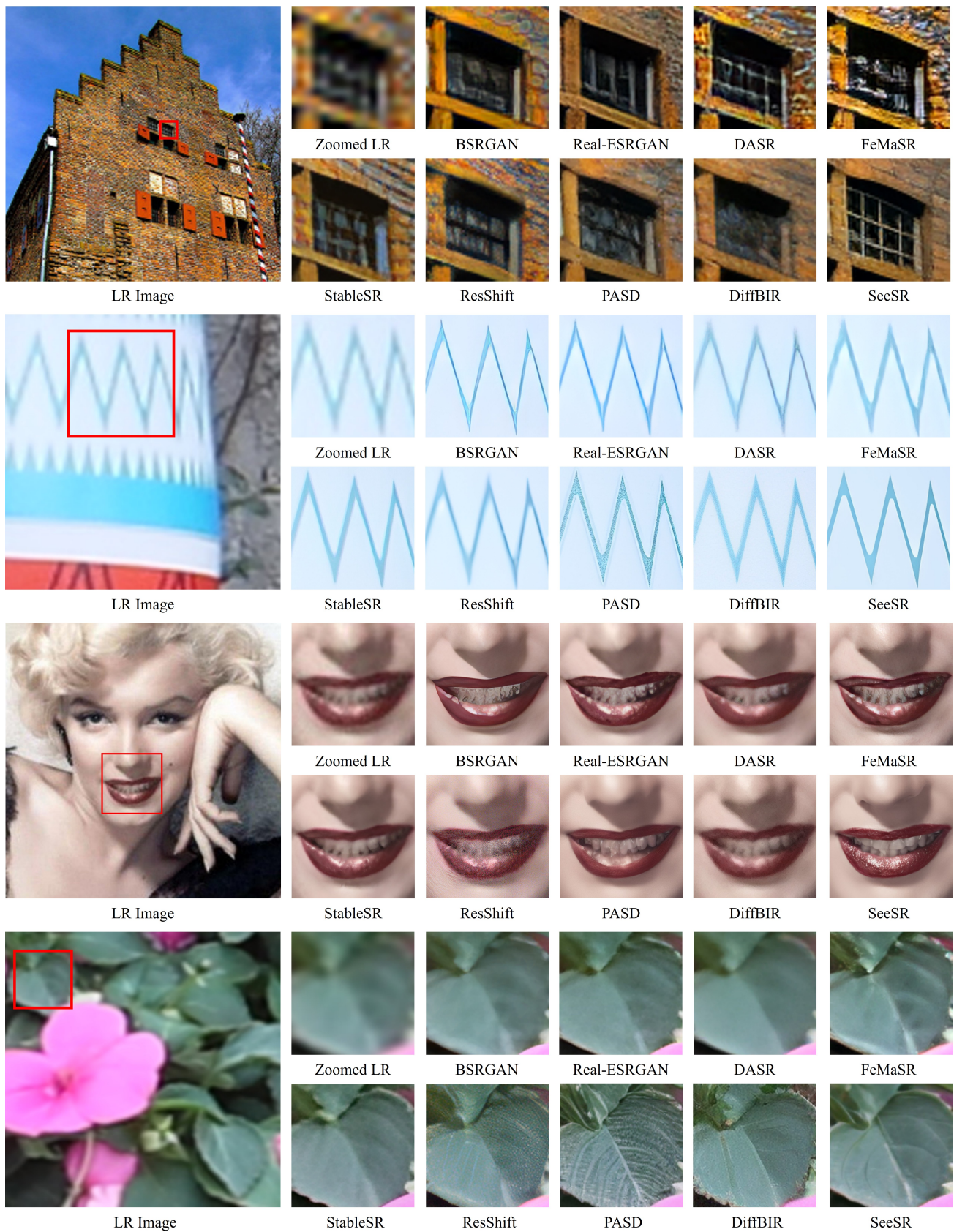
Figure 7. Qualitative comparisons of different methods on real-world examples. Please zoom in for a better view.

# References

[1] Stability.ai. https://stability.ai/stable-diffusion. 2

[2] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 126–135, 2017. 6

[3] Jianrui Cai, Hui Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3086–3095, 2019. 1, 6

[4] Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Investigating tradeoffs in real-world video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5962–5971, 2022. 6

[5] Chaofeng Chen, Xinyu Shi, Yipeng Qin, Xiaoming Li, Xiaoguang Han, Tao Yang, and Shihui Guo. Real-world blind super-resolution via feature matching with implicit high-resolution priors. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 1329–1338, 2022. 6, 10

[6] Du Chen, Jie Liang, Xindong Zhang, Ming Liu, Hui Zeng, and Lei Zhang. Human guided ground-truth generation for realistic image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14082–14091, 2023. 1

[7] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12299–12310, 2021. 1, 2

[8] Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22367–22377, 2023. 1

[9] Zheng Chen, Yulun Zhang, Jinjin Gu, Linghe Kong, Xiaokang Yang, and Fisher Yu. Dual aggregation transformer for image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12312–12321, 2023.

[10] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11065–11074, 2019. 2

[11] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 2

[12] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying structure and texture similarity. *IEEE transactions on pattern analysis and machine intelligence*, 44(5):2567–2581, 2020. 6

[13] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part IV 13*, pages 184–199. Springer, 2014. 1, 2

[14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. 1, 2

[15] Shuhang Gu, Andreas Lugmayr, Martin Danelljan, Manuel Fritsche, Julien Lamour, and Radu Timofte. Div8k: Diverse 8k resolution image dataset. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 3512–3516. IEEE, 2019. 6

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2, 3, 4

[17] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781*, 2019. 2

[18] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 3

[19] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6

[20] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 4

[21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2

[22] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 6

[23] Christopher Jarzynski. Equilibrium free-energy differences from nonequilibrium measurements: A master-equation approach. *Physical Review E*, 56(5):5018, 1997. 2

[24] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 6

[25] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022. 1, 2

[26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6

[27] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 3

[28] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1833–1844, 2021. 1, 2

[29] Jie Liang, Hui Zeng, and Lei Zhang. Details or artifacts: A locally discriminative learning approach to realistic image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5657–5666, 2022. 1, 2, 6

[30] Jie Liang, Hui Zeng, and Lei Zhang. Efficient and degradation-adaptive network for real-world image super-resolution. In *European Conference on Computer Vision*, pages 574–591. Springer, 2022. 1, 6

[31] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017. 2

[32] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 300–309, 2023. 2

[33] Shanchuan Lin, Bingchen Liu, Jiashi Li, and Xiao Yang. Common diffusion noise schedules and sample steps are flawed. *arXiv preprint arXiv:2305.08891*, 2023. 5

[34] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. 6, 8, 9

[35] Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Ben Fei, Bo Dai, Wanli Ouyang, Yu Qiao, and Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior. *arXiv preprint arXiv:2308.15070*, 2023. 2, 5, 6, 7, 10

[36] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2

[37] Radford M Neal. Annealed importance sampling. *Statistics and computing*, 11:125–139, 2001. 2

[38] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 6

[39] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 2

[40] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 4, 5, 6, 7, 10

[41] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 2

[42] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1, 2

[43] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4713–4726, 2022. 2, 7

[44] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 2

[45] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 2

[46] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 2

[47] Radu Timofte, Eirikur Agustsson, Luc Van Gool, Ming-Hsuan Yang, and Lei Zhang. Ntire 2017 challenge on single image super-resolution: Methods and results. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 114–125, 2017. 6

[48] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. *arXiv preprint arXiv:2305.07015*, 2023. 2, 5, 6, 7, 10

[49] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 606–615, 2018. 6

[50] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European conference on computer vision (ECCV) workshops*, pages 0–0, 2018. 1

[51] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1905–1914, 2021. 1, 2, 5, 6, 10

13

[52] Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. *arXiv preprint arXiv:2212.00490*, 2022. 1, 2

[53] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6

[54] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023. 2

[55] Pengxu Wei, Ziwei Xie, Hannan Lu, Zongyuan Zhan, Qixiang Ye, Wangmeng Zuo, and Liang Lin. Component divide-and-conquer for real-world image super-resolution. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VIII 16*, pages 101–117. Springer, 2020. 1, 6, 8, 10

[56] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7623–7633, 2023. 2

[57] Liangbin Xie, Xintao Wang, Xiangyu Chen, Gen Li, Ying Shan, Jiantao Zhou, and Chao Dong. Desra: Detect and delete the artifacts of gan-based real-world super-resolution models. 2023. 1, 2

[58] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1191–1200, 2022. 6

[59] Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Pixel-aware stable diffusion for realistic image super-resolution and personalized stylization. *arXiv preprint arXiv:2308.14469*, 2023. 2, 3, 4, 5, 6, 7, 10

[60] Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. *arXiv preprint arXiv:2307.12348*, 2023. 6, 7, 10

[61] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianfeng Gao, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. In *CVPR*, 2023. 8

[62] Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. *arXiv preprint arXiv:2305.15347*, 2023. 3

[63] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE transactions on image processing*, 26(7):3142–3155, 2017. 6

[64] Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4791–4800, 2021. 1, 2, 5, 6

[65] Lin Zhang, Lei Zhang, and Alan C Bovik. A feature-enriched completely blind image quality evaluator. *IEEE Transactions on Image Processing*, 24(8):2579–2591, 2015. 6

[66] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. 2, 4, 9

[67] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 6

[68] Xindong Zhang, Hui Zeng, Shi Guo, and Lei Zhang. Efficient long-range attention network for image super-resolution. In *European Conference on Computer Vision*, pages 649–667. Springer, 2022. 1, 2

[69] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European conference on computer vision (ECCV)*, pages 286–301, 2018. 1

[70] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2472–2481, 2018. 2

[71] Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. Recognize anything: A strong image tagging model. *arXiv preprint arXiv:2306.03514*, 2023. 3, 4, 5, 6, 9, 10