# Scaling Up to Excellence:
# Practicing Model Scaling for Photo-Realistic Image Restoration In the Wild
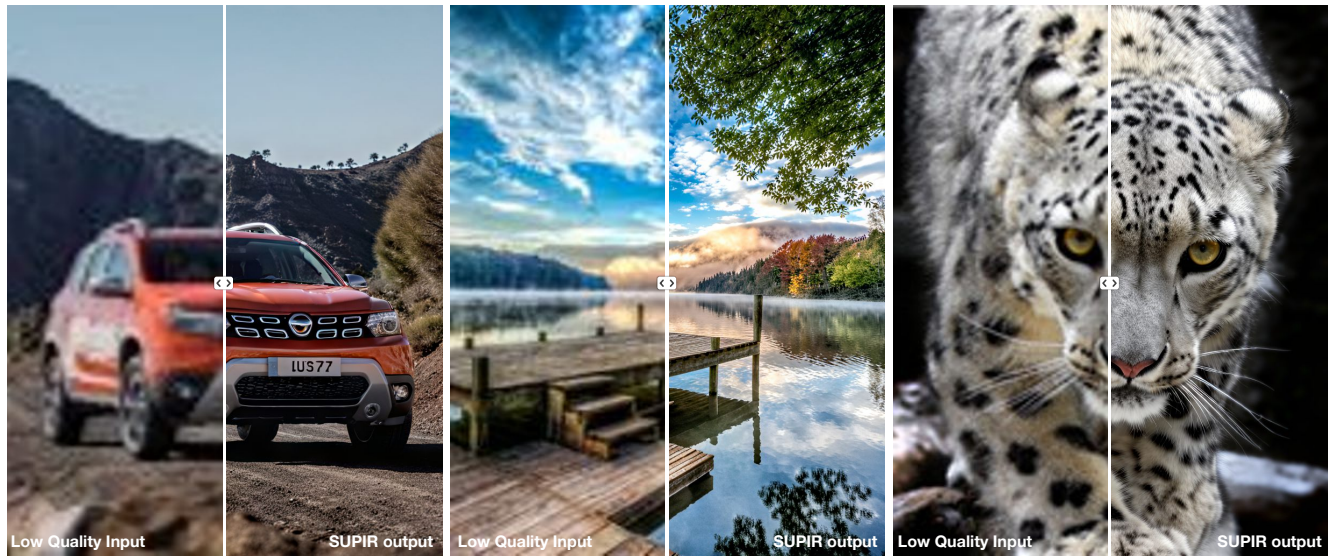
Fanghua Yu[1,*], Jinjin Gu[2,3,*], Zheyuan Li[1], Jinfan Hu[1], Xiangtao Kong[4],
Xintao Wang[5], Jingwen He[2,6], Yu Qiao[2], Chao Dong[1,2,†]

[1]Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences  [2]Shanghai AI Laboratory
[3]University of Sydney  [4]The Hong Kong Polytechnic University  [5]ARC Lab, Tencent PCG  [5]The Chinese University of Hong Kong

Project Page: https://supir.xpixel.group

**(a) Real-World Image Restoration Results**



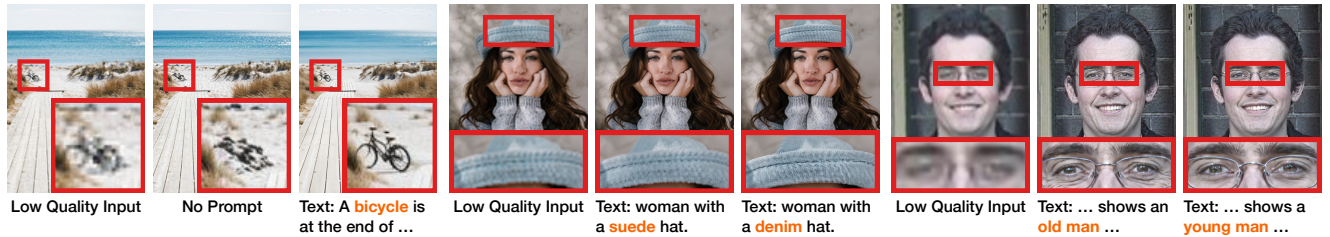**(b) Controllable Image Restoration with Textual Prompts**



Figure 1. Our SUPIR model demonstrates remarkable restoration effects on real-world low-quality images, as illustrated in (a). Additionally, SUPIR features targeted restoration capability driven by textual prompts. For instance, it can specify the restoration of blurry objects in the distance (case 1), define the material texture of objects (case 2), and adjust restoration based on high-level semantics (case 3).

## Abstract

We introduce SUPIR (Scaling-UP Image Restoration), a groundbreaking image restoration method that harnesses generative prior and the power of model scaling up. Leveraging multi-modal techniques and advanced generative prior, SUPIR marks a significant advance in intelligent and realistic image restoration. As a pivotal catalyst within SUPIR, model scaling dramatically enhances its capabilities and demonstrates new potential for image restoration. We collect a dataset comprising 20 million high-resolution, high-quality images for model training, each enriched with descriptive text annotations. SUPIR provides the capability to restore images guided by textual prompts, broadening its application scope and potential. Moreover,

---

* Contribute Equally.    † Corresponding Author.

*we introduce negative-quality prompts to further improve perceptual quality. We also develop a restoration-guided sampling method to suppress the fidelity issue encountered in generative-based restoration. Experiments demonstrate SUPIR's exceptional restoration effects and its novel capacity to manipulate restoration through textual prompts.*

## 1. Introduction

With the development of image restoration (IR), expectations for the perceptual effects and intelligence of IR results have significantly increased. IR methods based on generative priors [42, 49, 67, 82] leverage powerful pre-trained generative models to introduce high-quality generation and prior knowledge into IR, bringing significant progress in these aspects. Continuously enhancing the capabilities of the generative prior is key to achieving more intelligent IR results, with model scaling being a crucial and effective approach. There are many tasks that have obtained astonishing improvements from scaling, such as SAM [44] and large language models [7, 73, 74]. This further motivates our effort to build large-scale, intelligent IR models capable of producing ultra-high-quality images. However, due to engineering constraints such as computing resources, model architecture, training data, and the cooperation of generative models and IR, scaling up IR models is challenging.

In this work, we introduce SUPIR (Scaling-UP IR), the largest-ever IR method, aimed at exploring greater potential in visual effects and intelligence. Specifically, SUPIR employs StableDiffusion-XL (SDXL) [63] as a powerful generative prior, which contains 2.6 billion parameters. To effectively apply this model, we design and train a adaptor with more than 600 million parameters. Moreover, we have collected over 20 million high-quality, high-resolution images to fully realize the potential offered by model scaling. Each image is accompanied by detailed descriptive text, enabling the control of restoration through textual prompts. We also utilize a 13-billion-parameter multi-modal language model to provide image content prompts, greatly improving the accuracy and intelligence of our method. The proposed SUPIR model demonstrates exceptional performance in a variety of IR tasks, achieving the best visual quality, especially in complex and challenging real-world scenarios. Additionally, the model offers flexible control over the restoration process through textual prompts, vastly broadening the possibility of IR. Fig. 1 illustrates the effects by our model, showcasing its superior performance.

Our work goes far beyond simply scaling. While pursuing an increase in model scale, we face a series of complex challenges. First, when applying SDXL for IR, existing Adaptor designs either too simple to meet the complex requirements of IR [59] or are too large to train together with SDXL [95]. To solve this problem, we trim the ControlNet and designed a new connector called ZeroSFT to

work with the pre-trained SDXL, aiming to efficiently implement the IR task while reducing computing costs. In order to enhance the model's ability to accurately interpret the content of low-quality images, we fine-tune the image encoder to improve its robustness to variations in image degradation. These measures make scaling the model feasible and effective, and greatly improve its stability. Second, we amass a collection of 20 million high-quality, high-resolution images with descriptive text annotations, providing a solid foundation for the model's training. We adopt a counter-intuitive strategy by incorporating poor quality, negative samples into training. In this way, we can use negative quality prompts to further improve visual effects. Our results show that this strategy significantly improves image quality compared to using only high-quality positive samples. Finally, powerful generative prior is a double-edged sword. Uncontrolled generation may reduce restoration fidelity, making IR no longer faithful to the input image. To mitigate this low-fidelity issue, we propose a novel restoration-guided sampling method. All these strategies, coupled with efficient engineering implementation, are key to enabling the scaling up of SUPIR, pushing the boundaries of advanced IR. This comprehensive approach, encompassing everything from model architecture to data collection, positions SUPIR at the forefront of image restoration technology, setting a new benchmark for future advancements.

## 2. Related Work

**Image Restoration.** The goal of IR is to convert degraded images into high-quality degradation-free versions [22, 26, 89, 91, 98, 99]. In the early stage, researchers independently explored different types of image degradation, such as super-resolution (SR) [13, 19, 20], denoising [11, 90, 92], and deblurring [14, 60, 72]. However, these methods are often based on specific degradation assumptions [25, 50, 58] and therefore lack generalization ability to other degradations [29, 53, 97]. Over time, the need for blind restoration methods that are not based on specific degradation assumptions has grown [5, 10, 34, 35, 46–48, 78, 94]. In this trend, some methods [81, 93] approximate synthesize real-world degradation by more complex degradation models, and are well-known for handling multiple degradation with a single model. Recent research, such as DiffBIR [49], unifies different restoration problems into a single model. In this paper, we adopt a similar setting to DiffBIR and use a single model to achieve effective processing of various severe degradations.

**Generative Prior.** Generative priors are adept at capturing the inherent structures of the image, enabling the generation of images that follow natural image distribution. The emergence of GANs [23, 39, 40, 64] has underscored the significance of generative priors in IR. Various approaches employ these priors, including GAN inversion [2, 4, 27, 57,
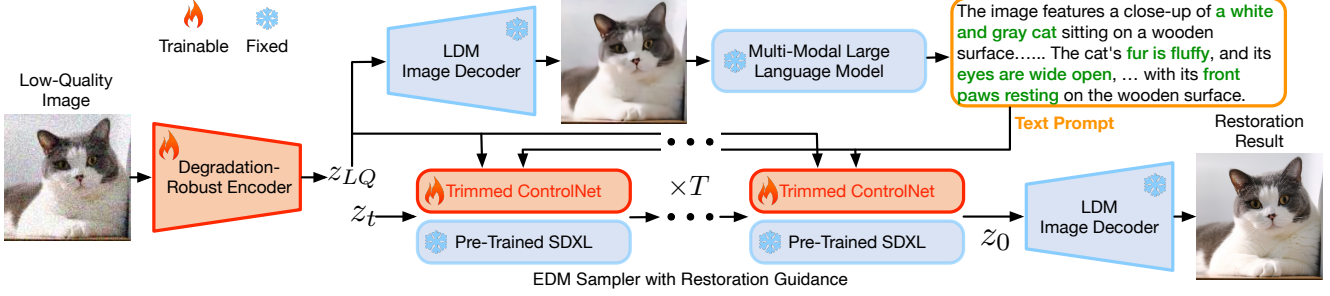
Figure 2. This figure briefly shows the workflow of the proposed SUPIR model.

62], GAN encoders [9, 103], or using GAN as the core module for IR [80, 87]. Beyond GANs, other generative models can also serve as priors [10, 36, 55, 75, 100–102]. Our work primarily focuses on generative priors derived from diffusion models [31, 61, 65, 67, 70, 71], which excel in controllable generation [15, 18, 32, 59, 95] and model scaling [63, 66, 68]. Diffusion models have also been effectively used as generative priors in IR [42, 49, 67, 77, 82]. However, these diffusion-based IR methods' performance is constrained by the scale of the used generative models, posing challenges in further enhancing their effectiveness.

**Model Scaling** is an important means to further improve the capabilities of deep-learning models. The most typical examples include the scaling of language models [7, 73, 74], text-to-image generation models [12, 37, 63, 67, 68, 85], and image segmentation models [44]. The scale and complexity of these models have increased dramatically, with billions or even hundreds of billions of parameters, but these parameters also lead to extraordinary performance improvements, demonstrating the potential of model scaling [38]. However, scaling up is a systematic problem, involving model design, data collection, computing resources, and other limitations. Many other tasks have not yet been able to enjoy the substantial performance improvements brought by scaling up. IR is one of them.

## 3. Method

An overview of the proposed SUPIR method is shown in Fig. 2. We introduce our method from three aspects: Sec. 3.1 introduces our network designs and training method; Sec. 3.2 introduces the collection of training data and the introduction of text modality; and Sec. 3.3 introduces the diffusion sampling method for image restoration.

### 3.1. Model Scaling Up

**Generative Prior.** There are not many choices for the large-scale generative models. The only ones to consider are Imagen [68], IF [16], and SDXL [63]. Our selection settled on SDXL for the following reasons. Imagen and IF prioritize text-to-image generation and rely on a hierarchical approach. They first generate small-resolution images and then hierarchically upsample them. SDXL directly

generates a high-resolution image without hierarchical design, which is more aligned with our objectives, as it utilizes its parameters effectively for image quality improvement rather than text interpretation. Additionally, SDXL employs a *Base-Refine* strategy. In the *Base* model, diverse but lower-quality images are generated. Subsequently, the *Refine* model enhances the perceptual quality of these images. Compared to the *Base* model, the *Refine* model uses training images with significantly higher quality but less diverse. Considering our strategy to train with an extensive dataset of high-quality images, the two-phase design of SDXL becomes superfluous for our needs. We opt for the *Base* model, which has a greater number of parameters, making it an ideal backbone for our generative prior.

**Degradation-Robust Encoder.** In SDXL, the diffusion generation process is performed in the latent space. The image is first mapped to the latent space through a pre-trained encoder. To effectively utilize the pre-trained SDXL, our LQ image $x_{LQ}$ should also be mapped to the same latent space. However, since the original encoder has not been trained on LQ images, using it for encoding will affect the model's judgment of LQ image content, and then misunderstand artifacts as image content [49]. To this end, we fine-tune the encoder to make it robust to the degradation by minimizing: $\mathcal{L}_{\mathcal{E}} = \|\mathcal{D}(\mathcal{E}_{dr}(x_{LQ})) - \mathcal{D}(\mathcal{E}_{dr}(x_{GT}))\|_2^2$, where $\mathcal{E}_{dr}$ is the degradation-robust encoder to be fine-tuned, $\mathcal{D}$ is the fixed decoder, $x_{GT}$ is the ground truth.

**Large-Scale Adaptor Design.** Considering the SDXL model as our chosen prior, we need an adaptor that can steer it to restore images according to the provided LQ inputs. The Adaptor is required to identify the content in the LQ image and to finely control the generation at the pixel level. LoRA [32], T2I adaptor [59], and ControlNet [95] are existing diffusion model adaptation methods, but none of them meet our requirements: LoRA limits generation but struggles with LQ image control; T2I lacks capacity for effective LQ image content identification; and ControlNet's direct copy is challenging for the SDXL model scale. To address this issue, we design a new adaptor with two key features, as shown in Fig. 3(a). First, we keep the high-level design of ControlNet but employ network trimming [33] to directly trim some blocks within the trainable copy, achiev-
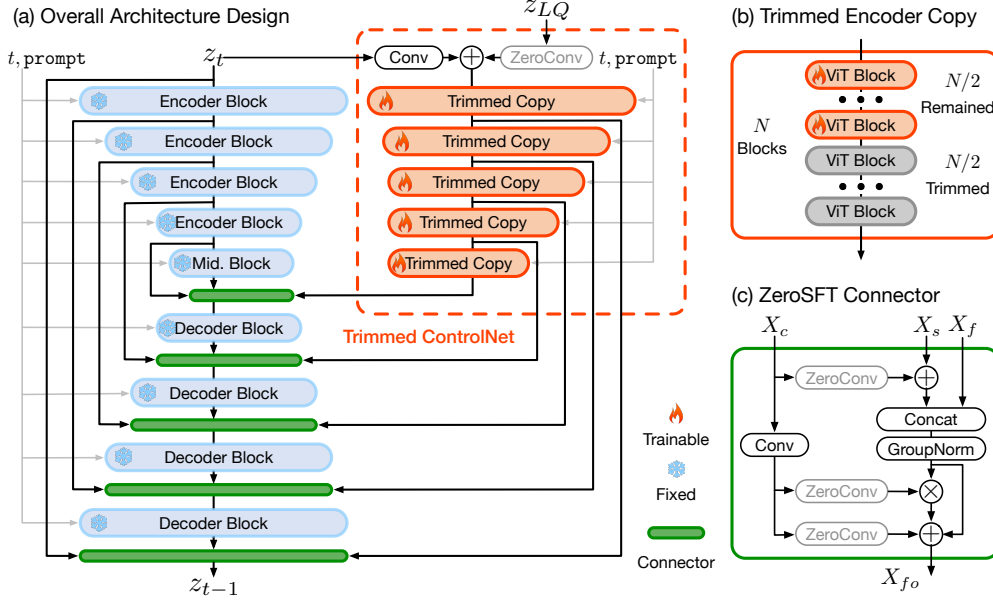
Figure 3. This figure illustrates (a) the overall architecture of the used SDXL and the proposed adaptor, (b) a trimmed trainable copy of the SDXL encoder with reduced ViT blocks for efficiency, and (c) a novel ZeroSFT connector for enhanced control in IR, where $X_f$ and $X_s$ denote the input feature maps from the Decoder and Encoder shortcut, respectively, $X_c$ is the input from the adaptor, and $X_{fo}$ is the output. The model is designed to effectively use the large-scale SDXL as a generative prior.

ing an engineering-feasible implementation. Each block within the encoder module of SDXL is mainly composed of several Vision Transformer (ViT) [21] blocks. We identified two key factors contributing to the effectiveness of ControlNet: large network capacity and efficient initialization of the trainable copy. Notably, even partial trimming of blocks in the trainable copy retains these crucial characteristics in the adaptor. Therefore, we simply trim half of the ViT blocks from each encoder block, as shown in Fig. 3(b). Second, we redesign the connector that links the adaptor to SDXL. While SDXL's generative capacity delivers excellent visual effects, it also renders pixel-level precise control challenging. ControlNet employs zero convolution for generation guidance, but relying solely on residuals is insufficient for the control required by IR. To amplify the influence of LQ guidance, we introduced a ZeroSFT module, as depicted in Fig. 3(c). Building based on zero convolution, ZeroSFT encompasses an additional spatial feature transfer (SFT) [79] operation and group normalization [84].

### 3.2. Scaling Up Training Data

**Image Collection.** The scaling of the model requires a corresponding scaling of the training data [38]. But there is no large-scale high-quality image dataset available for IR yet. Although DIV2K [3] and LSDIR [1] offer high image quality, they are limited in quantity. Larger datasets like ImageNet (IN) [17], LAION-5B[69], and SA-1B [44] contain more images, but their image quality does not meet our high standards. To this end, we collect a new large-scale dataset of high-resolution images, which includes 20 million $1024\times1024$ high-quality, texture-rich, and content-clear images. A comparison on scales of the collected dataset and the existing dataset is shown in Fig. 3. We also included an additional 70K unaligned high-resolution facial images from FFHQ-raw dataset [40] to improve the model's
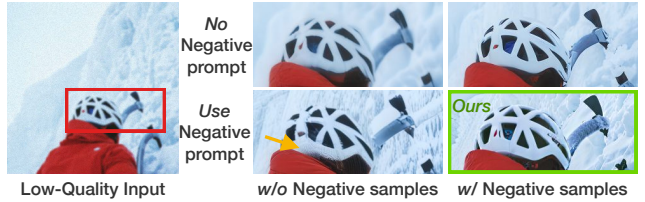


Figure 4. CFG introduces artifacts without negative training samples, hindering visual quality improvement. Adding negative samples allows further quality enhancement through CFG.

face restoration performance. In Fig. 5(a), we show the relative size of our data compared to other well-known datasets.

**Multi-Modality Language Guidance.** Diffusion models are renowned for their ability to generate images based on textual prompts. We believe that textual prompts can also significantly aid IR for the following reasons: (1) Understanding image content is crucial for IR. Existing frameworks often overlook or implicitly handle this understanding [24, 29]. By incorporating textual prompts, we explicitly convey the understanding of LQ images to the IR model, facilitating targeted restoration of missing information. (2) In cases of severe degradation, even the best IR models may struggle to recover completely lost information. In such cases, textual prompts can serve as a control mechanism, enabling targeted completion of missing information based on user preferences. (3) We can also describe the desired image quality through text, further enhancing the perceptual quality of the output. See Fig. 1(b) for some examples. To this end, we make two main modifications. First, we revise the overall framework to incorporate the LLaVA multi-modal large language model [52] into our pipeline, as shown in Fig. 2. LLaVA takes the degradation-robust processed LQ images $x'_{LQ} = \mathcal{D}(\mathcal{E}_{dr}(x_{LQ}))$ as input and explicitly understands the content within the images, out-

putting in the form of textual descriptions. These descriptions are then used as prompts to guide the restoration. This process can be automated during testing, eliminating the need for manual intervention. Secondly, following the approach of PixART [12], we also collect textual annotations for all the training images, to reinforce the role of textual control during the training of out model. These two changes endow SUPIR with the ability to understand image content and to restore images based on textual prompts.

**Negative-Quality Samples and Prompt.** Classifier-free guidance (CFG) [30] provides another way of control by using negative prompts to specify undesired content for the model. We can use this feature to specify the model NOT to produce low-quality images. Specifically, at each step of diffusion, we will make two predictions using positive prompts pos and negative prompts neg, and take the fusion of these two results as the final output $z_{t-1}$:

$$z_{t-1}^{\text{pos}} = \mathcal{H}(z_t, z_{LQ}, \sigma_t, \text{pos}), z_{t-1}^{\text{neg}} = \mathcal{H}(z_t, z_{LQ}, \sigma_t, \text{neg}),$$
$$z_{t-1} = z_{t-1}^{\text{pos}} + \lambda_{\text{cfg}} \times (z_{t-1}^{\text{pos}} - z_{t-1}^{\text{neg}}),$$

where $\mathcal{H}(\cdot)$ is our diffusion model with adaptor, $\sigma_t$ is the variance of the noise at time-step $t$, and $\lambda_{\text{cfg}}$ is a hyper-parameter. In our framework, pos can be the image description with positive words of quality, and neg is the negative words of quality, *e.g.*, "*oil painting, cartoon, blur, dirty, messy, low quality, deformation, low resolution, oversmooth*". Accuracy in predicting both positive and negative directions is crucial for the CFG technique. However, the absence of negative-quality samples and prompts in our training data may lead to a failure of the fine-tuned SUPIR in understanding negative prompts. Therefore, using negative-quality prompts during sampling may introduce artifacts, see Fig. 4 for an example. To address this problem, we used SDXL to generate 100K images corresponding to the negative-quality prompts. We counter-intuitively add these low-quality images to the training data to ensure that negative-quality concept can be learned by the proposed SUPIR model.

### 3.3. Restoration-Guided Sampling

Powerful generative prior is a double-edged sword, as too much generation capacity will in turn affect the fidelity of the recovered image. This highlights the fundamental difference between IR tasks and generation tasks. We need means to limit the generation to ensure that the image recovery is faithful to the LQ image. We modified the EDM sampling method [41] and proposed a restoration-guided sampling method to solve this problem. We hope to selectively guide the prediction results $z_{t-1}$ to be close to the LQ image $z_{LQ}$ in each diffusion step. The specific algorithm is shown in Algorithm 1, where $T$ is the total step number, $\{\sigma_t\}_{t=1}^T$ are the noise variance for $T$ steps, $c$ is the additional text prompt condition. $\tau_r, S_{\text{churn}}, S_{\text{noise}}, S_{\text{min}}, S_{\text{max}}$ are five hyper-parameters, but only $\tau_r$ is related to the restoration

---

**Algorithm 1** Restoration-Guided Sampling.

**Input:** $\mathcal{H}, \{\sigma_t\}_{t=1}^T, z_{LQ}, c$
**Hyper-parameter:** $\tau_r, S_{\text{churn}}, S_{\text{noise}}, S_{\text{min}}, S_{\text{max}}$
1: **sample** $z_T \sim \mathcal{N}(\mathbf{0}, \sigma_T^2 \mathbf{I})$
2: **for** $t \in \{T, \dots, 1\}$ **do**
3:     **sample** $\epsilon_t \sim \mathcal{N}(\mathbf{0}, S_{\text{noise}}^2 \mathbf{I})$
4:     $\gamma_t \leftarrow \begin{cases} \min\left(\frac{S_{\text{churn}}}{N}, \sqrt{2} - 1\right) & \text{if } \sigma_t \in [S_{\text{min}}, S_{\text{max}}] \\ 0 & \text{otherwise} \end{cases}$
5:     $k_t \leftarrow (\sigma_t/\sigma_T)^{\tau_r}, \hat{z}_t \leftarrow z_t + \sqrt{\hat{\sigma}_t^2 - \sigma_t^2}\epsilon_t, \hat{\sigma}_t \leftarrow \sigma_t + \gamma_t \sigma_t$
6:     $\hat{z}_{t-1} \leftarrow \mathcal{H}(\hat{z}_t, z_{LQ}, \hat{\sigma}_t, c)$
7:     $d_t \leftarrow (\hat{z}_t - (\hat{z}_{t-1} + k_t(z_{LQ} - \hat{z}_{t-1})))/\hat{\sigma}_t$
8:     $z_{t-1} \leftarrow \hat{z}_t + (\sigma_{t-1} - \hat{\sigma}_t) d_t$
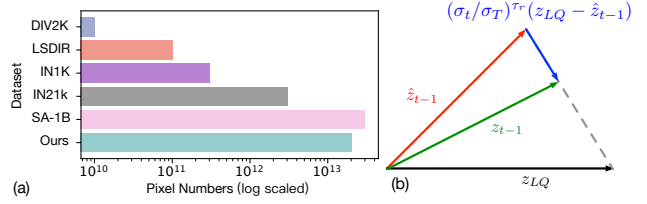9: **end for**

---



Figure 5. (a) We show the relative size of our data compared to other well-known datasets. Compared with SA-1B [44], our dataset has higher quality and more image diversity. (b) We demonstrate our restoration-guided sampling mechanism.

guidance, the others remain unchanged compared to the original EDM method [41]. For better understanding, a simple diagram is shown in Fig. 5(b). We perform weighted interpolation between the predicted output $\hat{z}_{t-1}$ and the LQ latent $z_{LQ}$ as the restoration-guided output $z_{t-1}$. Since the low-frequency information of the image is mainly generated in the early stage of diffusion prediction [67] (where $t$ and $\sigma_t$ are relatively large, and the weight $k = (\sigma_t/\sigma_T)^{\tau_r}$ is also large), the prediction result is closer to $z_{LQ}$ to enhance fidelity. In the later stages of diffusion prediction, mainly high-frequency details are generated. There should not be too many constraints at this time to ensure that detail and texture can be adequately generated. At this time, $t$ and $\sigma_t$ are relatively small, and weight $k$ is also small. Therefore, the predicted results will not be greatly affected Through this method, we can control the generation during the diffusion sampling process to ensure fidelity.

## 4. Experiments

### 4.1. Model Training and Sampling Settings

For training, the overall training data includes 20 million high-quality images with text descriptions, 70K face images and 100K negative-quality samples and the corresponding negative prompts. To enable a larger batch size, we crop them into 512×512 patches during training. We train our model using a synthetic degradation model, following the setting used by Real-ESRGAN [81], the only difference is that we resize the produced LQ images to 512×512 for training. We use the AdamW optimizer [54] with a learning rate of 0.00001. The training process spans 10 days and is
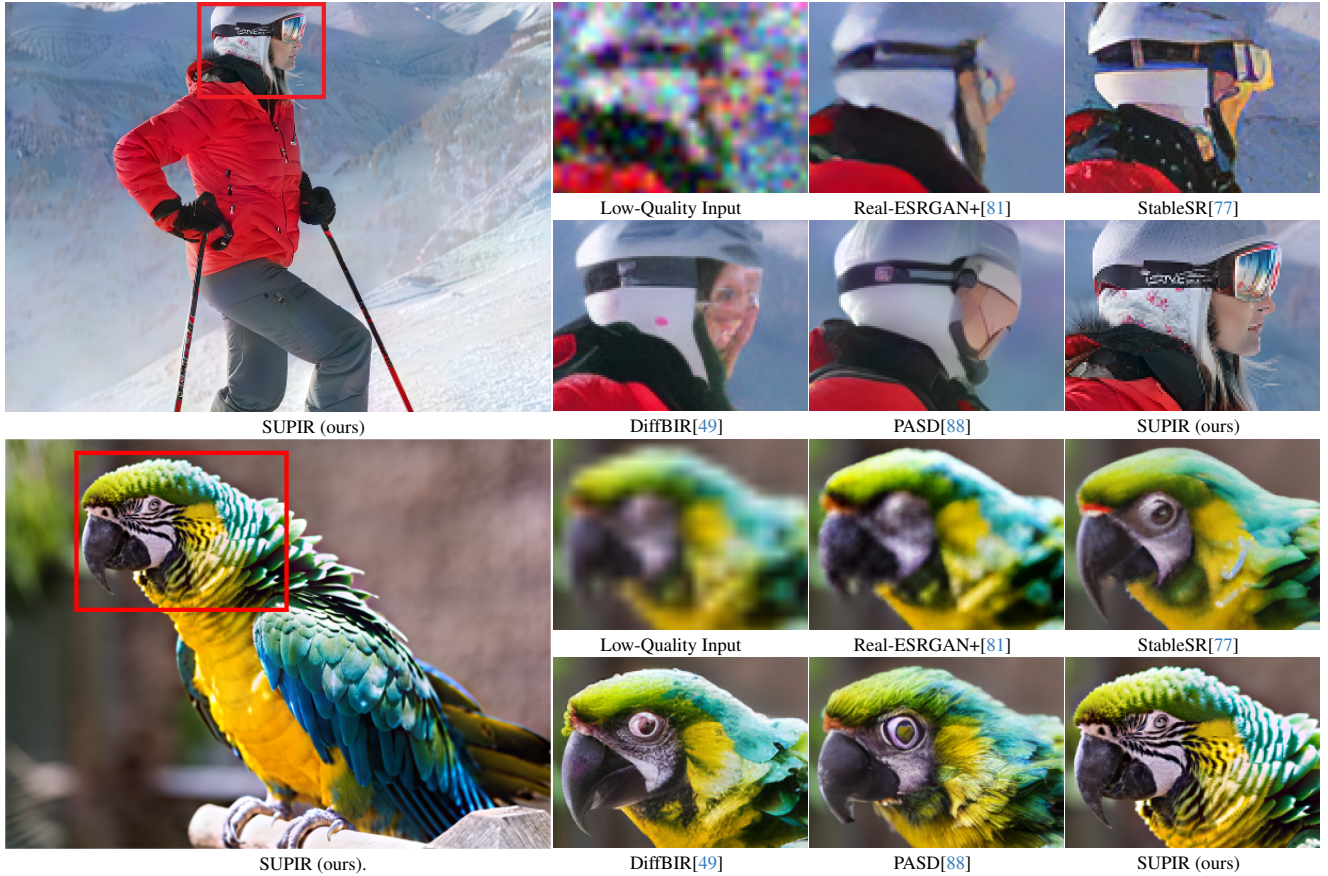
Figure 6. Qualitative comparison with different methods. Our method can accurately restore the texture and details of the corresponding object under challenging degradation. Other methods fail to recover semantically correct details such as broken beaks and irregular faces.
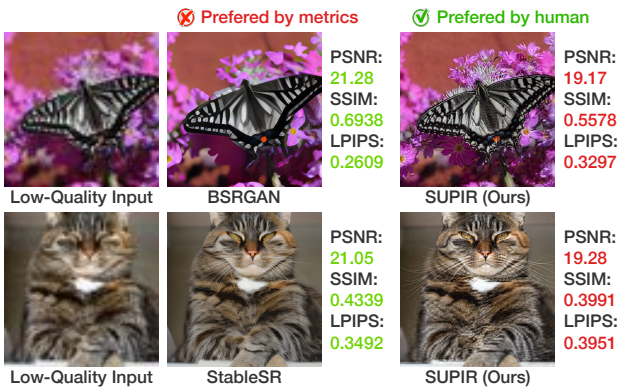


Figure 7. These examples show the misalignment between metric evaluation and human evaluation. SUPIR generates images with high-fidelity textures, but obtains lower metrics.

conducted on 64 Nvidia A6000 GPUs, with a batch size of 256. For testing, the hyper-parameters are $T$=100, $\lambda_{cfg}$=7.5, and $\tau_r = 4$. Our method is able to process images with the size of $1024\times1024$. We resize the short side of the input image to 1024 and crop a $1024\times1024$ sub-image for testing, and then resize it back to the original size after restoration. Unless stated otherwise, prompts will not be provided manually – the processing will be entirely automatic.

## 4.2. Comparison with Existing Methods

Our method can handle a wide range of degradations, and we compare it with the state-of-the-art methods with the same capabilities, including BSRGAN [93], Real-ESRGAN [81], StableSR [77], DiffBIR [49] and PASD [88]. Some of them are constrained to generating images of $512\times512$ size. In our comparison, we crop the test image to meet this requirement and downsample our results to facilitate fair comparisons. We conduct comparisons on both synthetic data and real-world data.

**Synthetic Data.** To synthesize LQ images for testing, we follow previous works [45, 97] and demonstrate our effects on several representative degradations, including both single degradations and complex mixture degradations. Specific details can be found in Tab. 1. We selected the following metrics for quantitative comparison: full-reference metrics PSNR, SSIM, LPIPS [96], and the non-reference metrics ManIQA [86], ClipIQA [76], MUSIQ [43]. It can be seen that our method achieves the best results on all non-reference metrics, which reflects the excellent image quality of our results. At the same time, we also note the disadvantages of our method in full-reference metrics. We present a simple experiment that highlights the limitations of these full-reference metrics, see Fig. 7. It can be seen

6

| Degradation | Method | PSNR | SSIM | LPIPS↓ | ManIQA | ClipIQA | MUSIQ |
|---|---|---|---|---|---|---|---|
| Single:<br>SR (×4) | BSRGAN | 25.06 | 0.6741 | 0.2159 | 0.2214 | 0.6169 | 70.38 |
| | Real-ESRGAN | 24.26 | 0.6657 | 0.2116 | 0.2287 | 0.5884 | 69.51 |
| | StableSR | 22.59 | 0.6019 | 0.2130 | 0.3304 | 0.7520 | 72.94 |
| | DiffBIR | 23.44 | 0.5841 | 0.2337 | 0.2879 | 0.7147 | 71.64 |
| | PASD | 24.90 | 0.6653 | 0.1893 | 0.2607 | 0.6466 | 71.39 |
| | SUPIR (ours) | 22.66 | 0.5763 | 0.2662 | 0.4738 | 0.8049 | 73.83 |
| Single:<br>SR (×8) | BSRGAN | 22.26 | 0.5212 | 0.3523 | 0.2069 | 0.5836 | 67.04 |
| | Real-ESRGAN | 21.79 | 0.5280 | 0.3276 | 0.2051 | 0.5349 | 63.80 |
| | StableSR | 21.27 | 0.4857 | 0.3118 | 0.3039 | 0.7333 | 71.74 |
| | DiffBIR | 21.86 | 0.4957 | 0.3106 | 0.2845 | 0.7080 | 70.26 |
| | PASD | 21.97 | 0.5149 | 0.3034 | 0.2412 | 0.6402 | 70.20 |
| | SUPIR (ours) | 20.68 | 0.4488 | 0.3749 | 0.4687 | 0.8009 | 73.16 |
| Mixture:<br>Blur (σ=2) +<br>SR (×4) | BSRGAN | 24.97 | 0.6572 | 0.2261 | 0.2127 | 0.5984 | 69.44 |
| | Real-ESRGAN | 24.08 | 0.6496 | 0.2208 | 0.2357 | 0.5853 | 69.27 |
| | StableSR | 22.26 | 0.5721 | 0.2301 | 0.3204 | 0.7488 | 72.87 |
| | DiffBIR | 23.28 | 0.5741 | 0.2395 | 0.2829 | 0.7055 | 71.22 |
| | PASD | 24.85 | 0.6560 | 0.1952 | 0.2500 | 0.6335 | 71.07 |
| | SUPIR (ours) | 22.43 | 0.5626 | 0.2771 | 0.4757 | 0.8110 | 73.55 |
| Mixture:<br>SR (×4)+<br>Noise (σ=40) | BSRGAN | 17.74 | 0.3816 | 0.5659 | 0.1006 | 0.4166 | 51.25 |
| | Real-ESRGAN | 21.46 | 0.5220 | 0.4636 | 0.1236 | 0.4536 | 52.23 |
| | StableSR | 20.88 | 0.4174 | 0.4668 | 0.2365 | 0.5833 | 63.54 |
| | DiffBIR | 22.08 | 0.4918 | 0.3738 | 0.2403 | 0.6435 | 65.97 |
| | PASD | 21.79 | 0.4983 | 0.3842 | 0.2590 | 0.5939 | 69.09 |
| | SUPIR (ours) | 20.77 | 0.4571 | 0.3945 | 0.4674 | 0.7840 | 73.35 |
| Mixture:<br>Blur (σ=2) +<br>SR (×4)+<br>Noise<br>(σ=20)+<br>JPEG (q=50) | BSRGAN | 22.88 | 0.5397 | 0.3445 | 0.1838 | 0.5402 | 64.81 |
| | Real-ESRGAN | 22.01 | 0.5332 | 0.3494 | 0.2115 | 0.5730 | 64.76 |
| | StableSR | 21.39 | 0.4744 | 0.3422 | 0.2974 | 0.7354 | 70.94 |
| | DiffBIR | 21.79 | 0.4895 | 0.3465 | 0.2821 | 0.7059 | 69.28 |
| | PASD | 21.90 | 0.5118 | 0.3493 | 0.2397 | 0.6326 | 70.43 |
| | SUPIR (ours) | 20.84 | 0.4604 | 0.3806 | 0.4688 | 0.8021 | 73.58 |

Table 1. Quantitative comparison. **Red** and blue colors represent the best and second best performance. ↓ represents the smaller the better, and for the others, the bigger the better.

| Metrics | BSRGAN | Real-ESRGAN | StableSR | DiffBIR | PASD | Ours |
|---|---|---|---|---|---|---|
| CLIP-IQA | 0.4119 | 0.5174 | 0.7654 | 0.6983 | 0.7714 | 0.8232 |
| MUSIQ | 55.64 | 59.42 | 70.70 | 69.69 | 71.87 | 73.00 |
| MANIQA | 0.1585 | 0.2262 | 0.3035 | 0.2619 | 0.3169 | 0.4295 |

(a) Quantitative comparison on 60 real-world LQ images.

| Negative Samples | Positive Prompts | Negative Prompts | PSNR | SSIM | LPIPS↓ | ManIQA | ClipIQA | MUSIQ |
|---|---|---|---|---|---|---|---|---|
| ✓ | | | 22.90 | 0.5519 | 0.3010 | 0.3129 | 0.7049 | 68.94 |
| ✓ | ✓ | | 22.31 | 0.5250 | 0.3108 | 0.4018 | 0.7937 | 72.00 |
| ✓ | | ✓ | 20.63 | 0.4747 | 0.3603 | 0.4678 | 0.7933 | 73.60 |
| ✓ | ✓ | ✓ | 20.66 | 0.4763 | 0.3412 | 0.4740 | 0.8164 | 73.66 |
| | ✓ | ✓ | 21.79 | 0.5119 | 0.3139 | 0.3180 | 0.7102 | 72.68 |

(b) Ablation study of quality prompts and negative training samples.

| Connector | PSNR | SSIM | LPIPS↓ | ManIQA | ClipIQA | MUSIQ |
|---|---|---|---|---|---|---|
| Zero Convolution [95] | 19.47 | 0.4261 | 0.3969 | 0.4845 | 0.8184 | 74.00 |
| ZeroSFT | 20.66 | 0.4763 | 0.3412 | 0.4740 | 0.8164 | 73.66 |

(c) Ablation study of zero convolution and the proposed ZeroSFT.

Table 2. Real-world comparison results and ablation studies.



Figure 8. (a) These plots illustrate the quantitative results as a function of the variable $\tau_r$. "No $\tau_r$" means not to use the proposed sampling method. (b) The results of our user study.



Low-Quality Input · Zero Convolution Connector · Our ZeroSFT Connector

Figure 9. We compare the proposed ZeroSFT with zero convolution. Directly using zero convolution results in redundant details. The low-fidelity details can be effectively mitigated by ZeroSFT.

imals, plants, faces, buildings, and landscapes. We show the qualitative results in Fig. 10, and the quantitative results are shown in Tab. 2a. These results indicate that the images produced by our method have the best perceptual quality. We also conduct a user study comparing our method on real-world LQ images, with 20 participants involved. For each set of comparison images, we instructed participants to choose the restoration result that was of the highest quality among these test methods. The results are shown in Fig. 8, revealing that our approach significantly outperformed state-of-the-art methods in perceptual quality.

### 4.3. Controlling Restoration with Textual Prompts

After training on a large dataset of image-text pairs and leveraging the feature of the diffusion model, our method can selectively restore images based on human prompts. Fig. 1(b) illustrates some examples. In the first case, the bike restoration is challenging without prompts, but upon receiving the prompt, the model reconstructs it accurately. In the second case, the material texture of the hat can be adjusted through prompts. In the third case, even high-level semantic prompts allow manipulation over face attributes. In addition to prompting the image content, we can also prompt the model to generate higher-quality images through negative-quality prompts. Fig. 11(a) shows two examples. It can be seen that the negative prompts are very effective in improving the overall quality of the output image. We also observed that prompts in our method are not always effective. When the provided prompts do not align with the LQ image, the prompts become ineffective, see Fig. 11(b). We consider this reasonable for an IR method to stay faithful to the provided LQ image. This reflects a significant distinction from text-to-image generation models and underscores the robustness of our approach.

that our results have better visual effects, but they do not have an advantage in these metrics. This phenomenon has also been noted in many studies as well [6, 26, 28]. We argue that with the improving quality of IR, there is a need to reconsider the reference values of existing metrics and suggest more effective ways to evaluate advanced IR methods. We also show some qualitative comparison results in Fig. 6. Even under severe degradation, our method consistently produces highly reasonable and high-quality images that faithfully represent the content of the LQ images.

**Restoration in the Wild.** We also test our method on real-world LQ images. We collect a total of 60 real-world LQ images from RealSR [8], DRealSR [83], Real47 [49], and online sources, featuring diverse content including an-
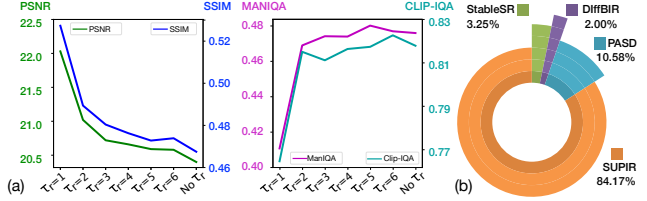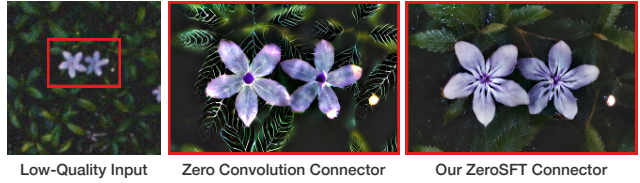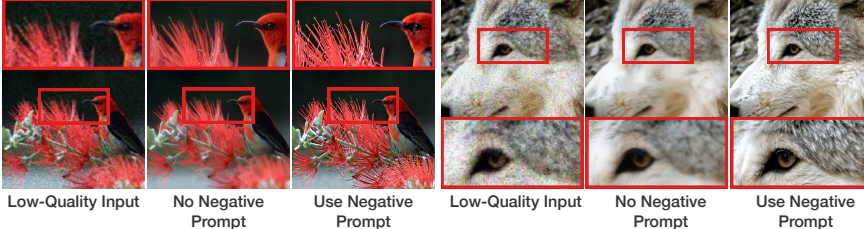
Figure 10. Qualitative comparison on real-world LQ images. SUPIR successfully recovers structured buildings and lifelike rivers. It also maintains the details existing in LQ, such as the horizontal planks in the beach chairs. Zoom in for better view.
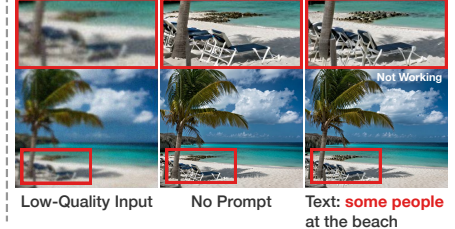


Figure 11. Influences of text prompts. (a) Negative prompts lead to detailed and sharp restoration results. (b) Given a positive prompt with hallucinations, SUPIR avoids generating content absent in the LQ images. Zoom in for better view.



Figure 12. Qualitative comparison for SUPIR training on datasets with different scales. Zoom in for better view.



Figure 13. The effect of the proposed restoration-guided sampling method. A smaller $\tau_r$ makes the result more biased toward the LQ image, which emphasizes the fidelity. A larger $\tau_r$ emphasizes perceived quality, but with lower fidelity. Zoom in for better view.

### 4.4. Ablation Study

**Connector.** We compare the proposed ZeroSFT connector with zero convolution [95]. Quantitative results are shown in Tab. 2c. Compared to ZeroSFT, zero convolution yields comparable performance on non-reference metrics and much lower full-reference performance. In Fig. 9, we find that the drop in non-reference metrics is caused by generating low-fidelity content. Therefore, for IR tasks, ZeroSFT ensures fidelity without losing the perceptual effect.

**Training data scaling.** We trained our large-scale model on two smaller datasets for IR, DIV2K [3] and LSDIR [1]. The qualitative results are shown in Fig. 12, which clearly demonstrate the importance and necessity of training on large-scale high-quality data.

**Negative-quality samples and prompt.** Tab. 2b shows some quantitative results under different settings. Here, we use positive words describing image quality as "positive prompt", and use negative quality words and the CFG methods described in Sec. 3.2 as negative prompt. It can be seen that adding positive prompts or negative prompts alone can improve the perceptual quality of the image. Using both of them simultaneously yields the best perceptual results. If negative samples are not included for training, these two prompts will not be able to improve the perceptual quality. Fig. 4 and Fig. 11(a) demonstrate the improvement in image quality brought by using negative prompts.

**Restoration-guided sampling method.** The proposed restoration-guided sampling method is mainly controlled by the hyper-parameter $\tau_r$. The larger $\tau_r$ is, the fewer corrections are made to the generation at each step. The smaller $\tau_r$ is, the more generated content will be forced to be closer to the LQ image. Please refer to Fig. 13 for a qualitative comparison. When $\tau_r = 0.5$, the image is blurry because its output is limited by the LQ image and cannot generate texture and details. When $\tau_r = 6$, there is not much guidance during generation. The model generates a lot of texture that is not present in the LQ image, especially in flat area. Fig. 8(a) illustrates the quantitative results of restoration as a function of the variable $\tau_r$. As shown in Fig. 8(a), decreasing $\tau_r$ from 6 to 4 does not result in a significant decline in visual quality, while fidelity performance improves. As restoration guidance continues to strengthen, although PSNR continues to improve, the images gradually become blurry with loss of details, as depicted in Fig. 13. Therefore, we choose $\tau_r = 4$ as the default parameter, as it doesn't significantly compromise image quality while effectively en-

hancing fidelity.

## 5. Conclusion

We propose SUPIR as a pioneering IR method, empowered by model scaling, dataset enrichment, and advanced design features, expanding the horizons of IR with enhanced perceptual quality and controlled textual prompts.

## References

[1] Lsdir dataset: A large scale dataset for image restoration. https://data.vision.ee.ethz.ch/yawli/index.html, 2023. Accessed: 2023-11-15. 4, 8

[2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4432–4441, 2019. 2

[3] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2017. 4, 8

[4] David Bau, Hendrik Strobelt, William Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. *arXiv preprint arXiv:2005.07727*, 2020. 2

[5] Sefi Bell-Kligler, Assaf Shocher, and Michal Irani. Blind super-resolution kernel estimation using an internal-gan. *Advances in Neural Information Processing Systems*, 32, 2019. 2

[6] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6228–6237, 2018. 7

[7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 2, 3

[8] 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. 7

[9] Kelvin CK Chan, Xintao Wang, Xiangyu Xu, Jinwei Gu, and Chen Change Loy. Glean: Generative latent bank for large-factor image super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14245–14254, 2021. 3

[10] 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. 2, 3

[11] Haoyu Chen, Jinjin Gu, Yihao Liu, Salma Abdel Magid, Chao Dong, Qiong Wang, Hanspeter Pfister, and Lei Zhu. Masked image training for generalizable deep image denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1692–1703, 2023. 2

[12] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart-$\alpha$: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 3, 5

[13] 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. 2

[14] Zheng Chen, Yulun Zhang, Ding Liu, Bin Xia, Jinjin Gu, Linghe Kong, and Xin Yuan. Hierarchical integration diffusion model for realistic image deblurring. *Advances in Neural Information Processing Systems*, 2023. 2

[15] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. *arXiv preprint arXiv:2108.02938*, 2021. 3

[16] DeepFloyd. Deepfloyd inference framework. https://www.deepfloyd.ai/deepfloyd-if, 2023. Accessed: 2023-11-14. 3

[17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 4

[18] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3

[19] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015. 2

[20] Chao Dong, Chen Change Loy, and Xiaoou Tang. Accelerating the super-resolution convolutional neural network. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 391–407. Springer, 2016. 2

[21] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 4

[22] Yuchen Fan, Jiahui Yu, Yiqun Mei, Yulun Zhang, Yun Fu, Ding Liu, and Thomas S Huang. Neural sparse representation for image restoration. *Advances in Neural Information Processing Systems*, 33:15394–15404, 2020. 2

[23] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 2

[24] Jinjin Gu and Chao Dong. Interpreting super-resolution networks with local attribution maps. In *Proceedings of*

the *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9199–9208, 2021. 4

[25] Jinjin Gu, Hannan Lu, Wangmeng Zuo, and Chao Dong. Blind super-resolution with iterative kernel correction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1604–1613, 2019. 2

[26] Jinjin Gu, Haoming Cai, Haoyu Chen, Xiaoxing Ye, Jimmy Ren, and Chao Dong. Pipal: a large-scale image quality assessment dataset for perceptual image restoration. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 633–651, 2020. 2, 7

[27] Jinjin Gu, Yujun Shen, and Bolei Zhou. Image processing using multi-code gan prior. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3012–3021, 2020. 2

[28] Jinjin Gu, Haoming Cai, Chao Dong, Jimmy S Ren, Radu Timofte, Yuan Gong, Shanshan Lao, Shuwei Shi, Jiahao Wang, Sidi Yang, et al. Ntire 2022 challenge on perceptual image quality assessment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 951–967, 2022. 7

[29] Jinjin Gu, Xianzheng Ma, Xiangtao Kong, Yu Qiao, and Chao Dong. Networks are slacking off: Understanding generalization problem in image deraining. *Advances in Neural Information Processing Systems*, 2023. 2, 4

[30] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 5, 13

[31] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3

[32] 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. 3

[33] Hengyuan Hu, Rui Peng, Yu-Wing Tai, and Chi-Keung Tang. Network trimming: A data-driven neuron pruning approach towards efficient deep architectures. *arXiv preprint arXiv:1607.03250*, 2016. 3

[34] Yan Huang, Shang Li, Liang Wang, Tieniu Tan, et al. Unfolding the alternating optimization for blind super resolution. *Advances in Neural Information Processing Systems*, 33:5632–5643, 2020. 2

[35] Zheng Hui, Jie Li, Xiumei Wang, and Xinbo Gao. Learning the non-differentiable optimization for blind super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2093–2102, 2021. 2

[36] Younghyun Jo, Sejong Yang, and Seon Joo Kim. Srflow-da: Super-resolution using normalizing flow with deep convolutional block. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 364–372, 2021. 3

[37] Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10124–10134, 2023. 3

[38] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. 3, 4

[39] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 2

[40] 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. 2, 4

[41] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022. 5

[42] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022. 2, 3

[43] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5148–5157, 2021. 6

[44] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 2, 3, 4, 5

[45] Xiangtao Kong, Xina Liu, Jinjin Gu, Yu Qiao, and Chao Dong. Reflash dropout in image super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6002–6012, 2022. 6

[46] 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. 2

[47] Jingyun Liang, Kai Zhang, Shuhang Gu, Luc Van Gool, and Radu Timofte. Flow-based kernel prior with application to blind super-resolution. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10601–10610, 2021.

[48] 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. 2

[49] 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, 3, 6, 7, 13, 18, 19, 20

[50] Anran Liu, Yihao Liu, Jinjin Gu, Yu Qiao, and Chao Dong. Blind image super-resolution: A survey and beyond. *IEEE transactions on pattern analysis and machine intelligence*, 45(5):5461–5480, 2022. 2

[51] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 13

[52] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 4, 13

[53] Yihao Liu, Anran Liu, Jinjin Gu, Zhipeng Zhang, Wenhao Wu, Yu Qiao, and Chao Dong. Discovering distinctive" semantics" in super-resolution networks. *arXiv preprint arXiv:2108.00406*, 2021. 2

[54] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5

[55] A Lugmayr, M Danelljan, L Van Gool, and R Timofte. Learning the super-resolution space with normalizing flow. *ECCV, Srflow*, 2020. 3

[56] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 14

[57] Sachit Menon, Alexandru Damian, Shijia Hu, Nikhil Ravi, and Cynthia Rudin. Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*, pages 2437–2445, 2020. 2

[58] Tomer Michaeli and Michal Irani. Nonparametric blind super-resolution. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 945–952, 2013. 2

[59] 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, 3

[60] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3883–3891, 2017. 2

[61] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3

[62] Xingang Pan, Xiaohang Zhan, Bo Dai, Dahua Lin, Chen Change Loy, and Ping Luo. Exploiting deep generative prior for versatile image restoration and manipulation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7474–7489, 2021. 3

[63] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 3

[64] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 2

[65] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3

[66] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 3

[67] 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. 2, 3, 5

[68] 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. 3

[69] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 4

[70] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019. 3

[71] 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. 3

[72] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Jiaya Jia. Scale-recurrent network for deep image deblurring. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8174–8182, 2018. 2

[73] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities. https://github.com/InternLM/InternLM, 2023. 2, 3

[74] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 2, 3

[75] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 3

[76] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2555–2563, 2023. 6

[77] 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. 3, 6, 18, 19, 20

[78] Longguang Wang, Yingqian Wang, Xiaoyu Dong, Qingyu Xu, Jungang Yang, Wei An, and Yulan Guo. Unsupervised degradation representation learning for blind super-resolution. In *Proceedings of the IEEE/CVF Conference on*

*Computer Vision and Pattern Recognition*, pages 10581–10590, 2021. 2

[79] 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. 4

[80] Xintao Wang, Yu Li, Honglun Zhang, and Ying Shan. Towards real-world blind face restoration with generative facial prior. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9168–9178, 2021. 3

[81] 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. 2, 5, 6, 18, 19, 20

[82] Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. *arXiv preprint arXiv:2212.00490*, 2022. 2, 3

[83] 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. 7

[84] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 4

[85] Zeyue Xue, Guanglu Song, Qiushan Guo, Boxiao Liu, Zhuofan Zong, Yu Liu, and Ping Luo. Raphael: Text-to-image generation via large mixture of diffusion paths. *arXiv preprint arXiv:2305.18295*, 2023. 3

[86] 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

[87] Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Gan prior embedded network for blind face restoration in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 672–681, 2021. 3

[88] 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. 6, 18, 19, 20

[89] Jiale Zhang, Yulun Zhang, Jinjin Gu, Yongbing Zhang, Linghe Kong, and Xin Yuan. Accurate image restoration with attention retractable transformer. In *International Conference on Learning Representations (ICLR)*, 2023. 2

[90] 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. 2

[91] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. Learning deep cnn denoiser prior for image restoration. In

[92] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Transactions on Image Processing*, 27(9): 4608–4622, 2018. 2

[93] 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. 2, 6

[94] Kai Zhang, Yawei Li, Jingyun Liang, Jiezhang Cao, Yulun Zhang, Hao Tang, Deng-Ping Fan, Radu Timofte, and Luc Van Gool. Practical blind image denoising via swin-conv-unet and data synthesis. *Machine Intelligence Research*, pages 1–14, 2023. 2

[95] 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, 3, 7, 8

[96] 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

[97] Ruofan Zhang, Jinjin Gu, Haoyu Chen, Chao Dong, Yulun Zhang, and Wenming Yang. Crafting training degradation distribution for the accuracy-generalization trade-off in real-world super-resolution. *International Conference on Machine Learning (ICML)*, 2023. 2, 6

[98] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. In *International Conference on Learning Representations (ICLR)*, 2019. 2

[99] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image restoration. *IEEE transactions on pattern analysis and machine intelligence*, 43(7):2480–2495, 2020. 2

[100] Yi Zhang, Xiaoyu Shi, Dasong Li, Xiaogang Wang, Jian Wang, and Hongsheng Li. A unified conditional framework for diffusion-based image restoration. *arXiv preprint arXiv:2305.20049*, 2023. 3

[101] Yang Zhao, Yu-Chuan Su, Chun-Te Chu, Yandong Li, Marius Renn, Yukun Zhu, Changyou Chen, and Xuhui Jia. Rethinking deep face restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7652–7661, 2022.

[102] Shangchen Zhou, Kelvin Chan, Chongyi Li, and Chen Change Loy. Towards robust blind face restoration with codebook lookup transformer. *Advances in Neural Information Processing Systems*, 35:30599–30611, 2022. 3

[103] Jiapeng Zhu, Deli Zhao, Bo Zhang, and Bolei Zhou. Disentangled inference for gans with latently invertible autoencoder. *International Journal of Computer Vision*, 130(5): 1259–1276, 2022. 3
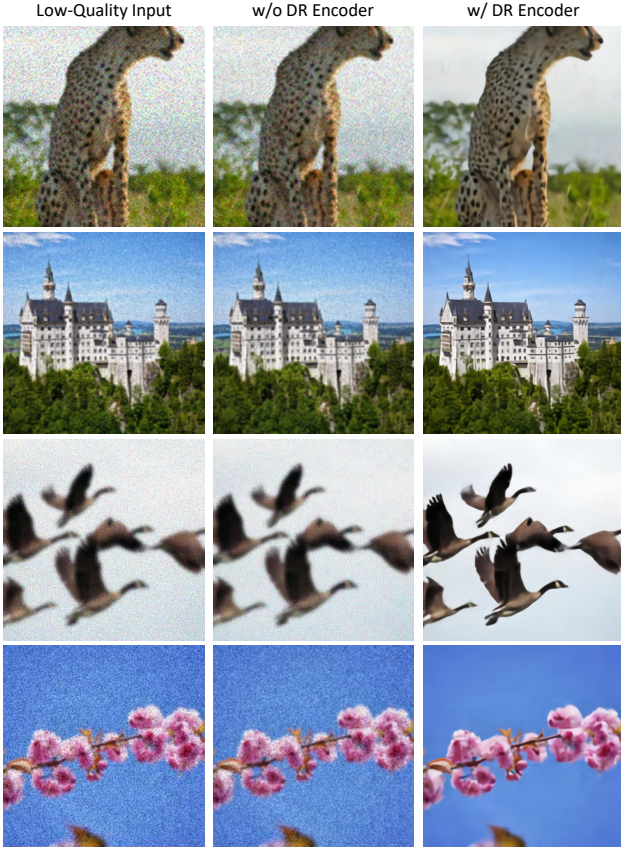
Figure 14. The effectiveness of the degradation-robust encoder (DR Encoder) is demonstrated by the results, which are achieved by initially encoding with various encoders and subsequently decoding. This process effectively reduces the degradations in low-quality inputs before they are introduced into the diffusion models.



Figure 15. Negative prompt causes artifacts when low-quality inputs do not have clear semantics.

# Appendix

# A. Discussions

## A.1. Degradation-Robust Encoder

As shown in Fig. 2 of the main text, a degradation-robust encoder is trained and deployed prior to feeding the low-quality input into the adaptor. We conduct experiments using synthetic data to demonstrate the effectiveness of the proposed degradation-robust encoder. In Fig. 14, we show the results of using the same decoder to decode the latent representations from different encoders. It can be seen that the original encoder has no ability to resist degradation and its decoded images still contain noise and blur. The proposed degradation-robust encoder can reduce the impact of degradation, which further prevents generative models from misunderstanding artifacts as image content [49].

## A.2. LLaVA Annotation

Our diffusion model is capable of accepting textual prompts during the restoration process. The prompt strategy we employ consists of two components: one component is automatically annotated by LLaVA-v1.5-13B [51], and the other is a standardized default positive quality prompt. The fixed portion of the prompt strategy provides a positive description of quality, including words like "`cinematic, High Contrast, highly detailed, unreal engine, taken using a Canon EOS R camera, hyper detailed photo-realistic maximum detail, 32k, Color Grading, ultra HD, extreme meticulous detailing, skin pore detailing, hyper sharpness, perfect without deformations, Unreal Engine 5, 4k render`". For the LLaVA component, we use the command "`Describe this image and its style in a very detailed manner`" to generate detailed image captions, as exemplified in Fig. 18. While occasional inaccuracies may arise, LLaVA-v1.5-13B generally captures the essence of the low-quality input with notable precision. Using the reconstructed version of the input proves effective in correcting these inaccuracies, allowing LLaVA to provide an accurate description of the majority of the image's content. Additionally, SUPIR is effective in mitigating the impact of potential hallucination prompts, as detailed in [52].

## A.3. Limitations of Negative Prompt

Figure 23 presents evidence that the use of negative quality prompts [30] substantially improves the image quality of restored images. However, as observed in Fig. 15, the negative prompt may introduce artifacts when the restoration target lacks clear semantic definition. This issue likely stems from a misalignment between low-quality inputs and language concepts.

## A.4. Negative Samples Generation

While negative prompts are highly effective in enhancing quality, the lack of negative-quality samples and prompts in

the training data results in the fine-tuned SUPIR's inability to comprehend these prompts effectively. To address this problem, in Sec. 3.2 of the main text, we introduce a method to distill negative concepts from the SDXL model. The process for generating negative samples is illustrated in Fig. 19. Direct sampling of negative samples through a text-to-image approach often results in meaningless images. To address this issue, we also utilize training samples from our dataset as source images. We create negative samples in an image-to-image manner as proposed in [56], with a strength setting of $0.5$.

## B. More Visual Results

We provide more results in this section. Fig. 16 presents additional cases where full-reference metrics do not align with human evaluation. In Fig. 17, we show that using negative-quality prompt without including negative samples in training may cause artifacts. In Figs. 20 to 22, we provide more visual caparisons with other methods. Plenty of examples prove the strong restoration ability of SUPIR and the most realistic of restored images. More examples of controllable image restoration with textual prompts can be found in Fig. 23.

Figure 16. Additional samples highlight the misalignment between metric evaluations and human assessments. While SUPIR produces images with high-fidelity textures, it tends to receive lower scores in metric evaluations.
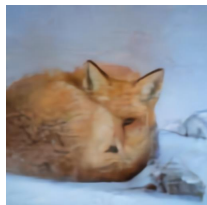


Figure 17. More visual results for Fig. 4 of the main text. CFG introduces artifacts if we do not include negative-quality samples in training. Adding negative-quality samples allows further quality enhancement through CFG.
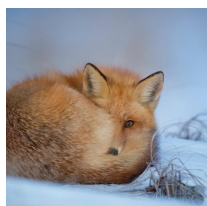
Case 1



Low-Quality Input

The image features a small, brown, furry animal, possibly a **hamster** or a **cat**, laying on a white surface. The animal appears to be relaxed and comfortable, as it is resting on a **bed** or a **cushion**. The white background provides a contrast to the brown fur of the animal, making it the focal point of the image.



Reconstructed Input

The image features a **fox** lying down in the snow, appearing to be sleeping or resting. The fox is positioned in the center of the scene, with **its body occupying a significant portion of the image**. The **snowy** background provides a serene and natural setting for the fox, emphasizing its natural habitat.



Ground-Truth

The image features a **fox** curled up in the snow, resting comfortably. The fox is positioned in the center of the scene, with **its body facing the left side of the image**. The **snowy** environment provides a natural and serene backdrop for the fox, which appears to be enjoying its time in the snow.
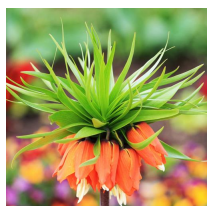
Case 2



Low-Quality Input

The image features a beautiful garden with a variety of colorful flowers. The flowers are arranged in a visually appealing manner, creating a vibrant and lively atmosphere. The garden is filled with different types of flowers, each with unique colors and shapes, contributing to the overall beauty of the scene…



Reconstructed Input

The image features a beautiful flower garden with a variety of colorful flowers. The **main focus is on a large, bright orange flower with a green center**, surrounded by other vibrant flowers. The orange flower is situated in the middle of the garden, drawing attention to its striking color and unique shape. The garden is filled with a diverse array of flowers, creating a visually stunning and lively atmosphere.



Ground-Truth

The image features a vibrant garden with **a large, colorful flower in the center**. The flower is surrounded by a variety of other flowers, creating a beautiful and lively scene. The main flower is **orange** and **yellow**, with a **green stem**, and it stands out among the other flowers in the garden. The garden is filled with a diverse assortment of flowers…
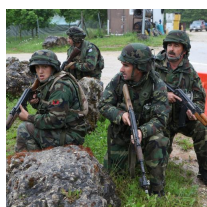
Case 3



Low-Quality Input

The image features a group of **three** men dressed in military uniforms, sitting on the ground and holding guns. They are positioned in a line, with one man on the left, another in the middle, and the third on the right. Each of them is holding a rifle, with one rifle located on the left side, another in the middle, and the third on the right side of the group. The men appear to be soldiers…
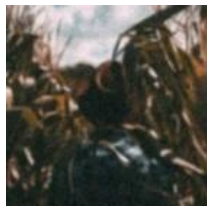


Reconstructed Input

The image features a group of **four** men dressed in military uniforms, standing in a **grassy area**. They are all holding guns, with some of them also wearing **backpacks**. The men appear to be soldiers, possibly on a mission or training exercise. The scene captures their readiness and focus as they prepare for their task.
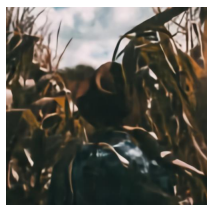


Ground-Truth

The image features a group of **four** men dressed in military uniforms, standing in a **field** and holding guns. They appear to be soldiers, possibly engaged in a training exercise or a mission. The men are positioned in a line, with one soldier on the left, another in the middle, and two more on the right side of the image. Each soldier is holding a gun…
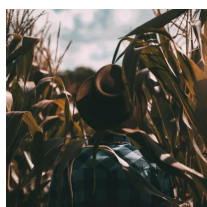
Case 4



Low-Quality Input

The image features a person standing in a field of tall grass, surrounded by tall corn plants. The person appears to be wearing a **backpack**, and their head is partially hidden by the tall grass. The scene gives off a sense of adventure and exploration, as the person seems to be navigating through the field. The tall grass and corn plants create a sense of depth and natural…



Reconstructed Input

The image features a person standing in a field of tall corn, with the corn surrounding them on all sides. The person is wearing a **hat** and appears to be looking down, possibly observing the corn or the ground. The field is vast, with the corn reaching up to the person's shoulders, creating a sense of being engulfed by the tall plants. The scene captures the essence of being in a cornfield…
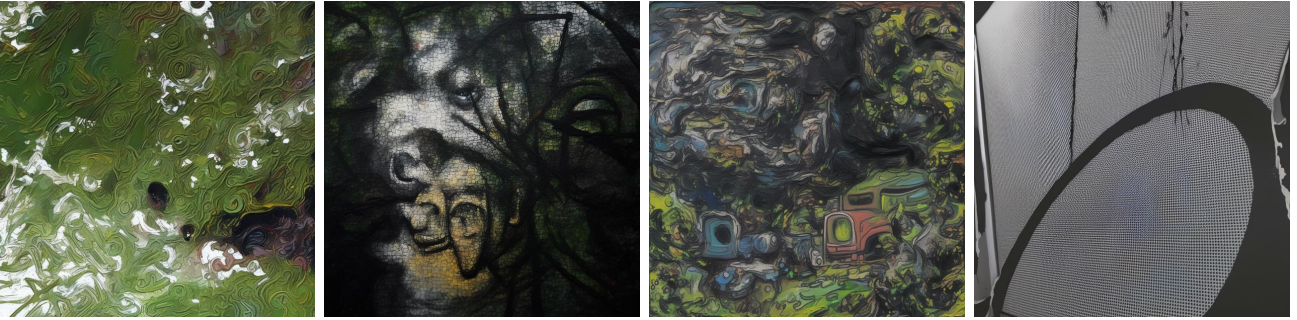


Ground-Truth

The image features a man wearing a **hat** and a **plaid shirt**, standing in a field of tall corn. He appears to be looking over the corn, possibly observing the surroundings or searching for something. The corn is quite tall, reaching up to the man's shoulders, and the field extends in the background, creating a sense of depth and vastness. The man's presence in the field, along with his attire…

Figure 18. Snapshots showcasing LLaVA annotations demonstrate that LLaVA accurately predicts most content even with low-quality inputs. Please zoom in for a more detailed view.

16

(a) Noise to Image      Prompt = {oil painting, cartoon, blurring, dirty, messy, low quality, frames, deformed, lowres, over-smooth}



(b) Image to Image      Prompt = {oil painting, cartoon, blurring, dirty, messy, low quality, frames, deformed, lowres, over-smooth}
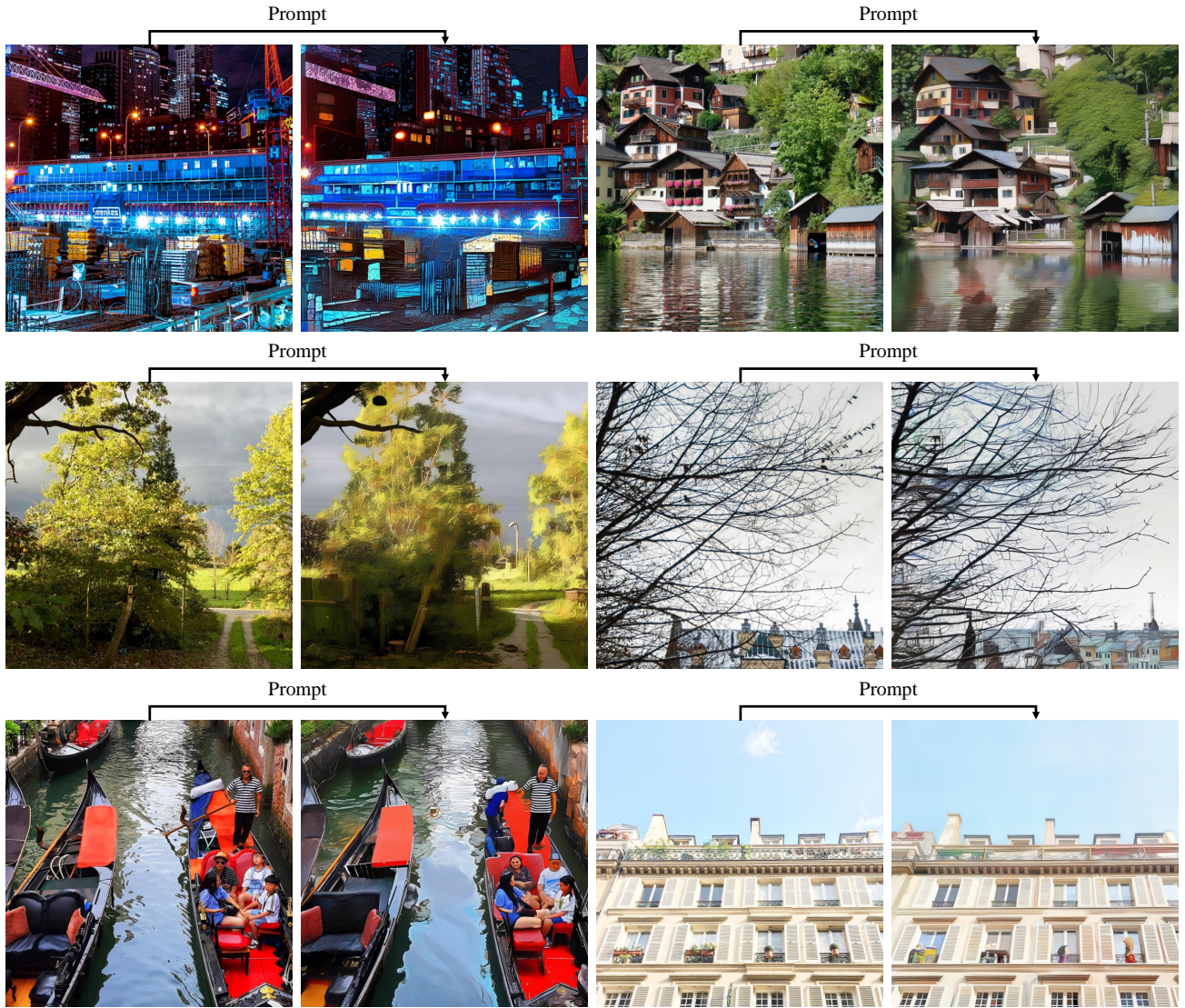


Figure 19. Pipeline of negative sample generation. (a) Sampling in a noise-to-image approach leads to meaningless outputs. (b) We synthetic negative samples from high quality images. Zoom in for better view.
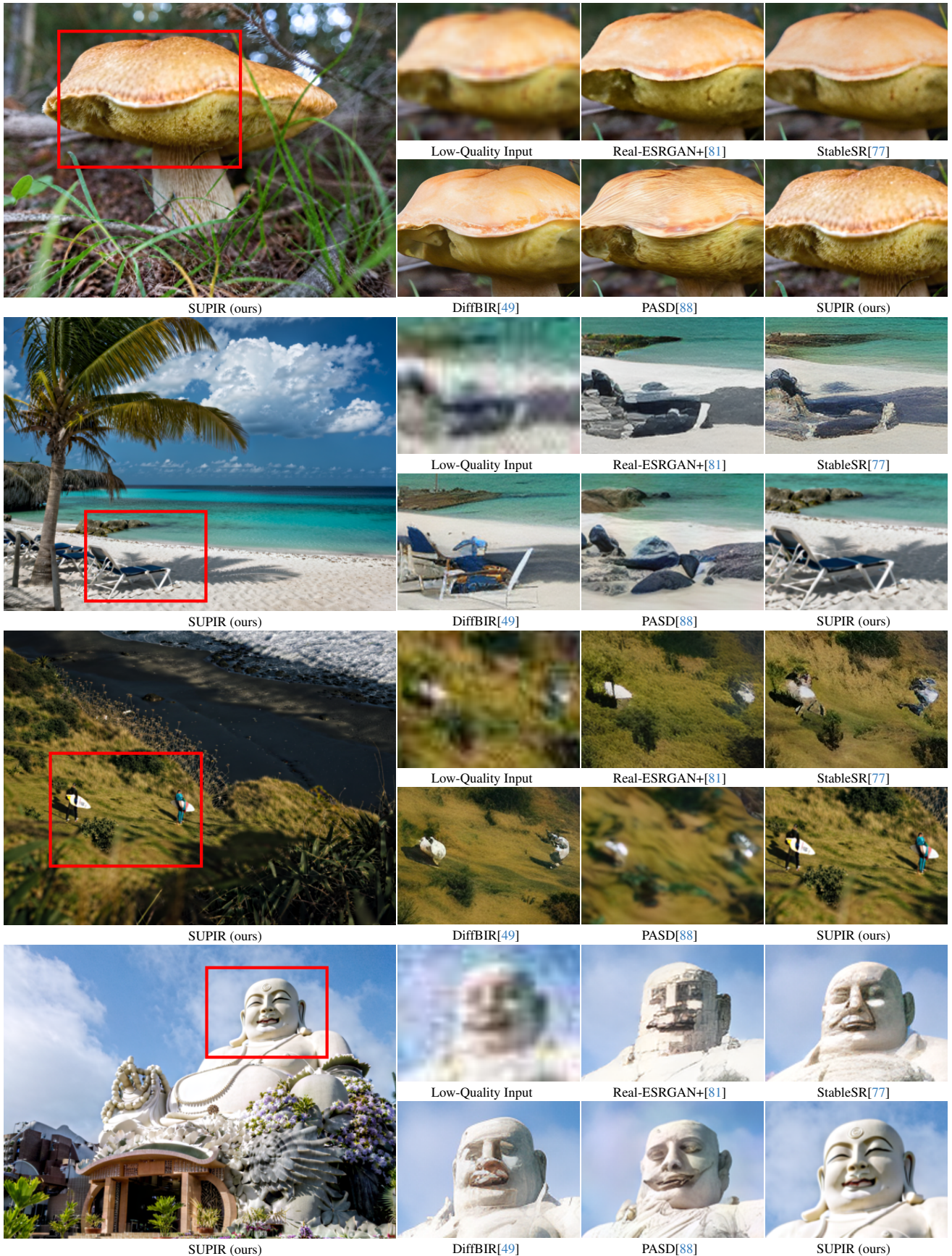
Figure 20. Qualitative comparison with different methods. Our method can accurately restore the texture and details of the corresponding object under challenging degradation. Zoom in for better view.

Figure 21. Qualitative comparison with different methods. Our method can accurately restore the texture and details of the corresponding object under challenging degradation. Zoom in for better view.
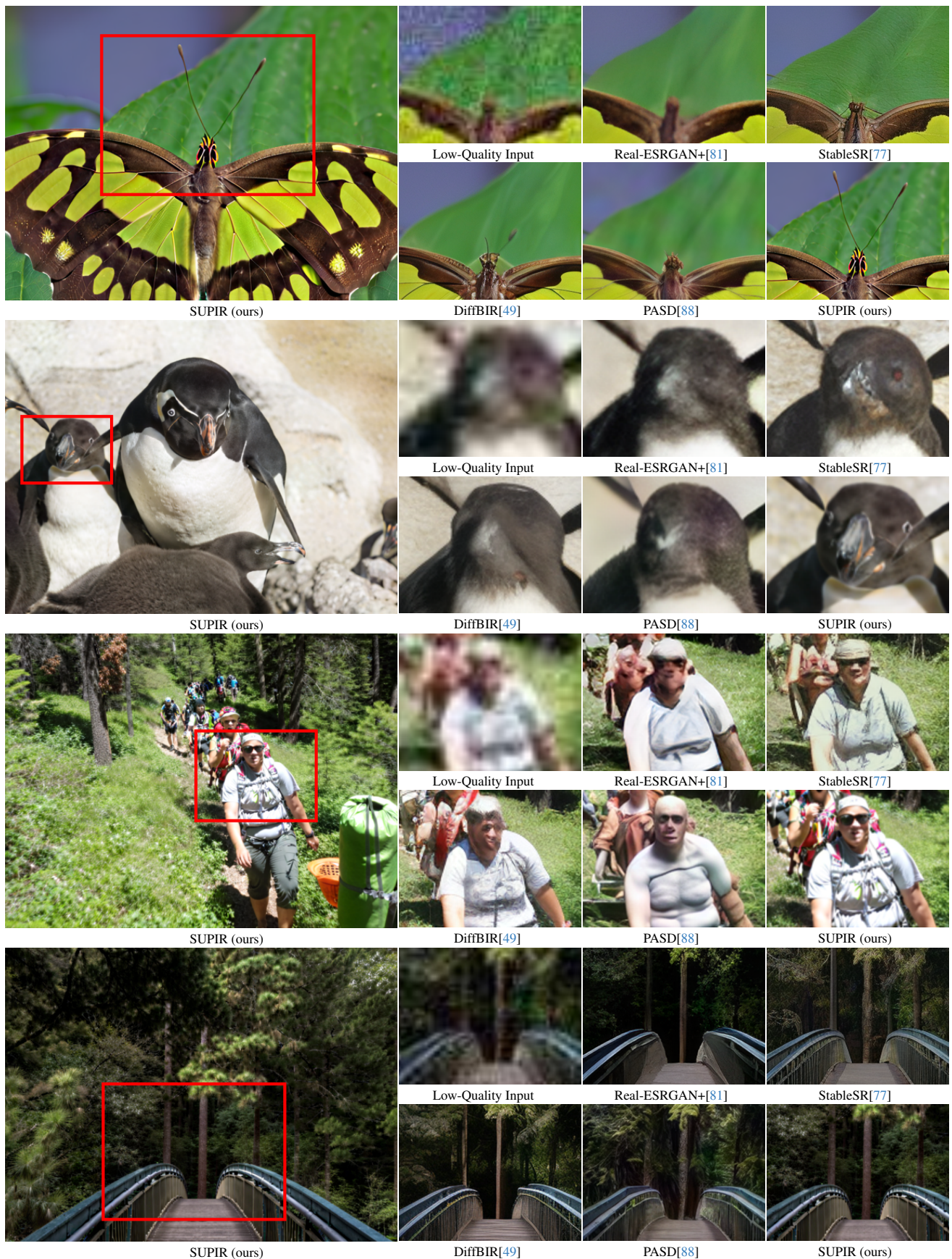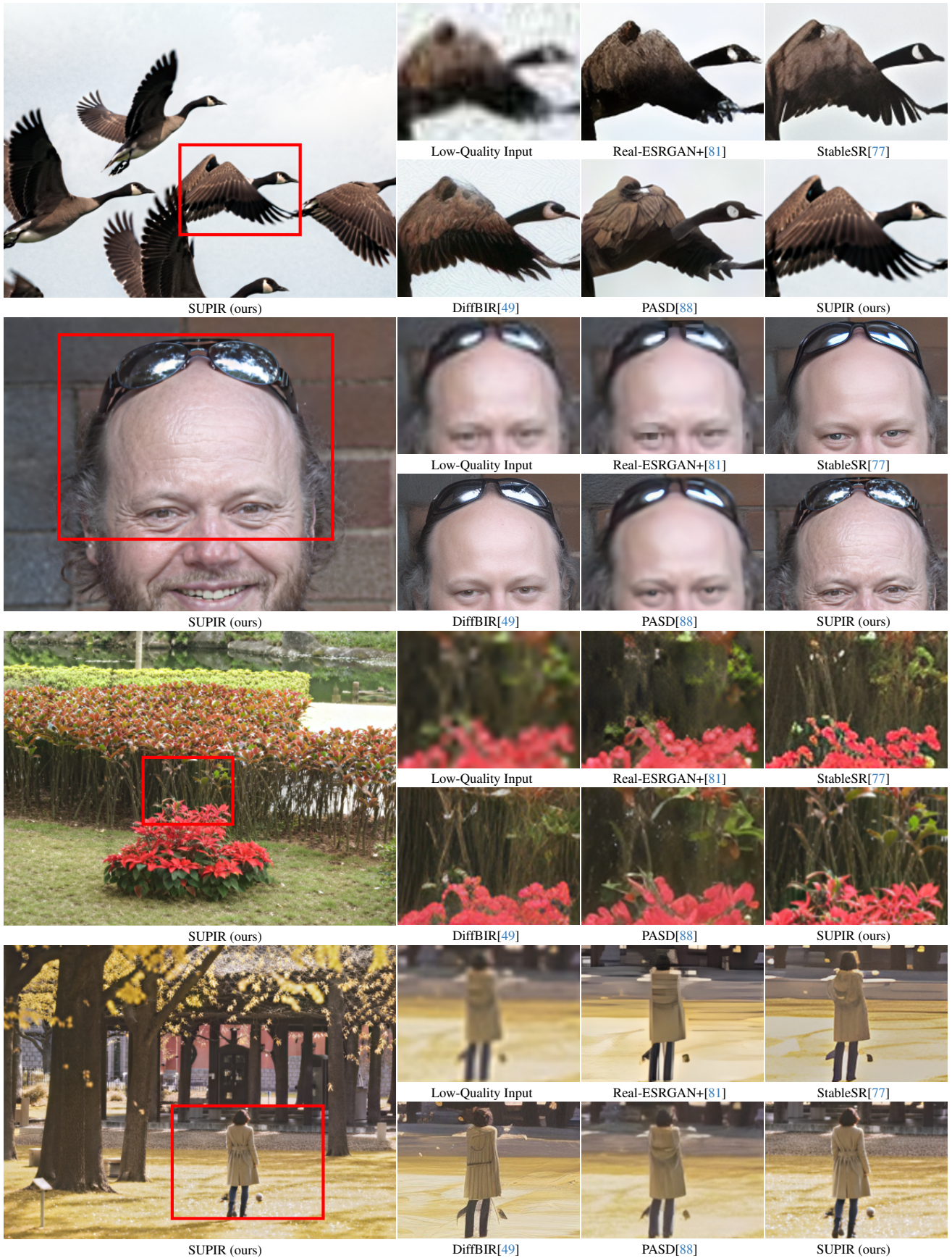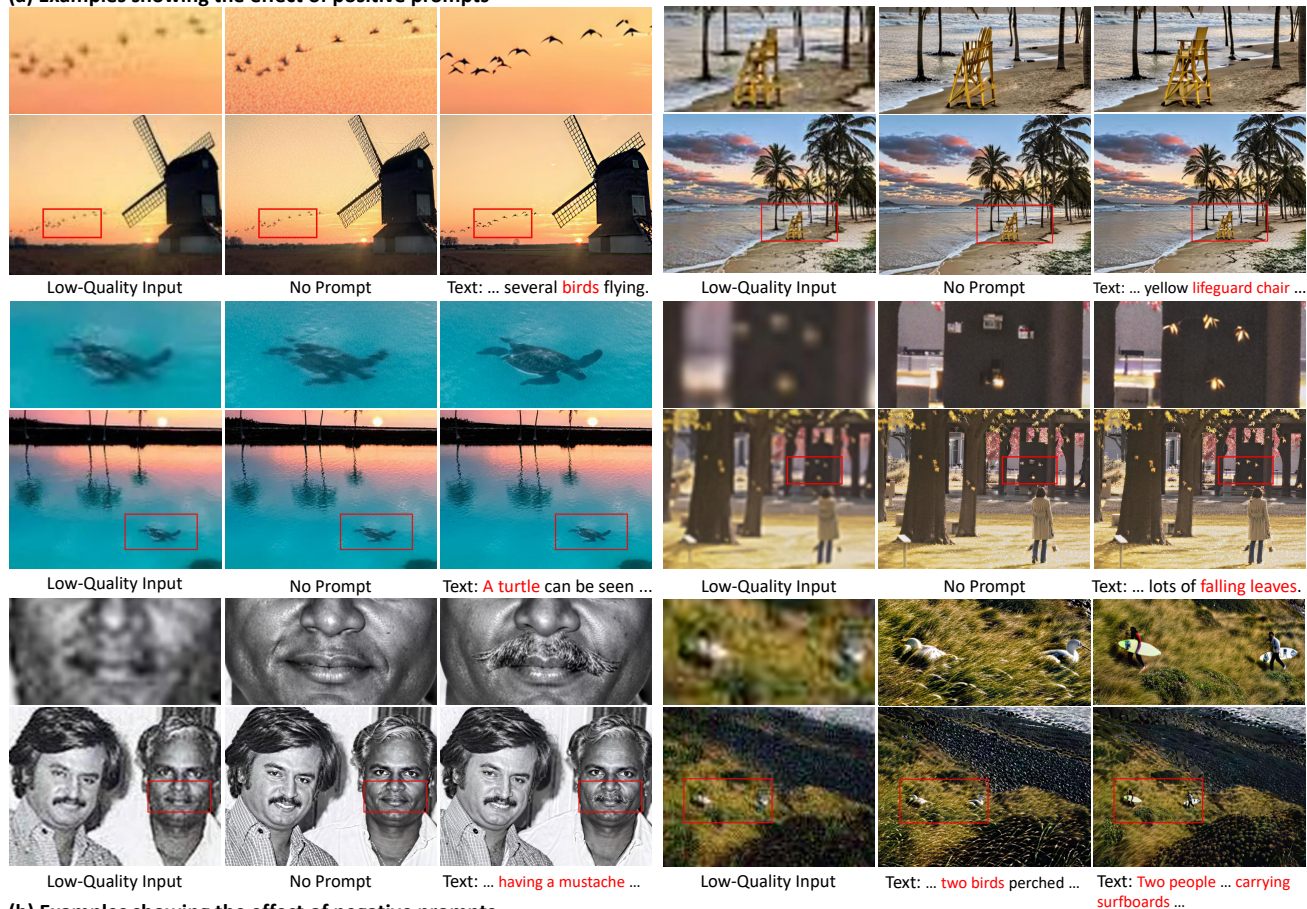
Figure 22. Qualitative comparison with different methods. Our method can accurately restore the texture and details of the corresponding object under challenging degradation. Zoom in for better view.

**(a) Examples showing the effect of positive prompts**



| Low-Quality Input | No Prompt | Text: … several birds flying. | Low-Quality Input | No Prompt | Text: … yellow lifeguard chair … |

| Low-Quality Input | No Prompt | Text: A turtle can be seen … | Low-Quality Input | No Prompt | Text: … lots of falling leaves. |

| Low-Quality Input | No Prompt | Text: … having a mustache … | Low-Quality Input | Text: … two birds perched … | Text: Two people … carrying surfboards … |

**(b) Examples showing the effect of negative prompts**



| Low-Quality Input | No Negative Prompt | Use Negative Prompt | Low-Quality Input | No Negative Prompt | Use Negative Prompt |

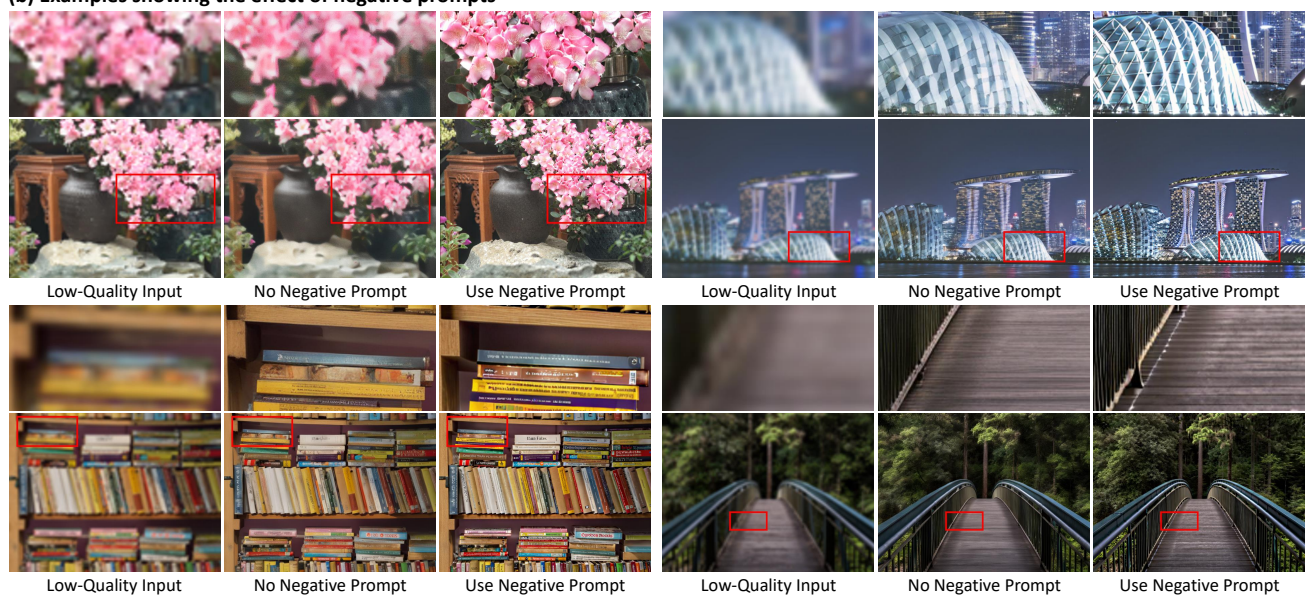| Low-Quality Input | No Negative Prompt | Use Negative Prompt | Low-Quality Input | No Negative Prompt | Use Negative Prompt |

Figure 23. More visual results of the text prompts' influences. (a) and (b) show the examples of positive prompts and negative prompts, respectively. Zoom in for better view.