

# Generative Diffusion Prior for Unified Image Restoration and Enhancement

Ben Fei<sup>1,2,\*</sup>, Zhaoyang Lyu<sup>2,\*</sup>, Liang Pan<sup>3</sup>, Junzhe Zhang<sup>3</sup>,  
Weidong Yang<sup>1,†</sup>, Tianyue Luo<sup>1</sup>, Bo Zhang<sup>2</sup>, Bo Dai<sup>2,†</sup>

<sup>1</sup>Fudan University, <sup>2</sup>Shanghai AI Laboratory, <sup>3</sup>S-Lab, Nanyang Technological University

bfei21@m.fudan.edu.cn, wdyang@fudan.edu.cn, (lvzhaoyang, daibo)@pjlab.org.cn



Figure 1. **Generative Diffusion Prior (GDP) is capable of generating high-fidelity restoration across various tasks.** GDP gives faithful image recovery on (a) **linear and multi-linear restoration**. In addition, GDP also enables novel applications of (b) **blind, non-linear, multiple-guidance, or any-size image**, including low-light enhancement and HDR recovery.

## Abstract

Existing image restoration methods mostly leverage the posterior distribution of natural images. However, they often assume known degradation and also require supervised training, which restricts their adaptation to complex real applications. In this work, we propose the Generative Diffusion Prior (GDP) to effectively model the posterior distributions in an unsupervised sampling manner. GDP utilizes a pre-train denoising diffusion generative model (DDPM) for solving linear inverse, non-linear, or blind problems. Specifically, GDP systematically explores a protocol of conditional guidance, which is verified more practical than the commonly used guidance way. Furthermore, GDP is strength at optimizing the parameters of degradation model during the denoising process, achieving blind image restoration. Besides, we devise hierarchical guidance and patch-based methods, enabling the GDP to generate images of arbitrary resolutions. Experimentally, we demonstrate GDP’s versatility on several image datasets for linear problems, such as super-resolution, deblurring, inpainting, and colorization, as well as non-linear and blind issues, such as low-light enhancement and HDR image re-

covery. GDP outperforms the current leading unsupervised methods on the diverse benchmarks in reconstruction quality and perceptual quality. Moreover, GDP also generalizes well for natural images or synthesized images with arbitrary sizes from various tasks out of the distribution of the ImageNet training set. The project page is available at <https://generativediffusionprior.github.io/>

## 1. Introduction

Image quality often degrades during capture, storage, transmission, and rendering. Image restoration and enhancement [44] aim to inverse the degradation and improve the image quality. Typically, restoration and enhancement tasks can be divided into two main categories: 1) **Linear inverse problems**, such as image super-resolution (SR) [24, 39], deblurring [37, 80], inpainting [93], colorization [38, 100], where the degradation model is usually linear and known; 2) **Non-linear or blind problems** [1], such as image low-light enhancement [41] and HDR image recovery [10, 84], where the degradation model is non-linear and unknown. For a specific linear degradation model, image restoration can be tackled through end-to-end supervised training of neural networks [16, 100]. Nonetheless, corrupted images in the real world often have multiple complex degradations [60], where fully supervised approaches suffer

\*Equal contribution, †Corresponding author.

to generalize.

There is a surge of interest to seek for more general image priors through generative models [1, 21, 74], and tackle image restoration in an unsupervised setting [8, 19], where multiple restoration tasks of different degradation models can be addressed during inference without re-training. For instance, Generative Adversarial Networks (GANs) [20] that are trained on a large dataset of clean images learn rich knowledge of the real-world scenes have succeeded in various linear inverse problems through GAN inversion [21, 54, 62]. In parallel, Denoising Diffusion Probabilistic Models (DDPMs) [2, 7, 36, 72, 77, 82] have demonstrated impressive generative capabilities, level of details, and diversity on top of GAN [27, 66, 67, 76, 78, 81]. As an early attempt, Kawar *et al.* [32] explore pre-trained DDPMs with variational inference, and achieve satisfactory results on multiple restoration tasks, but their Denoising Diffusion Restoration Model (DDRM) leverages the singular value decomposition (SVD) on a known linear degradation matrix, making it still limited to linear inverse problems.

In this study, we take a step further and propose an efficient approach named Generative Diffusion Prior (GDP). It exploits a well-trained DDPM as effective prior for general-purpose image restoration and enhancement, using degraded image as guidance. As a unified framework, GDP not only works on various linear inverse problems, but also generalizes to non-linear, and blind image restoration and enhancement tasks for the first time. However, solving the blind inverse problem is not trivial, as one would need to concurrently estimate the degradation model and recover the clean image with high fidelity. Thanks to the generative prior in a pre-trained DDPM, denoising within the DDPM manifold naturally regularizes the realness and fidelity of the recovered image. Therefore, we adopt a blind degradation estimation strategy, where the degradation model parameters of GDP are randomly initialized and optimized during the denoising process. Moreover, to further improve the photorealism and image quality, we systematically investigate an effective way to guide the diffusion models. Specifically, in the sampling process, the pre-trained DDPM first predicts a clean image  $\tilde{x}_0$  from the noisy image  $x_t$  by estimating the noise in  $x_t$ . We can add guidance on this intermediate variable  $\tilde{x}_0$  to control the generation process of the DDPMs. In addition, with the help of the proposed hierarchical guidance and patch-based generation strategy, GDP is able to recover images of arbitrary resolutions, where low-resolution images and degradation models are first predicted to guide the generation of high-resolution images.

We demonstrate the empirical effectiveness of GDP by comparing it with various competitive unsupervised methods under the linear or multi-linear inverse problem on ImageNet [14], LSUN [94], and CelebA [31] datasets in terms of consistency and FID. Over the low-light [41] and

NTIRE [63] datasets, we further show GDP results on non-linear and blind issues, including low-light enhancement and HDR recovery, superior to other zero-shot baselines both qualitatively and quantitatively, manifesting that GDP trained on ImageNet also works on images out of its training set distribution.

Our contributions are fourfold: (1) To our best knowledge, GDP is the first unified problem solver that can effectively use **a single unconditional DDPM pre-trained on ImageNet** provide by [15] to produce diverse and high-fidelity outputs for unified image restoration and enhancement in an unsupervised manner. (2) GDP is capable of optimizing randomly initiated parameters of degradation that are unknown, resulting in a powerful framework that can tackle any blind image restoration. (3) Further, to achieve arbitrary size image generation, we propose hierarchical guidance and patch-based methods, greatly promoting GDP on natural image enhancement. (4) Moreover, the comprehensive experiments are carried out, different from the conventional guidance way, where GDP directly predicts the temporary output given the noisy image in every step, which will be leveraged to guide the generation of images in the next step.

## 2. Related works

**Linear Inverse Image Restoration.** Most diffusion models toward linear inverse problems have employed unconditional models for the conditional tasks [53, 79], where only one model needs to be trained. However, unconditional tasks tend to be more difficult than conditional tasks. Moreover, the multi-linear task is also a relatively under-explored subject in image restoration. For instance, [65, 95] train simultaneously on multiple tasks, but they mainly concentrate on the enhancement tasks like deblurring and so on. Some works have also handled the multi-scale super-resolution by simultaneous training over multiple degradations [35]. Here, we propose GDP as a single model for dealing with single linear inverse or multiple linear inverse tasks.

**Non-linear Image Restoration.** The non-linear image formation model provides an accurate description of several imaging systems, including camera response functions in high-dynamic-range imaging [68]. The non-linear image restoration model is more accurate but is often more computationally intractable. Recently, great attention has been paid to non-linear image restoration problems. For example, HDR-GAN [59] is proposed for synthesizing HDR images from multi-exposed LDR images, while Enlighten-GAN [29] is devised as an unsupervised GAN to generalize very well on various real-world test images. The diffusion models are rarely studied for non-linear image restoration.

**Blind Image Restoration.** Early supervised attempts [5, 25] tend to estimate the unknown point spread function. As an example, [34] designs a class of structured denoisers, and

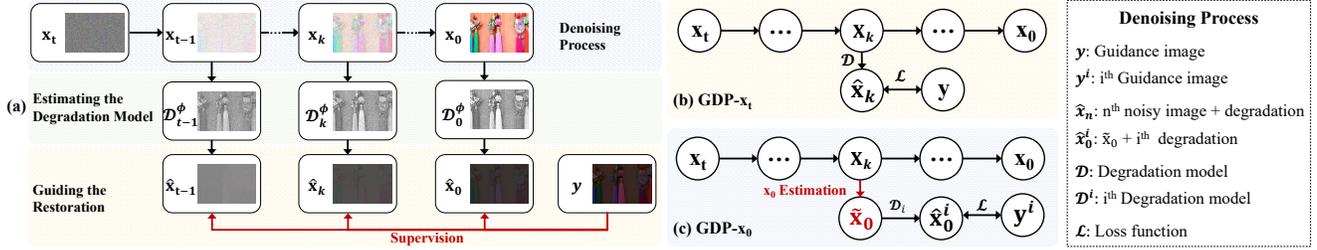


Figure 2. **Overview of our GDP for unified image recovery.** (a) Given a corrupted image  $\mathbf{y}$  during inference, GDP systematically studies the reverse process from  $\mathbf{x}_T$  to  $\mathbf{x}_0$  guided by the  $\mathbf{y}$ . The guidance can be added on  $\hat{\mathbf{x}}_0$  or  $\hat{\mathbf{x}}_t$ , leading to two variants of GDP. And  $\hat{\mathbf{x}}_0$  and  $\hat{\mathbf{x}}_t$  can be collectively expressed as  $\hat{\mathbf{x}}_t$ . The supervision signal (Sec. 5) is applied between  $\hat{\mathbf{x}}_t$  and  $\mathbf{y}$ . GDP looks for an intermediate variable  $\mathbf{x}_t$  and optimizes the degradation model  $\{\mathcal{D}_\phi^i \mid i = 1, 2, \dots, n\}$  that best reconstruct the image corresponding to  $\mathbf{y}$  via gradient descent. Note that GDP is a generic image restoration method. We illustrate it with the low light enhancement example. (b)  $\text{GDP-}\mathbf{x}_t$  adds guidance  $\mathbf{x}_t$  in every time step (Algo. 1), while (c)  $\text{GDP-}\mathbf{x}_0$  estimates the  $\tilde{\mathbf{x}}_0$  given  $\mathbf{x}_t$ , then adds guidance on  $\tilde{\mathbf{x}}_0$  to obtain  $\hat{\mathbf{x}}_0$  (Algo. 2). The number of guidance images  $\{\mathbf{y}^i \mid i = 1, 2, \dots, n\}$  and the degradation models  $\{\mathcal{D}_\phi^i \mid i = 1, 2, \dots, n\}$  are dependent on the tasks. For instance,  $n = 3$  for HDR recovery, while  $n = 1$  for other tasks.

[75] employs a fixed down-sampling operation to generate synthetic pairs during testing. However, these methods often are incapable of obtaining the parameters or distribution of the observed data due to the complicated degradation types. Another way to solve the blind image restoration is to utilize unsupervised learning methods [18, 107]. Following the CycleGAN [107], CinCGAN [96] and MCinCGAN [103] employ a pre-trained SR model together with cycle consistency loss to learn a mapping from the input image to high-quality image space. However, it still remains a challenge to exploit a unified architecture for blind image restoration. By the merits of the powerful GDP, these blind problems could also be solved by simultaneously estimating the recovered image and a specific degradation model.

### 3. Preliminary

Diffusion models [2, 22, 36, 71, 89] transform complex data distribution  $\mathbf{x}_0 \sim p_{\text{data}}$  into simple noise distribution  $\mathbf{x}_T \sim p_{\text{latent}} = \mathcal{N}(\mathbf{0}, \mathbf{I})$  and recover data from noise, where  $\mathcal{N}$  is the Gaussian distribution. DDPMs mainly comprise the diffusion process and the reverse process.

*The Diffusion Process* is a Markov chain that gradually corrupts data  $\mathbf{x}_0$  until it approaches Gaussian noise  $p_{\text{latent}}$  at  $T$  diffusion time steps. Corrupted data  $\mathbf{x}_1, \dots, \mathbf{x}_T$  are sampled from data  $p_{\text{data}}$ , with a diffusion process, which is defined as Gaussian transition:

$$q(\mathbf{x}_1, \dots, \mathbf{x}_T \mid \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1}), \quad (1)$$

where  $t$  denotes as diffusion step,  $q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$ , and  $\beta_t$  are fixed or learned variance schedule. An important property of the forward noising process is that any step  $\mathbf{x}_t$  may be sampled directly from  $\mathbf{x}_0$  through the following equation:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad (2)$$

where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ .

Proved by Ho *et al.* [27], there is a closed form expression for  $q(\mathbf{x}_t \mid \mathbf{x}_0)$ . We can obtain  $q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$ . Herein,  $\bar{\alpha}_t$  goes to 0 with large  $T$ , and  $q(\mathbf{x}_T \mid \mathbf{x}_0)$  is close to the latent distribution  $p_{\text{latent}}$ . *The Reverse Process* is a Markov chain that iteratively denoises a sampled Gaussian noise to a clean image. Starting from noise  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , the reverse process from latent  $\mathbf{x}_T$  to clean data  $\mathbf{x}_0$  is defined as:

$$p_\theta(\mathbf{x}_0, \dots, \mathbf{x}_{T-1} \mid \mathbf{x}_T) = \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t), \quad (3)$$

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \Sigma_\theta \mathbf{I})$$

According to Ho *et al.* [27], the mean  $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$  is the target we want to estimate by a neural network  $\theta$ . The variance  $\Sigma_\theta$  can be either time-dependent constants [27] or learnable parameters [58].  $\epsilon_\theta$  is a function approximator intended to predict  $\epsilon$  from  $\mathbf{x}_t$  as follow:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) \quad (4)$$

In practice,  $\tilde{\mathbf{x}}_0$  is usually predicted from  $\mathbf{x}_t$ , then  $\mathbf{x}_{t-1}$  is sampled using both  $\tilde{\mathbf{x}}_0$  and  $\mathbf{x}_t$  computed as:

$$\tilde{\mathbf{x}}_0 = \frac{\mathbf{x}_t}{\sqrt{\alpha_t}} - \frac{\sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}} \quad (5)$$

$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \tilde{\mathbf{x}}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \tilde{\mathbf{x}}_0), \tilde{\beta}_t \mathbf{I}),$$

$$\text{where } \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \tilde{\mathbf{x}}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \tilde{\mathbf{x}}_0 + \frac{\sqrt{\alpha_t} (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t$$

$$\text{and } \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \quad (6)$$

### 4. Generative Diffusion Prior

In this study, we aim to exploit a well-trained DDPM as an effective prior for unified image restoration and enhancement, in particular, to handle degraded images of a

Table 1. Comparison of different generative priors and regularization priors for image restoration and enhancement.

Methods	DGP [62]	SNIPS [33]	RED [69]	DDRM [32]	GDP (Ours)
Prior	GAN	MMSE Gaussian denoiser	Laplacian-based regularization function	DDPM	DDPM
Linear	✓	✓	✓	✓	✓
Non-linear	✗	✗	✗	✗	✓
Blind	✗	✗	✗	✗	✓

wide range of varieties. In detail, assume degraded image  $\mathbf{y}$  is captured via  $\mathbf{y} = \mathcal{D}(\mathbf{x})$ , where  $\mathbf{x}$  is the original natural image, and  $\mathcal{D}$  is a degradation model. We employ statistics of  $\mathbf{x}$  stored in some prior and search in the space of  $\mathbf{x}$  for an optimal  $\mathbf{x}$  that best matches  $\mathbf{y}$ , regarding  $\mathbf{y}$  as corrupted observations of  $\mathbf{x}$ . Due to the limited GAN inversion performance and the restricted applications of previous works [32, 33, 62, 69] in Table 1, in this paper, we focus on studying a more generic image prior, *i.e.*, the diffusion models trained on large-scale natural images for image synthesis. Inspired by the [4, 9, 12, 70, 73], the reverse denoising process of the DDPM can be conditioned on the degraded image  $\mathbf{y}$ . Specifically, the reverse denoising distribution  $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$  in Eq. 3 is adopted to a conditional distribution  $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$ . [15, 76] prove that

$$\begin{aligned} \log p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) &= \log(p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) p(\mathbf{y}|\mathbf{x}_t)) + K_1 \\ &\approx \log p(\mathbf{r}) + K_2, \end{aligned} \quad (7)$$

where  $\mathbf{r} \sim \mathcal{N}(\mathbf{r}; \mu_{\theta}(\mathbf{x}_t, t) + \Sigma \mathbf{g}, \Sigma)$  and  $\mathbf{g} = \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t)$ , where  $\Sigma = \Sigma_{\theta}(\mathbf{x}_t)$  for conciseness.  $K_1$  and  $K_2$  are constants,  $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$  is defined by Eq. 3.  $p(\mathbf{y} | \mathbf{x}_t)$  can be regarded as the probability that  $\mathbf{x}_t$  will be denoised to a high-quality image consistent to  $\mathbf{y}$ . We propose a heuristic approximation of it:

$$p(\mathbf{y} | \mathbf{x}_t) = \frac{1}{Z} \exp(-[s\mathcal{L}(\mathcal{D}(\mathbf{x}_t), \mathbf{y}) + \lambda\mathcal{Q}(\mathbf{x}_t)]) \quad (8)$$

where  $\mathcal{L}$  is some image distance metric,  $Z$  is a normalization factor, and  $s$  is a scaling factor controlling the magnitude of guidance. Intuitively, this definition encourages  $\mathbf{x}_t$  to be consistent with the corrupted image  $\mathbf{y}$  to obtain a high probability of  $p(\mathbf{y} | \mathbf{x}_t)$ .  $\mathcal{Q}$  is the optional quality enhancement loss to enhance the flexibility of GDP, which can be used to control some properties (such as brightness) or enhance the quality of the denoised image.  $\lambda$  is the scale factor for adjusting the quality of images. The gradients of both sides are computed as:

$$\begin{aligned} \log p(\mathbf{y} | \mathbf{x}_t) &= -\log Z - s\mathcal{L}(\mathcal{D}(\mathbf{x}_t), \mathbf{y}) - \lambda\mathcal{Q}(\mathbf{x}_t) \\ \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) &= -s\nabla_{\mathbf{x}_t} \mathcal{L}(\mathcal{D}(\mathbf{x}_t), \mathbf{y}) - \lambda\nabla_{\mathbf{x}_t} \mathcal{Q}(\mathbf{x}_t). \end{aligned} \quad (9)$$

where the distance metric  $\mathcal{L}$  and the optional quality loss  $\mathcal{Q}$  can be found in Sec. 5.

In this way, the conditional transition  $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{y})$  can be approximately obtained through the unconditional transition  $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$  by shifting the mean of the

---

**Algorithm 1: GDP- $\mathbf{x}_t$**  with fixed degradation model: Conditioner guided diffusion sampling on  $\mathbf{x}_t$ , given a diffusion model  $(\mu_{\theta}(\mathbf{x}_t), \Sigma_{\theta}(\mathbf{x}_t))$ , corrupted image conditioner  $\mathbf{y}$ .

---

**Input:** Corrupted image  $\mathbf{y}$ , gradient scale  $s$ , degradation model  $\mathcal{D}$ , distance measure  $\mathcal{L}$ , optional quality enhancement loss  $\mathcal{Q}$ , quality enhancement scale  $\lambda$ .  
**Output:** Output image  $\mathbf{x}_0$  conditioned on  $\mathbf{y}$   
 Sample  $\mathbf{x}_T$  from  $\mathcal{N}(0, \mathbf{I})$   
**for**  $t$  from  $T$  to  $1$  **do**  
      $\mu, \Sigma = \mu_{\theta}(\mathbf{x}_t), \Sigma_{\theta}(\mathbf{x}_t)$   
      $\mathcal{L}_{\mathbf{x}_t}^{total} = \mathcal{L}(\mathbf{y}, \mathcal{D}(\mathbf{x}_t)) + \mathcal{Q}(\mathbf{x}_t)$   
     Sample  $\mathbf{x}_{t-1}$  by  $\mathcal{N}(\mu + s\nabla_{\mathbf{x}_t} \mathcal{L}_{\mathbf{x}_t}^{total}, \Sigma)$   
**end**  
**return**  $\mathbf{x}_0$

---



---

**Algorithm 2: GDP- $\mathbf{x}_0$ :** Conditioner guided diffusion sampling on  $\tilde{\mathbf{x}}_0$ , given a diffusion model  $(\mu_{\theta}(\mathbf{x}_t), \Sigma_{\theta}(\mathbf{x}_t))$ , corrupted image conditioner  $\mathbf{y}$ .

---

**Input:** Corrupted image  $\mathbf{y}$ , gradient scale  $s$ , degradation model  $\mathcal{D}_{\phi}$  with randomly initiated parameters  $\phi$ , learning rate  $l$  for optimizable degradation model, distance measure  $\mathcal{L}$ , optional quality enhancement loss  $\mathcal{Q}$ , quality enhancement scale  $\lambda$ .  
**Output:** Output image  $\mathbf{x}_0$  conditioned on  $\mathbf{y}$   
 Sample  $\mathbf{x}_T$  from  $\mathcal{N}(0, \mathbf{I})$   
**for**  $t$  from  $T$  to  $1$  **do**  
      $\mu, \Sigma = \mu_{\theta}(\mathbf{x}_t), \Sigma_{\theta}(\mathbf{x}_t)$   
      $\tilde{\mathbf{x}}_0 = \frac{\mathbf{x}_t}{\sqrt{\alpha_t}} - \frac{\sqrt{1-\alpha_t}\epsilon_{\theta}(\mathbf{x}_t, t)}{\sqrt{\alpha_t}}$   
      $\mathcal{L}_{\phi, \tilde{\mathbf{x}}_0}^{total} = \mathcal{L}(\mathbf{y}, \mathcal{D}_{\phi}(\tilde{\mathbf{x}}_0)) + \mathcal{Q}(\tilde{\mathbf{x}}_0)$   
      $\phi \leftarrow \phi - l\nabla_{\phi} \mathcal{L}_{\phi, \tilde{\mathbf{x}}_0}^{total}$   
     Sample  $\mathbf{x}_{t-1}$  by  $\mathcal{N}(\mu + s\nabla_{\tilde{\mathbf{x}}_0} \mathcal{L}_{\phi, \tilde{\mathbf{x}}_0}^{total}, \Sigma)$   
**end**  
**return**  $\mathbf{x}_0$

---

unconditional distribution by  $-(s\Sigma\nabla_{\mathbf{x}_t} \mathcal{L}(\mathcal{D}(\mathbf{x}_t), \mathbf{y}) + \lambda\Sigma\nabla_{\mathbf{x}_t} \mathcal{Q}(\mathbf{x}_t))$ . However, we find that the way of adding guidance [3] and the variance  $\Sigma$  negatively influence the reconstructed images.

#### 4.1. Single Image Guidance

The super-resolution, inpainting, colorization, deblurring, and enlightening tasks use single-image guidance.

**The Influence of Variance  $\Sigma$  on the Guidance.** In previous conditional diffusion models [15, 83], the variance  $\Sigma$  is applied for the mean shift in the sampling process, which is theoretically proved in the Appendix. In our work, we find that the variance  $\Sigma$  might exert a negative influence on the quality of the generated images in our experiments. Therefore, we remove the variance during the guided denoising process to improve our performance. With the absence of  $\Sigma$  and the fixed guidance scale  $s$ , the guided denoising process can be controlled by the variable scale  $\hat{s}$ .

**Guidance on  $\mathbf{x}_t$ .** Further, as vividly shown in Fig. 2b, Algo. 1 and Algo. 3 in Appendix, this class of guided dif-

fusion models is the commonly used one [15, 48, 89], where the guidance is conditioned on  $x_t$  but with the absence of  $\Sigma$ , named GDP- $x_t$ . However, this variant that applies the guidance on  $x_t$  may still yield less satisfactory quality images. The intuition is  $x_t$  is a noisy image with a specific noise magnitude, but  $y$  is in general a corrupted image with no noise or noises of different magnitude. We lack reliable ways to define the distance between  $x_t$  and  $y$ . A naive MSE loss or perceptual loss will make  $x_t$  deviate from its original noise magnitude and result in low-quality image generation.

**Guidance on  $\tilde{x}_0$ .** To tackle the problem as mentioned above, we systematically study the conditional signal applied on  $\tilde{x}_0$ . Detailly, in the sampling process, the pre-trained DDPM usually first predicts a clean image  $\tilde{x}_0$  from the noisy image  $x_t$  by estimating the noise in  $x_t$ , which can be directly inferred when given  $x_t$  by the Eq. 6 in every timestep  $t$ . Then the predicted  $\tilde{x}_0$  together with  $x_t$  are utilized to sample the next step latent  $x_{t-1}$ . We can add guidance on this intermediate variable  $\tilde{x}_0$  to control the generation process of the DDPM. The detailed sampling process can be found in Fig. 2c and Algo. 2, where there is only one corrupted image.

**Known Degradation.** Several tasks [24, 37, 38, 93] can be categorized into the class that the degradation function is known. In detail, the degradation model for image deblurring and super-resolution can be formulated as  $y = (x \otimes k) \downarrow_s$ . It assumes the low-resolution (LR) image is obtained by first convolving the high-resolution (HR) image with a Gaussian kernel (or point spread function)  $k$  to get a blurry image  $x \otimes k$ , followed by a down-sampling operation  $\downarrow_s$  with scale factor  $s$ . The goal of image inpainting is to recover the missing pixels of an image. The corresponding degradation transform is to multiply the original image with a binary mask  $m$ :  $\psi(x) = x \odot m$ , where  $\odot$  is Hadamard's product. Further, image colorization aims at restoring a gray-scale image  $y \in \mathbb{R}^{H \times W}$  to a colorful image with RGB channels  $x \in \mathbb{R}^{3 \times H \times W}$ . To obtain  $y$  from the colorful image  $x$ , the degradation transform  $\psi$  is a graying transform that only preserves the brightness of  $x$ .

**Unknown Degradation.** In the real world, many images undergo complicated degradations [98], where the degradation models or the parameters of degradation models are unknown [45, 86]. In this case, the original images and the parameters of degradation models should be estimated simultaneously. For instance, in our work, the low-light image enhancement and the HDR recovery can be regarded as tasks with unknown degradation models. Here, we devise a simple but effective degradation model to simulate the complicated degradation, which can be formulated as follows:

$$y = fx + \mathcal{M}, \quad (10)$$

where the light factor  $f$  is a scalar and the light mask  $\mathcal{M}$  is a vector of the same dimension as  $x$ .  $f$  and  $\mathcal{M}$  are unknown parameters of the degradation model. The reason that we

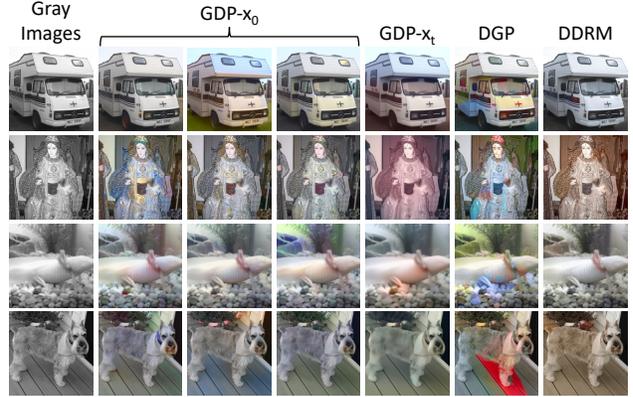


Figure 3. **Qualitative comparison of colorization results on ImageNet validation images.** GDP- $x_0$  generates various samples on the same input.

can use this simple degradation model is that the transform between any pair of corrupted images and the corresponding high-quality image can be captured by  $f$  and  $\mathcal{M}$  as long as they have the same size. If they do not have the same size, we can first resize  $x$  to the same size as  $y$  and then apply this transform. It is worth noting that this degradation model is non-linear in general, since  $f$  and  $\mathcal{M}$  depend on  $x$  and  $y$ . We need to estimate  $f$  and  $\mathcal{M}$  for every individual corrupted image. We achieve this by randomly initializing them and synchronously optimizing them in the reverse process of DDPMs as shown in Algo. 2.

## 4.2. Extended version

**Multi-images Guidance.** Under specific circumstances, there are several images could be utilized to guide the generation of a single image [59, 105], which is merely studied and much more challenging than single-image guidance. To this end, we propose the HDR-GDP for the HDR image recovery with multiple images as guidance, consisting of three input LDR images, *i.e.* short, medium, and long exposures. Similar to low-light enhancement, the degradation models are also treated as Eq. 10, where the parameters remain unknown that determine the HDR recovery is the blind problem. However, as shown in Fig. 2c and Algo. 5 in Appendix, in the reverse process, there are three corrupted images ( $n = 3$ ) to guide the generation so that **three pairs of blind parameters** for three LDR images are randomly initiated and optimized.

**Restore Any-size Image.** Furthermore, the pre-trained diffusion models provide by [15] with the size of 256 are only able to generate the fixed size of images, while the sizes of images from various image restoration are diverse. Herein, we employ the patch-based method as [47] to tackle this problem. By the merits of this patch-based strategy (Fig. 13 and Algo. 6 in the Appendix), GDP can be extended to recover the images of arbitrary resolution to promote the versatility of the GDP.

Table 2. **Quantitative comparison of linear image restoration tasks on ImageNet 1k [62].** GDP outperforms other methods in terms of FID and Consistency across all tasks.

Method	4× Super-resolution				Deblur				25% Inpainting				Colorization			
	PSNR ↑	SSIM ↑	Consistency ↓	FID ↓	PSNR ↑	SSIM ↑	Consistency ↓	FID ↓	PSNR ↑	SSIM ↑	Consistency ↓	FID ↓	PSNR ↑	SSIM ↑	Consistency ↓	FID ↓
DGP [62]	21.65	0.56	158.74	152.85	26.00	0.54	475.10	136.53	27.59	0.82	414.60	60.65	18.42	0.71	305.59	94.59
SNIPS [33]	22.38	0.66	21.38	154.43	24.73	0.69	60.11	17.11	17.55	0.74	587.90	103.50	-	-	-	-
RED [69]	24.18	0.71	27.57	98.30	21.30	0.58	63.20	69.55	-	-	-	-	-	-	-	-
DDRM [32]	<b>26.53</b>	<b>0.78</b>	19.39	40.75	<b>35.64</b>	<b>0.98</b>	50.24	4.78	34.28	0.95	<b>4.08</b>	24.09	<b>22.12</b>	0.91	37.33	47.05
GDP- $x_t$	24.27	0.67	80.32	64.67	25.86	0.75	54.08	5.00	31.06	0.93	8.80	20.24	21.30	0.86	75.24	66.43
GDP- $x_0$	24.42	0.68	<b>6.49</b>	<b>38.24</b>	25.98	0.75	<b>41.27</b>	<b>2.44</b>	<b>34.40</b>	<b>0.96</b>	5.29	<b>16.58</b>	21.41	<b>0.92</b>	<b>36.92</b>	<b>37.60</b>

Table 3. **Quantitative comparison of image enlighten task on LOL [88], VE-LOL-L [47], and LoLi-phone [41] benchmarks.** Bold font indicates the best performance in zero-shot learning, and the underlined font denotes the best results in all models.

Learning	Methods	LOL [88]					VE-LOL-L [47]					LoLi-Phone [41]	
		PSNR ↑	SSIM ↑	FID ↓	LOE ↓	PI ↓	PSNR ↑	SSIM ↑	FID ↓	LOE ↓	PI ↓	LOE ↓	PI ↓
Supervised learning	LLNet [50]	<u>17.91</u>	0.76	169.20	384.21	<u>4.10</u>	17.38	0.73	124.98	291.59	<u>5.54</u>	343.34	<u>5.36</u>
	LightenNet [43]	10.29	0.45	90.91	273.21	7.09	13.26	0.57	82.26	199.45	7.29	500.22	6.63
	Retinex-Net [88]	17.24	0.55	129.99	513.28	8.63	16.41	0.64	135.20	421.41	8.62	542.29	8.23
	MBLLEN [52]	17.90	0.77	122.69	175.10	8.39	15.95	0.70	105.74	114.91	7.45	137.34	6.46
	KinD [104]	17.57	0.82	<u>74.52</u>	377.59	7.41	18.07	0.78	80.12	253.79	7.51	265.47	6.84
	KinD++ [102]	17.60	0.80	100.15	712.12	7.96	16.80	0.74	101.23	421.97	7.98	382.51	7.71
	TBFEN [51]	17.25	<u>0.83</u>	90.59	367.66	8.29	<u>18.91</u>	<u>0.81</u>	91.30	276.65	8.02	214.30	7.34
	DSLRL [46]	14.98	0.67	183.92	272.68	7.09	15.70	0.68	124.80	271.63	7.27	281.25	6.99
Unsupervised learning	EnlightenGAN [29]	17.44	0.74	82.60	379.23	8.78	17.45	0.75	86.51	311.85	8.27	373.41	7.26
Self-supervised learning	DRBN [92]	15.15	0.52	94.96	692.99	5.53	18.47	0.78	88.10	268.70	6.15	285.06	5.31
Zero-shot learning	ExCNet [99]	<b>16.04</b>	0.62	111.18	220.38	8.70	16.20	0.66	115.24	225.15	8.62	359.96	7.95
	Zero-DCE [23]	14.91	<b>0.70</b>	81.11	245.54	8.84	<b>17.84</b>	<b>0.73</b>	85.72	194.10	8.12	214.30	7.34
	Zero-DCE++ [42]	14.86	0.62	86.22	302.06	7.08	16.12	0.45	86.96	313.50	7.92	308.15	7.18
	RRDNet [106]	11.37	0.53	89.09	127.22	8.17	13.99	0.58	83.41	94.23	7.36	92.73	7.20
	GDP- $x_t$	7.32	0.57	238.92	364.15	8.26	9.45	0.50	152.68	194.49	7.12	508.73	8.06
	GDP- $x_0$	13.93	0.63	<b>75.16</b>	<b>110.39</b>	<b>6.47</b>	13.04	0.55	<b>78.74</b>	<b>79.08</b>	<b>6.47</b>	<b>75.29</b>	<b>6.35</b>

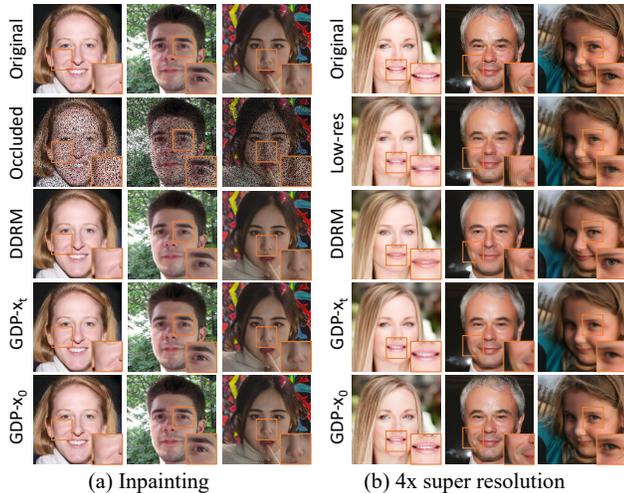


Figure 4. **Qualitative results of (a) 25 % inpainting and (b) 4 × super-resolution on CelebA [31].**

## 5. Loss Function

In GDP, the loss function can be divided into two main parts: Reconstruction loss and quality enhancement loss, where the former aims to recover the information contained in the conditional signal while the latter is integrated to promote the quality of the final outputs.

**Reconstruction Loss.** The reconstruction loss can be MSE, structural similarity index measure (SSIM), perceptual loss,



Figure 5. **Results of image deblurring task on  $256 \times 256$  USC-SIPI images [87] using an ImageNet model.**

or other reconstructive loss. Here, we primarily choose MSE loss as our reconstruction loss.

**Quality Enhancement Loss.** 1) *Exposure Control Loss:* To enhance the versatility of GDP, an exposure control loss  $L_{exp}$  [23] is employed to control the exposure level for low-light image enhancement, which is written as:

$$L_{exp} = \frac{1}{U} \sum_{k=1}^U |R_k - E|, \quad (11)$$

where  $U$  stands for the number of non-overlapping local regions of size  $8 \times 8$ , and  $R$  represents the average intensity value of a local region in the reconstructed image. Follow-

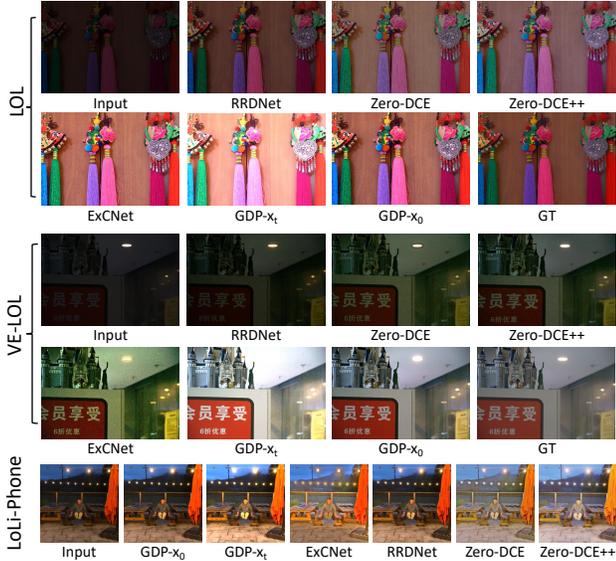


Figure 6. **Qualitative results of low-light enhancement on the LOL [88], VE-LOL [47], and LoLi-Phone [41] datasets.**

ing the previous works [55, 56],  $E$  is set as the gray level in the RGB color space. As expected,  $E$  can be adjusted to control the brightness in our experiments.

2) *Color Constancy Loss*: Following the Gray-World color constancy hypothesis [6], a color constancy loss  $L_{\text{col}}$  is exploited to correct the potential color deviations in the restored image and bridge the relations among the three adjusted channels in the colorization task, formulated as:

$$L_{\text{col}} = \sum_{\forall(m,n) \in \varepsilon} (Y^m - Y^n)^2, \varepsilon = \{(R, G), (R, B), (G, B)\} \quad (12)$$

where  $Y^m$  is the average intensity value of  $m$  channel in the recovered image,  $(m, n)$  represents a pair of channels.

3) *Illumination Smoothness Loss*: To maintain the monotonicity relations between neighboring pixels in the optimized light mask  $\mathcal{M}$ , an illumination smoothness loss [23] is utilized for each light variance  $\mathcal{M}$ . The illumination smoothness loss  $L_{tv, \mathcal{M}}$  is defined as:

$$L_{tv, \mathcal{M}} = \frac{1}{N} \sum_{n=1}^N \sum_{c \in \xi} (|\nabla_h \mathcal{M}_n^c| + |\nabla_v \mathcal{M}_n^c|)^2, \xi = \{R, G, B\}, \quad (13)$$

where  $N$  is iteration times,  $\nabla_h$  and  $\nabla_v$  are the horizontal and vertical gradient operations, respectively.

Specifically, the image colorization task uses color constancy loss to obtain more natural colors. The low-light enhancement requires color constancy loss for the same reason. In addition, the low-light enhancement task uses illumination smoothness loss to make the estimated light mask  $\mathcal{M}$  smoother. Exposure control loss enables us to manually control the brightness of the restored image. The weights of the losses can be found in Appendix.

## 6. Experiments

In this section, we systematically compare GDP, which uses a **single unconditional DDPM pre-trained on Im-**

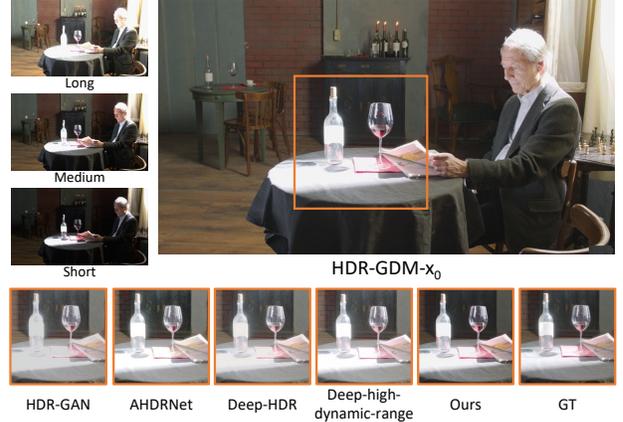


Figure 7. **Example from the NTIRE dataset [63].** We compare a set of patches cropped from the tone-mapped HDR images generated by state-of-the-art methods.

geNet provide by [15], with other methods of various image restoration and enhancement tasks, and ablate the effectiveness of the proposed design. We furthermore list details on implementation, datasets, evaluation, and more qualitative results for all tasks in Appendix.

### 6.1. Linear and Multi-linear Degradation Tasks

Aiming at quantifying the performance of GDP, we focus on the ImageNet dataset for its diversity. For each experiment, we report the average peak signal-to-noise ratio (PSNR), SSIM, and Consistency to measure faithfulness to the original image and the FID to measure the resulting image quality. GDP is compared with other unsupervised methods that can operate on ImageNet, including RED [69], DGP [62], SNIPS [33], and DDRM [32]. We evaluate all methods on the tasks of  $4\times$  super-resolution, deblurring, inpainting, and colorization on one validation set from each of the 1000 ImageNet classes, following [62]. Table 2 shows that GDP- $x_0$  outperforms other methods in Consistency and FID. The only exception is that DDRM achieves better PSNR and SSIM than GDP, but it requires higher Consistency and FID [11, 13, 15, 17, 28, 73]. GDP produces high-quality reconstructions across all the tested datasets and problems, which can be seen in Appendix. As a posterior sampling algorithm, GDP can produce multiple outputs for the same input, as demonstrated in the colorization task in Fig. 3. Moreover, the unconditional ImageNet DDPMs can be used to solve inverse problems on out-of-distribution images with general content. In Figs. 4 and 5, and more illustrations in Appendix, we show GDP successfully restores  $256 \times 256$  images from USC-SIPI [87], LSUN [94], and CelebA [31], which do not necessarily belong to any ImageNet class. GDP can also restore the images under multi-degradation (Fig. 1 and Appendix).

Table 4. Quantitative comparison on the NTIRE dataset [63].

Methods	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
AHDRNet [91]	18.72	0.58	0.39	81.98
HDR-GAN [59]	21.67	0.74	0.26	52.71
Deep-HDR [90]	21.66	0.76	0.26	57.52
Deep-high-dynamic-range [30]	21.33	0.71	0.26	51.92
GDP- $x_t$	19.36	0.65	0.30	63.89
GDP- $x_0$	<b>24.88</b>	<b>0.86</b>	<b>0.13</b>	<b>50.05</b>

## 6.2. Exposure Correction Tasks

Encouraged by the excellent performance on the linear inverse problem, we further evaluate our GDP on the low-light image enhancement, which is categorized into non-linear and blind issues. Following the previous works [41], the three datasets LOL [88], VE-LOL-L [47], and the most challenging LoLi-phone [41] are leveraged to test the capability of GDP on low-light enhancement. As shown in Table 3, our GDP- $x_0$  fulfills the best FID, lightness order error (LOE) [85], and perceptual index (PI) [57] across all the zero-shot methods under three datasets. The lower LOE demonstrates better preservation for the naturalness of lightness, while the lower PI indicates better perceptual quality. In Fig. 6 and Appendix, our GDP- $x_0$  yields the most reasonable and satisfactory results across all methods. For more control, by the merits of Exposure Control Loss, the brightness of the generated images can be adjusted by the well-exposedness level  $E$  (Fig. 1 and Appendix)

## 6.3. HDR Image Recovery

To evaluate our model on the HDR recovery [41], we compare HDR-GDP- $x_0$  with the state-of-the-art HDR methods on the test images in the HDR dataset from the NTIRE2021 Multi-Frame HDR Challenge [63], from which we randomly select 100 different scenes as the validation. Each scene consists of three LDR images with various exposures and corresponding HDR ground truth. The state-of-the-art methods used for comparison include AHDRNet [91], HDR-GAN [59], DeepHDR [90] and deep-high-dynamic-range [30]. The quantitative results are provided in Table 4, where HDR-GDP- $x_0$  performs best in PSNR, SSIM, LPIPS, and FID. As shown in Fig. 7 and Appendix, HDR-GDP- $x_0$  achieves a better quality of reconstructed images, where the low-light parts can be enhanced, and the over-exposure regions are adjusted. Moreover, HDR-GDP- $x_0$  recovers the HDR images with more clear details.

## 6.4. Ablation Study

**The Effectiveness of the Variance  $\Sigma$  and the Guidance Protocol.** The ablation studies on the variance  $\Sigma$  and two ways of guidance are performed to unveil their effectiveness. As shown in Table 5, the performance of GDP- $x_t$  and GDP- $x_0$  is superior to GDP- $x_t$  with  $\Sigma$  and GDP- $x_0$  with  $\Sigma$ , respectively, verifying the absence of variance  $\Sigma$  can yield

Table 5. The ablation study on the variance  $\Sigma$  and the way of the guidance.

Task	4 $\times$ Super resolution				Deblur			
	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID
GDP- $x_t$	22.86	0.60	88.37	68.04	22.06	0.57	69.46	80.39
with $\Sigma$								
GDP- $x_0$	22.09	0.58	93.19	41.22	23.49	0.65	68.67	50.29
with $\Sigma$								
GDP- $x_t$	24.27	0.67	80.32	64.67	25.86	0.73	54.08	5.00
GDP- $x_0$	<b>24.42</b>	<b>0.68</b>	<b>6.49</b>	<b>38.24</b>	<b>25.98</b>	<b>0.75</b>	<b>41.27</b>	<b>2.44</b>

Task	25% Inpainting				Colorization			
	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID
GDP- $x_t$	25.28	0.70	171.44	73.32	17.67	0.70	246.26	145.20
with $\Sigma$								
GDP- $x_0$	24.58	0.75	65.59	22.77	21.28	0.91	66.57	38.39
with $\Sigma$								
GDP- $x_t$	31.06	0.93	8.80	20.24	21.30	0.86	75.24	66.43
GDP- $x_0$	<b>34.40</b>	<b>0.96</b>	<b>5.29</b>	<b>16.58</b>	<b>21.41</b>	<b>0.92</b>	<b>36.92</b>	<b>37.60</b>

Table 6. The ablation study on the optimizable degradation and patch-based tactic.

Methods	LOL					NTIRE			
	PSNR	SSIM	FID	LOE	PI	PSNR	SSIM	LPIPS	FID
Model A	11.05	0.49	156.51	707.57	8.61	24.12	0.67	0.32	86.69
Model B	9.01	0.37	355.99	969.89	9.04	9.83	0.04	1.02	253.11
GDP- $x_t$	7.32	0.57	238.92	364.15	8.26	19.36	0.65	0.30	63.89
GDP- $x_0$	<b>13.93</b>	<b>0.63</b>	<b>75.16</b>	<b>110.39</b>	<b>6.47</b>	<b>24.88</b>	<b>0.86</b>	<b>0.13</b>	<b>50.05</b>

better quality of images. Moreover, the results of GDP- $x_0$  and GDP- $x_0$  with  $\Sigma$  are better than GDP- $x_t$  and GDP- $x_t$  with  $\Sigma$ , respectively, demonstrating the superiority of the guidance on  $x_0$  protocol.

**The Effectiveness of the Trainable Degradation and the Patch-based Tactic.** Moreover, to validate the influence of trainable parameters of the degradation model and our patch-based methods, further experiments are carried out on the LOL [88] and NTIRE [63] datasets. Model A is devised to naively restore the images from patches and patches where the parameters are not related. ModelB is designed with fixed parameters for all patches in the images. As shown in Table 6, our GDP- $x_0$  ranks first across all models and obtains the best visualization results (Fig. 7 and Appendix), revealing the strength of our proposed hierarchical guidance and patch-based method.

## 7. Conclusion

In this paper, we propose the Generative Diffusion Prior for unified image restoration that can be employed to tackle the linear inverse, non-linear and blind problems. Our GDP is able to restore any-size images via hierarchical guidance and patch-based methods. We systematically studied the way of guidance to exploit the strength of the DDPM. The GDP is comprehensively utilized on various tasks such as super-resolution, deblurring, inpainting, colorization, low-light enhancement, and HDR recovery, demonstrating the capabilities of GDP on unified image restoration.

**Acknowledgement.** This project is funded in part by Shanghai AI Laboratory

## References

- [1] Muhammad Asim, Fahad Shamshad, and Ali Ahmed. Blind image deconvolution using deep generative priors. *IEEE Transactions on Computational Imaging*, 6:1493–1506, 2020. 1, 2
- [2] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarrow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021. 2, 3
- [3] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. 4
- [4] Georgios Batzolis, Jan Stanczuk, Carola-Bibiane Schönlieb, and Christian Etmann. Conditional image generation with score-based diffusion models. *arXiv preprint arXiv:2111.13606*, 2021. 4
- [5] Isabelle Begin and FR Ferrie. Blind super-resolution using a learning-based approach. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, volume 2, pages 85–89. IEEE, 2004. 2
- [6] Gershon Buchsbaum. A spatial processor model for object colour perception. *Journal of the Franklin institute*, 310(1):1–26, 1980. 7
- [7] Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snavely, and Bharath Hariharan. Learning gradient fields for shape generation. In *European Conference on Computer Vision*, pages 364–381. Springer, 2020. 2
- [8] Jingwen Chen, Jiawei Chen, Hongyang Chao, and Ming Yang. Image blind denoising with generative adversarial network based noise modeling. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3155–3164, 2018. 2
- [9] Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*, 2020. 4
- [10] Xiangyu Chen, Yihao Liu, Zhengwen Zhang, Yu Qiao, and Chao Dong. Hdrunet: Single image hdr reconstruction with denoising and dequantization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 354–363, 2021. 1
- [11] Yu Chen, Ying Tai, Xiaoming Liu, Chunhua Shen, and Jian Yang. Fsrnet: End-to-end learning face super-resolution with facial priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2492–2501, 2018. 7, 19
- [12] Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11472–11481, 2022. 4
- [13] Ryan Dahl, Mohammad Norouzi, and Jonathon Shlens. Pixel recursive super resolution. In *Proceedings of the IEEE international conference on computer vision*, pages 5439–5448, 2017. 7, 19
- [14] 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. 2
- [15] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 2, 4, 5, 7
- [16] 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. 1
- [17] Alexey Dosovitskiy and Thomas Brox. Generating images with perceptual similarity metrics based on deep networks. *Advances in neural information processing systems*, 29, 2016. 7, 19
- [18] Wenchao Du, Hu Chen, and Hongyu Yang. Learning invariant representation for unsupervised image restoration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14483–14492, 2020. 3
- [19] Majed El Helou and Sabine Süsstrunk. Bigprior: Toward decoupling learned prior hallucination and data fidelity in image restoration. *IEEE Transactions on Image Processing*, 31:1628–1640, 2022. 2
- [20] 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. 2
- [21] 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
- [22] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022. 3
- [23] Chunle Guo, Chongyi Li, Jichang Guo, Chen Change Loy, Junhui Hou, Sam Kwong, and Runmin Cong. Zero-reference deep curve estimation for low-light image enhancement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1780–1789, 2020. 6, 7
- [24] Muhammad Haris, Gregory Shakhnarovich, and Norimichi Ukita. Deep back-projection networks for super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1664–1673, 2018. 1, 5
- [25] Yu He, Kim-Hui Yap, Li Chen, and Lap-Pui Chau. A soft map framework for blind super-resolution image reconstruction. *Image and Vision Computing*, 27(4):364–373, 2009. 2
- [26] 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 equilib-

- rium. *Advances in neural information processing systems*, 30, 2017. 16
- [27] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2, 3
- [28] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, 23:47–1, 2022. 7
- [29] Yifan Jiang, Xinyu Gong, Ding Liu, Yu Cheng, Chen Fang, Xiaohui Shen, Jianchao Yang, Pan Zhou, and Zhangyang Wang. Enlighten: Deep light enhancement without paired supervision. *IEEE Transactions on Image Processing*, 30:2340–2349, 2021. 2, 6
- [30] Nima Khademi Kalantari, Ravi Ramamoorthi, et al. Deep high dynamic range imaging of dynamic scenes. *ACM Trans. Graph.*, 36(4):144–1, 2017. 8, 23
- [31] 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, 6, 7
- [32] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *arXiv preprint arXiv:2201.11793*, 2022. 2, 4, 6, 7, 16, 21
- [33] Bahjat Kawar, Gregory Vaksman, and Michael Elad. Snips: Solving noisy inverse problems stochastically. *Advances in Neural Information Processing Systems*, 34:21757–21769, 2021. 4, 6, 7
- [34] Rihuan Ke and Carola-Bibiane Schönlieb. Unsupervised image restoration using partially linear denoisers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5796–5812, 2021. 2
- [35] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1637–1645, 2016. 2
- [36] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021. 2, 3
- [37] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8878–8887, 2019. 1, 5
- [38] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Learning representations for automatic colorization. In *European conference on computer vision*, pages 577–593. Springer, 2016. 1, 5
- [39] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4681–4690, 2017. 1
- [40] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5549–5558, 2020. 16
- [41] Chongyi Li, Chunle Guo, Ling-Hao Han, Jun Jiang, Ming-Ming Cheng, Jinwei Gu, and Chen Change Loy. Low-light image and video enhancement using deep learning: A survey. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 44(01):1–1, 2021. 1, 2, 6, 7, 8, 16
- [42] Chongyi Li, Chunle Guo, and Chen Change Loy. Learning to enhance low-light image via zero-reference deep curve estimation. *arXiv preprint arXiv:2103.00860*, 2021. 6
- [43] Chongyi Li, Jichang Guo, Fatih Porikli, and Yanwei Pang. Lightnet: A convolutional neural network for weakly illuminated image enhancement. *Pattern recognition letters*, 104:15–22, 2018. 6
- [44] Jingyun Liang, Jiezhong 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
- [45] 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. 5
- [46] Seokjae Lim and Wonjun Kim. Dslr: deep stacked laplacian restorer for low-light image enhancement. *IEEE Transactions on Multimedia*, 23:4272–4284, 2020. 6
- [47] Jiaying Liu, DeJia Xu, Wenhan Yang, Minhao Fan, and Haofeng Huang. Benchmarking low-light image enhancement and beyond. *International Journal of Computer Vision*, 129(4):1153–1184, 2021. 5, 6, 7, 8, 16
- [48] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu, Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with semantic diffusion guidance. *arXiv preprint arXiv:2112.05744*, 2021. 5
- [49] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015. 16
- [50] Kin Gwn Lore, Adedotun Akintayo, and Soumik Sarkar. Llnet: A deep autoencoder approach to natural low-light image enhancement. *Pattern Recognition*, 61:650–662, 2017. 6
- [51] Kun Lu and Lihong Zhang. Tbefn: A two-branch exposure-fusion network for low-light image enhancement. *IEEE Transactions on Multimedia*, 23:4093–4105, 2020. 6
- [52] Feifan Lv, Feng Lu, Jianhua Wu, and Chongsoon Lim. Mblen: Low-light image/video enhancement using cnns. In *BMVC*, volume 220, page 4, 2018. 6
- [53] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 2

- [54] 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](#)
- [55] Tom Mertens, Jan Kautz, and Frank Van Reeth. Exposure fusion. In *15th Pacific Conference on Computer Graphics and Applications (PG'07)*, pages 382–390. IEEE, 2007. [7](#)
- [56] Tom Mertens, Jan Kautz, and Frank Van Reeth. Exposure fusion: A simple and practical alternative to high dynamic range photography. In *Computer graphics forum*, volume 28, pages 161–171. Wiley Online Library, 2009. [7](#)
- [57] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Making a “completely blind” image quality analyzer. *IEEE Signal processing letters*, 20(3):209–212, 2012. [8](#), [17](#)
- [58] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. [3](#), [19](#)
- [59] Yuzhen Niu, Jianbin Wu, Wenxi Liu, Wenzhong Guo, and Rynson WH Lau. Hdr-gan: Hdr image reconstruction from multi-exposed ldr images with large motions. *IEEE Transactions on Image Processing*, 30:3885–3896, 2021. [2](#), [5](#), [8](#), [23](#)
- [60] Gregory Ongie, Ajil Jalal, Christopher A Metzler, Richard G Baraniuk, Alexandros G Dimakis, and Rebecca Willett. Deep learning techniques for inverse problems in imaging. *IEEE Journal on Selected Areas in Information Theory*, 1(1):39–56, 2020. [1](#)
- [61] Ozan Özdenizci and Robert Legenstein. Restoring vision in adverse weather conditions with patch-based denoising diffusion models. *arXiv preprint arXiv:2207.14626*, 2022. [18](#)
- [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*, 2021. [2](#), [4](#), [6](#), [7](#), [16](#), [19](#)
- [63] Eduardo Pérez-Pellitero, Sibi Catley-Chandar, Ales Leonardis, and Radu Timofte. Ntire 2021 challenge on high dynamic range imaging: Dataset, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 691–700, 2021. [2](#), [7](#), [8](#)
- [64] Eduardo Pérez-Pellitero, Sibi Catley-Chandar, Ales Leonardis, and Radu Timofte. Ntire 2021 challenge on high dynamic range imaging: Dataset, methods and results. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 691–700, 2021. [16](#), [19](#)
- [65] Guocheng Qian, Jinjin Gu, Jimmy S Ren, Chao Dong, Furong Zhao, and Juan Lin. Trinity of pixel enhancement: a joint solution for demosaicking, denoising and super-resolution. *arXiv preprint arXiv:1905.02538*, 1(3):4, 2019. [2](#)
- [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*, 2022. [2](#)
- [67] 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. [2](#)
- [68] Renu M Rameshan, Subhasis Chaudhuri, and Rajbabu Velmurugan. High dynamic range imaging under noisy observations. In *2011 18th IEEE International Conference on Image Processing*, pages 1333–1336. IEEE, 2011. [2](#)
- [69] Yaniv Romano, Michael Elad, and Peyman Milanfar. The little engine that could: Regularization by denoising (red). *SIAM Journal on Imaging Sciences*, 10(4):1804–1844, 2017. [4](#), [6](#), [7](#)
- [70] 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. [4](#)
- [71] Dohoon Ryu and Jong Chul Ye. Pyramidal denoising diffusion probabilistic models. *arXiv preprint arXiv:2208.01864*, 2022. [3](#)
- [72] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. [2](#), [15](#)
- [73] 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*, 2022. [4](#), [7](#), [17](#), [19](#)
- [74] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. Singan: Learning a generative model from a single natural image. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4570–4580, 2019. [2](#)
- [75] Assaf Shocher, Nadav Cohen, and Michal Irani. “zero-shot” super-resolution using deep internal learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3118–3126, 2018. [3](#)
- [76] 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](#), [4](#), [20](#)
- [77] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. [2](#), [19](#)
- [78] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020. [2](#)
- [79] 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](#)
- [80] Maitreya Suin, Kuldeep Purohit, and AN Rajagopalan. Spatially-attentive patch-hierarchical network for adaptive

- motion deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3606–3615, 2020. **1**
- [81] Anwaar Ulhaq, Naveed Akhtar, and Ganna Pogrebna. Efficient diffusion models for vision: A survey. *arXiv preprint arXiv:2210.09292*, 2022. **2**
- [82] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021. **2**
- [83] Jinyi Wang, Zhaoyang Lyu, Dahua Lin, Bo Dai, and Hongfei Fu. Guided diffusion model for adversarial purification. *arXiv preprint arXiv:2205.14969*, 2022. **4**
- [84] Lin Wang and Kuk-Jin Yoon. Deep learning for hdr imaging: State-of-the-art and future trends. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. **1**
- [85] Shuhang Wang, Jin Zheng, Hai-Miao Hu, and Bo Li. Naturalness preserved enhancement algorithm for non-uniform illumination images. *IEEE transactions on image processing*, 22(9):3538–3548, 2013. **8, 17**
- [86] 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. **5**
- [87] Allan G Weber. The usc-sipi image database: Version 5. <http://sipi.usc.edu/database/>, 2006. **6, 7, 16**
- [88] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. *arXiv preprint arXiv:1808.04560*, 2018. **6, 7, 8, 16, 19**
- [89] Quanlin Wu, Hang Ye, and Yuntian Gu. Guided diffusion model for adversarial purification from random noise. *arXiv preprint arXiv:2206.10875*, 2022. **3, 5**
- [90] Shangzhe Wu, Jiarui Xu, Yu-Wing Tai, and Chi-Keung Tang. Deep high dynamic range imaging with large foreground motions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 117–132, 2018. **8, 23**
- [91] Qingsen Yan, Dong Gong, Qinfeng Shi, Anton van den Hengel, Chunhua Shen, Ian Reid, and Yanning Zhang. Attention-guided network for ghost-free high dynamic range imaging. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1751–1760, 2019. **8, 23**
- [92] Wenhan Yang, Shiqi Wang, Yuming Fang, Yue Wang, and Jiaying Liu. From fidelity to perceptual quality: A semi-supervised approach for low-light image enhancement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3063–3072, 2020. **6**
- [93] Raymond A Yeh, Chen Chen, Teck Yian Lim, Alexander G Schwing, Mark Hasegawa-Johnson, and Minh N Do. Semantic image inpainting with deep generative models. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5485–5493, 2017. **1, 5**
- [94] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. **2, 7, 16**
- [95] Ke Yu, Chao Dong, Liang Lin, and Chen Change Loy. Crafting a toolchain for image restoration by deep reinforcement learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2443–2452, 2018. **2**
- [96] Yuan Yuan, Siyuan Liu, Jiawei Zhang, Yongbing Zhang, Chao Dong, and Liang Lin. Unsupervised image super-resolution using cycle-in-cycle generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 701–710, 2018. **3**
- [97] Junzhe Zhang, Xinyi Chen, Zhongang Cai, Liang Pan, Haiyu Zhao, Shuai Yi, Chai Kiat Yeo, Bo Dai, and Chen Change Loy. Unsupervised 3d shape completion through gan inversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1768–1777, 2021. **15**
- [98] 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. **5**
- [99] Lin Zhang, Lijun Zhang, Xiao Liu, Ying Shen, Shaoming Zhang, and Shengjie Zhao. Zero-shot restoration of back-lit images using deep internal learning. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 1623–1631, 2019. **6**
- [100] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European conference on computer vision*, pages 649–666. Springer, 2016. **1**
- [101] 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. **17**
- [102] Yonghua Zhang, Xiaojie Guo, Jiayi Ma, Wei Liu, and Jiawan Zhang. Beyond brightening low-light images. *International Journal of Computer Vision*, 129(4):1013–1037, 2021. **6**
- [103] Yongbing Zhang, Siyuan Liu, Chao Dong, Xinfeng Zhang, and Yuan Yuan. Multiple cycle-in-cycle generative adversarial networks for unsupervised image super-resolution. *IEEE transactions on Image Processing*, 29:1101–1112, 2019. **3**
- [104] Yonghua Zhang, Jiawan Zhang, and Xiaojie Guo. Kindling the darkness: A practical low-light image enhancer. In *Proceedings of the 27th ACM international conference on multimedia*, pages 1632–1640, 2019. **6**
- [105] Zhuoran Zheng, Wenqi Ren, Xiaochun Cao, Xiaobin Hu, Tao Wang, Fenglong Song, and Xiuyi Jia. Ultra-high-definition image dehazing via multi-guided bilateral learning. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16180–16189. IEEE, 2021. **5**
- [106] Anqi Zhu, Lin Zhang, Ying Shen, Yong Ma, Shengjie Zhao, and Yicong Zhou. Zero-shot restoration of underexposed images via robust retinex decomposition. In *2020 IEEE*

*International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2020. 6

- [107] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 3

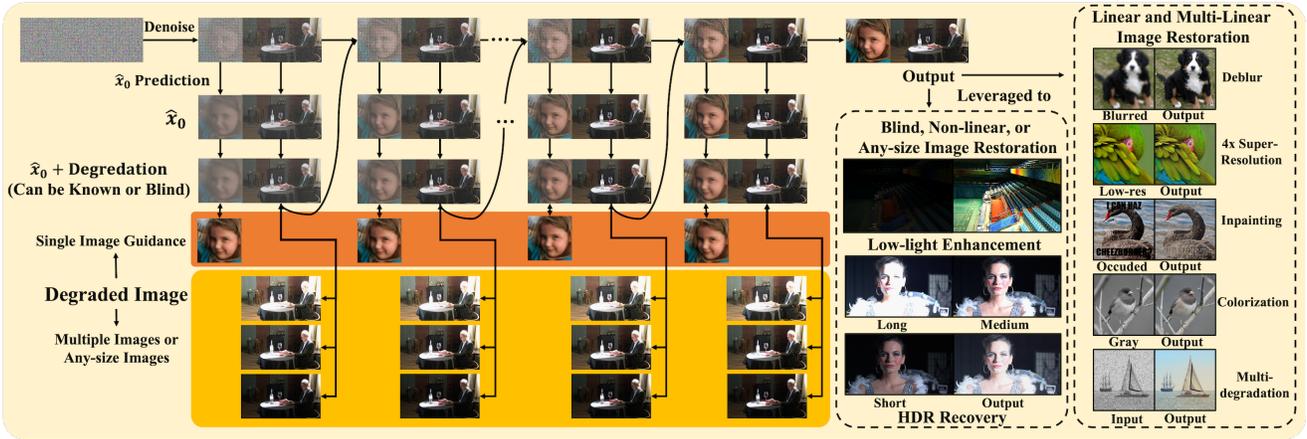


Figure 8. **Illustration of our GDP method for unified image recovery**, including linear inverse problems (Deblurring,  $4\times$  super-resolution, inpainting, and colorization), multi-degradation (*i.e.* Colorization + inpainting), non-linear and blind problems (Low-light enhancement and HDR recovery). Note that GDP can restore images of arbitrary sizes, and can accept multiple low-quality images as guidance as in the case of HDR recovery. GDP fulfills all the tasks using a single unconditional DDPM pre-trained on ImageNet.

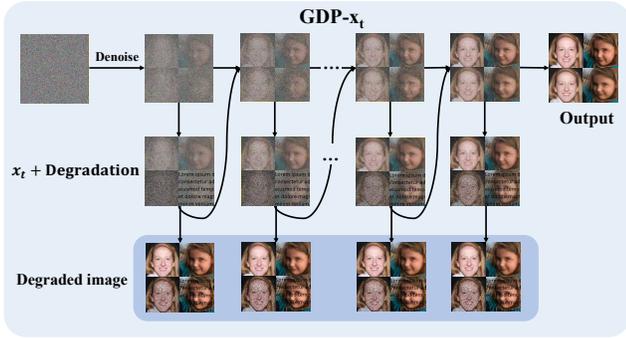


Figure 9. **Overview of the GDP- $x_t$** . The guidance will be added on the noisy image  $x_t$  in every time step.

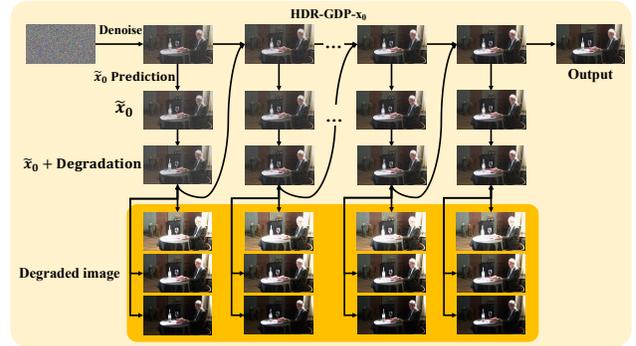


Figure 11. **Overview of the HDR-GDP- $x_0$** . The guidance will also be applied to a clean image  $\tilde{x}_0$ . Unlike the GDP- $x_0$ , three degraded images are utilized to guide the reverse process, and three sets of degradation models are optimized along the reverse process.

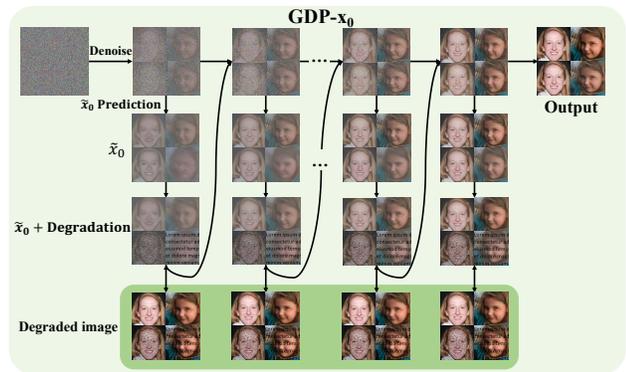


Figure 10. **Overview of the GDP- $x_0$** . The guidance will be applied to a clean image  $\tilde{x}_0$  predicted from the noisy image  $x_t$ .

## A. Limitations and Future works

**Limitations.** The main limitation of our work is its inference time. Since we might add several guidance steps in every time step  $t$ , the sampling time is extended. This limits the applicability of our method to real-time applications and weak end-user devices such as mobile devices. To address this issue, further research into accelerated diffusion sampling techniques is required.

In addition, the choice of the guidance scale is also obtained through experiments, which means that for samples with different distributions, it is necessary to manually select the optimal guidance scale. However, we found that for the same distribution of data, an approximate degradation model may lead to close guidance scales. This phenomenon

---

**Algorithm 3: GDP- $x_t$ :** Conditioner guided diffusion sampling on  $x_t$ , given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , corrupted image conditioner  $\mathbf{y}$ .

---

**Input:** Corrupted image  $\mathbf{y}$ , gradient scale  $s$ , degradation model  $\mathcal{D}_\phi$  with randomly initiated parameters  $\phi$ , learning rate  $l$  for optimizable degradation model, distance measure  $\mathcal{L}$ .

**Output:** Output image  $x_0$  conditioned on  $\mathbf{y}$   
Sample  $x_T$  from  $\mathcal{N}(0, \mathbf{I})$

**for**  $t$  from  $T$  to  $1$  **do**

$\mu, \Sigma = \mu_\theta(x_t), \Sigma_\theta(x_t)$   
 $\mathcal{L}_{\phi, x_t}^{total} = \mathcal{L}(\mathbf{y}, \mathcal{D}_\phi(x_t)) + \mathcal{Q}(x_t)$   
 $\phi \leftarrow \phi - l \nabla_\phi \mathcal{L}_{\phi, x_t}^{total}$   
 Sample  $x_{t-1}$  by  $\mathcal{N}(\mu + s \nabla_{x_t} \mathcal{L}_{\phi, x_t}^{total}, \Sigma)$

**end**

**return**  $x_0$

---

**Algorithm 4: GDP- $x_0$ :** Conditioner guided diffusion sampling on  $\tilde{x}_0$ , given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , corrupted image conditioner  $\mathbf{y}$ .

---

**Input:** Corrupted image  $\mathbf{y}$ , gradient scale  $s$ , degradation model  $\mathcal{D}$ , distance measure  $\mathcal{L}$ .

**Output:** Output image  $x_0$  conditioned on  $\mathbf{y}$   
Sample  $x_T$  from  $\mathcal{N}(0, \mathbf{I})$

**for**  $t$  from  $T$  to  $1$  **do**

$\mu, \Sigma = \mu_\theta(x_t), \Sigma_\theta(x_t)$   
 $\tilde{x}_0 = \frac{x_t}{\sqrt{\alpha_t}} - \frac{\sqrt{1-\alpha_t} \epsilon_\theta(x_t, t)}{\sqrt{\alpha_t}}$   
 $\mathcal{L}_{\tilde{x}_0}^{total} = \mathcal{L}(\mathbf{y}, \mathcal{D}(\tilde{x}_0)) + \mathcal{Q}(\tilde{x}_0)$   
 Sample  $x_{t-1}$  by  $\mathcal{N}(\mu + s \nabla_{\tilde{x}_0} \mathcal{L}_{\tilde{x}_0}^{total}, \Sigma)$

**end**

**return**  $x_0$

---

may be proved mathematically in future work.

**Future works.** In future work, in addition to further optimizing the time step and variance schedules, it would be interesting to investigate the following:

(i) The Guided Diffusion Prior can also theoretically be applied to 3D data restoration. For instance, point cloud completion and upsampling can be regarded as linear inverse problems in 3D vision. Shapeinversion [97] tackles the point cloud completion by GAN inversion, where the GDP can hopefully be integrated.

(ii) Moreover, since LiDAR is affected by various kinds of weather in the real world and also produces various non-linear degradations, GDP should also be explored for the recovery of these point clouds.

(iii) Self-supervised training techniques inspired by our GDP and techniques used in supervised techniques [72]

---

**Algorithm 5: GDP- $x_0$ :** Conditioner guided diffusion sampling on  $x_0$ , given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , corrupted images conditioner  $\{\mathbf{y}^i \mid i = 1, 2, \dots, n\}$ .

---

**Input:** Corrupted image  $\{\mathbf{y}^i \mid i = 1, 2, \dots, n\}$  ( $n = 3$  for HDR recovery (LDR-long image  $\mathbf{y}^1$ , LDR-medium image  $\mathbf{y}^2$ , LDR-short image  $\mathbf{y}^3$ ) and  $n = 1$  for other tasks), gradient scale  $s$ , degradation models  $\{\mathcal{D}_{\phi^i} \mid i = 1, 2, \dots, n\}$  with randomly initiated parameters  $\{\phi^i \mid i = 1, 2, \dots, n\}$ , learning rate  $l$  for optimizable degradation model, distance measure  $\mathcal{L}$ .

**Output:** Output image  $x_0$  conditioned on  $\{\mathbf{y}^i \mid i = 1, 2, \dots, n\}$

Sample  $x_T$  from  $\mathcal{N}(0, \mathbf{I})$

**for**  $t$  from  $T$  to  $1$  **do**

$\mu, \Sigma = \mu_\theta(x_t), \Sigma_\theta(x_t)$   
 $\tilde{x}_0 = \frac{x_t}{\sqrt{\alpha_t}} - \frac{\sqrt{1-\alpha_t} \epsilon_\theta(x_t, t)}{\sqrt{\alpha_t}}$   
 $\mathcal{L}_{\phi, \tilde{x}_0}^{total} = 0$

**for**  $j$  from  $1$  to  $n$  **do**

$\mathcal{L}_{\phi^j, \tilde{x}_0} = \mathcal{L}(\mathbf{y}^j, \mathcal{D}_{\phi^j}(\tilde{x}_0)) + \mathcal{Q}(\tilde{x}_0)$   
 $\phi^j = \phi^j - l \nabla_{\phi^j} \mathcal{L}_{\phi^j, \tilde{x}_0}$   
 $\mathcal{L}_{\phi, \tilde{x}_0}^{total} = \mathcal{L}_{\phi, \tilde{x}_0}^{total} + \mathcal{L}_{\phi^j, \tilde{x}_0}$

**end**

Sample  $x_{t-1}$  by  $\mathcal{N}(\mu + s \nabla_{\tilde{x}_0} \mathcal{L}_{\phi, \tilde{x}_0}^{total}, \Sigma)$

**end**

**return**  $x_0$

---

that further improve the performance of unsupervised image restoration models.

## B. Implementation Details

We apply GDP to a suite of challenging image restoration tasks: (1) **Colorization** transforms an input gray-scale image to a plausible color image. (2) **Inpainting** fills in user-specified masked regions of an image with realistic content. (3) **Super-resolution** extends a low-resolution image into a higher one. (4) **Deblurring** corrects the blurred images, restoring plausible image detail. (5) **Enlightening** enables the dark images turned into normal images. (6) **HDR image recovery** aims to obtain HDR images with the aid of three LDR images. Inputs and outputs of the first four tasks are represented as  $256 \times 256$  RGB images, while the last two tasks are various ( $1900 \times 1060$  for HDR image recovery and  $600 \times 400$  for image enlightening, respectively). We do so without task-specific hyperparameter tuning and architecture customization.

Colorization requires the representation of objects, segmentation, and layouts with long-range image dependen-

---

**Algorithm 6: Restore Any-size Image**

---

**Input:** Conditioner guided diffusion sampling on  $\tilde{\mathbf{x}}_0$ , given a diffusion model  $(\mu_\theta(\mathbf{x}_t), \Sigma_\theta(\mathbf{x}_t))$ , corrupted image conditioner  $\mathbf{y}$ , degradation model  $\mathcal{D}_\phi: \mathbf{y} = f\mathbf{x} + \mathcal{M}$  with randomly initiated parameters  $\phi$ , learning rate  $l$  for optimizable degradation model. Dictionary of  $K$  overlapping patch locations, and a binary patch mask  $\mathbf{P}^k$ .

**Output:** Output image  $\mathbf{x}_0$  conditioned on  $\mathbf{y}$

Sample  $\mathbf{x}_T$  from  $\mathcal{N}(0, \mathbf{I})$

**for**  $t$  from  $T$  to  $1$  **do**

$\mu, \Sigma = \mu_\theta(\mathbf{x}_t), \Sigma_\theta(\mathbf{x}_t)$   
Mean vector  $\Omega_t = \mathbf{0}$  and variance vector  $\psi_t = \mathbf{0}$  and weight vector  $\mathbf{G} = \mathbf{0}$  and  $f = \mathbf{0}$  and  $\mathcal{M} = \mathbf{0}$

**for**  $k = 1, \dots, K$  **do**

$\mathbf{x}_t^k = \text{Crop}(\mathbf{P}^k \circ \mathbf{x}_t)$   
 $\mathbf{y}^k = \text{Crop}(\mathbf{P}^k \circ \mathbf{y})$   
 $\mathcal{M}^k = \text{Crop}(\mathbf{P}^k \circ \mathcal{M})$   
 $\tilde{\mathbf{x}}_0^k = \frac{\mathbf{x}_t^k}{\sqrt{\alpha_t}} - \frac{\sqrt{1-\alpha_t}\epsilon_\theta(\mathbf{x}_t^k, t)}{\sqrt{\alpha_t}}$   
 $\mathcal{L}_{\phi, \tilde{\mathbf{x}}_0^k}^{\text{total}} = \mathcal{L}(\mathbf{y}^k, \mathcal{D}_\phi(\tilde{\mathbf{x}}_0^k)) + \mathcal{Q}(\tilde{\mathbf{x}}_0^k)$   
 $f^k \leftarrow f^k - l \nabla_{f^k} \mathcal{L}_{f^k, \tilde{\mathbf{x}}_0^k}^{\text{total}}$   
 $\mathcal{M}^k \leftarrow \mathcal{M}^k - l \nabla_{\mathcal{M}^k} \mathcal{L}_{\mathcal{M}^k, \tilde{\mathbf{x}}_0^k}^{\text{total}}$   
 $\mu^k = \mu + s \nabla_{\tilde{\mathbf{x}}_0^k} \mathcal{L}_{\phi, \tilde{\mathbf{x}}_0^k}^{\text{total}}$   
 $f = f + f^k$   
 $\Omega_t = \Omega_t + \mathbf{P}^k \cdot \mu^k$   
 $\psi_t = \psi_t + \mathbf{P}^k \cdot \sigma^k$   
 $\mathcal{M} = \mathcal{M} + \mathbf{P}^k \cdot \mathcal{M}^k$   
 $\mathbf{G} = \mathbf{G} + \mathbf{P}^k$

**end**

$\Omega_t = \Omega_t \oslash \mathbf{G}$  //  $\oslash$ : element-wise division  
 $\psi_t = \psi_t \oslash \mathbf{G}$   
 $\mathcal{M} = \mathcal{M} \oslash \mathbf{G}$   
 $f = f / K$   
Sample  $\mathbf{x}_{t-1}$  by  $\mathcal{N}(\Omega_t, \psi_t)$

**end**

**return** Restored any-size image  $\mathbf{x}_0$

---

cies. Inpainting is challenging due to large masks, image diversity, and cluttered scenes. Super-resolution and deblurring are also not trivial because the degradation might damage the content of the images. While the other tasks are linear in nature, low-light enhancement and HDR recovery are non-linear inverse problems; they require a good model of natural image statistics to detect and correct over-exposed and under-exposed areas. Although previous works have studied these problems extensively, it is rare that a model with no task-specific engineering achieves strong performance in all tasks, beating strong task-specific GAN and regression baselines. Our GDP is devised to achieve this goal.

### B.1. Dataset briefs

**ImageNet, LSUN, CelebA, and USC-SIPI Datasets.** To quantitatively evaluate GDP on linear image restoration tasks, we test on 1k images from the ImageNet validation set following [62]. The CelebA-HQ [40] dataset is a high-quality subset of the Large-Scale CelebFaces At-

tributes (CelebA) dataset [49]. LSUN dataset [94] contains around one million labeled images for each of 10 scene categories and 20 object categories. And the USC-SIPI dataset [87] is a collection of various digitized images. We utilize the images from CelebA, LSUN, and USC-SIPI provided by [32].

**LOL Dataset.** The LOL dataset [88] is composed of 500 low-light and normal-light image pairs and divided into 485 training pairs and 15 testing pairs. The low-light images contain noise produced during the photo capture process. Most of the images are indoor scenes. All the images have a resolution of  $400 \times 600$ .

**VE-LOL-L Dataset.** For underexposure correction experiments, we use the paired data of the VE-LOL-L dataset [47], in which each captured well-exposed image has its underexposed version with different underexposure levels. Note that the VE-LOL-L dataset, consisting of VE-LOL-Cap and VE-LOL-Syn, is also carried out. Due to the different distribution of the two sub-set, we solve them under different guidance scales.

**LoLi-Phone Dataset.** LoLi-Phone [41] is a large-scale low-light image and video dataset for low-light image enhancement. The images and videos are taken by different mobile phone cameras under diverse illumination conditions.

**NTIRE Dataset [64].** In the NTIRE dataset, there are 1494 LDRs/HDR for training, 60 images for validation, and 201 images for testing. The 1494 frames consist of 26 long shots. Each scene contains three LDR images, their corresponding exposure and alignment information, and HDR ground truth. The size of an image is  $1060 \times 1900$ . Since the ground truth of the validation and test sets are not available, we only do experiments on the training set. We select 100 images as the test set.

### B.2. Experimental Setup

In each inverse problem, the pixel values are in the range  $[0, 1]$ , and the resulting degradation measures are as follows: (i) For super-resolution, a block averaging filter is utilized to downscale the image on each axis 4 times; (ii) In terms of deblurring, the image is blurred by a  $9 \times 9$  unified kernel. (iii) For colorization, the gray-scale image is the average of the red, green, and blue channels of the original image; (iv) For inpainting, we cover parts of the original image with text overlays or randomly delete 25% pixels.

In the non-linear and blind problem, the images from the low-light dataset and NTIRE dataset are naturally over-exposed or under-exposed. Therefore, no additional operations are required for the images.

### C. Evaluation Metrics

Apart from the commonly used PSNR and SSIM, other metrics are also utilized for evaluation: (i) **FID** [26] is an

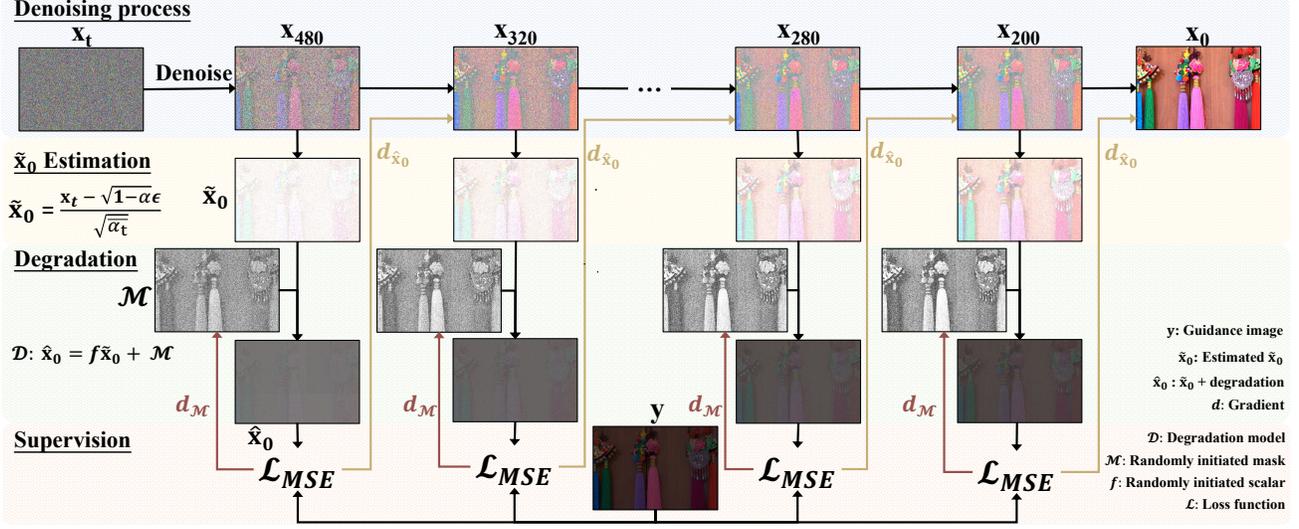


Figure 12. Illustration of the patch-based method for any-size image restoration.

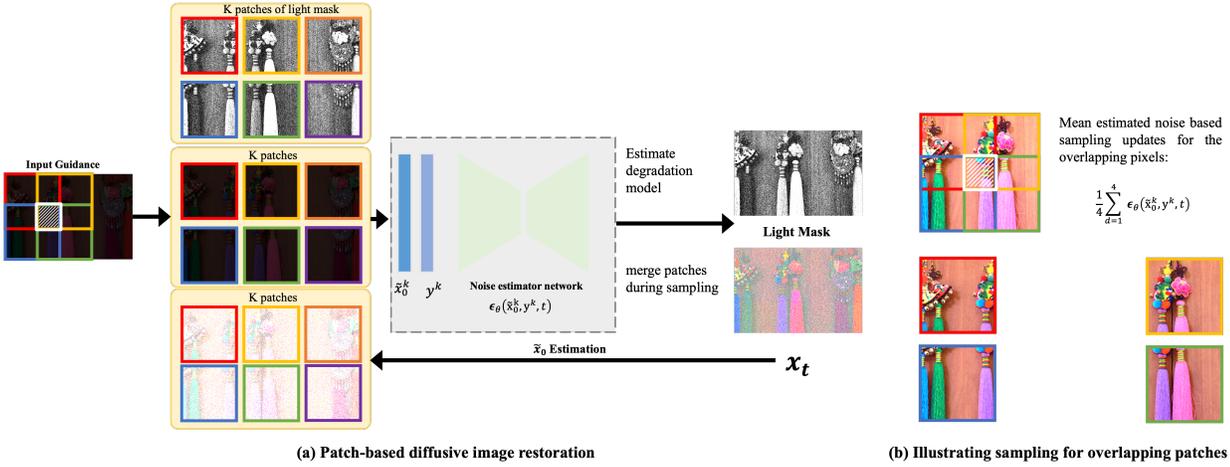


Figure 13. (a) Illustration of the patch-based image restoration pipeline detailed in Algorithm 6. (b) Illustrating mean estimated noise-guided sampling updates for overlapping pixels across patches. We demonstrate a simplified example where  $r = p/2$ ,  $r$  is the stride and  $p$  is the patch size of images. And there are only four overlapping patches sharing the grid cell marked with the white border and gratings. The pixels in this region would be updated at each denoising step  $t$  using the mean estimated noise over the four overlapping patches.

objective metric used to assess the quality of synthesized images. (ii) **Consistency** [73] measures MSE between the degraded inputs and the outputs undergoing the same degradation. (iii) **Learned perceptual image patch similarity (LPIPS)** [101] is also adopted, a deep feature-based perceptual distance metric to further assess the image quality. (iv) The non-reference **perceptual index (PI)** [57] is also employed to evaluate perceptual quality. The PI metric is originally utilized to measure perceptual quality in image super-resolution. It has also been used to assess the performance of other image restoration tasks. A lower PI value indicates better perceptual quality. (v) The **lightness order**

**error (LOE)** [85] is employed as our objective metric to measure the performance. The definition of LOE is as follows:

$$LOE = \frac{1}{m} \sum_{x=1}^m \sum_{y=1}^m (U(\mathbf{T}(x), \mathbf{T}(y)) \oplus U(\mathbf{T}_r(x), \mathbf{T}_r(y))) \quad (14)$$

where  $m$  is the pixel number. The function  $U(p, q)$  returns 1 if  $p \geq q$ , 0 otherwise.  $\oplus$  stands for the exclusive-or operator. In addition,  $\mathbf{T}(x)$  and  $\mathbf{T}_r(x)$  are the maximum values among R, G and B channels at location  $x$  of the enhanced and reference images, respectively. The lower the

LOE is, the better the enhancement preserves the naturalness of lightness.

## D. Further elaboration of the models

**GDP- $x_t$ .** As shown in Fig. 9, the guidance is conditioned on  $x_t$  but with the absence of  $\Sigma$ . The noisy images are gradually denoised during the reverse process. And the  $x_t$  undergoing the degradation model is more similar to the corrupted image. The gradients  $\nabla$  of the loss function are utilized to control the mean of the conditional distribution.

**GDP- $x_0$ .** To make a clear comparison, we also illustrate the GDP- $x_0$  in Fig. 10, and Algorithm 2 in the main paper. Different from the GDP- $x_t$ , GDP- $x_0$  will predict the intermediate variance  $\tilde{x}_0$  from the noisy image  $x_t$  by estimating the noise in  $x_t$ , which can be directly inferred when given the  $x_t$  in every time steps  $t$ . Then the predicted  $\tilde{x}_0$  goes through the same degradation as input to obtain  $\hat{x}_0$ . Note that the degradation might be unknown. Then the loss between the  $\hat{x}_0$  and the corrupted image  $y$ , the gradients will be applied to optimize the unknown degradation models and sample the next step latent  $x_{t-1}$ .

**HDR-GDP- $x_0$ .** As depicted in Fig. 11, and Algorithm 5, there are three images to guide the reverse process. As a blind problem, we randomly initiate three sets of the parameters of the degradation models. At every time step,  $\tilde{x}_0$  will undergo the three degradation models  $\mathcal{D}^i$ , respectively. Unlike GDP- $x_0$ , the gradients of the three losses are used to optimize the corresponding degradation model and all leveraged to sample the next step latent  $x_{t-1}$ .

**Hierarchical Guidance and Patch-based Methods.** As vividly illustrated in Fig. 12 and 13, we resize the corrupted images  $y \in \mathbb{R}^{3 \times H \times W}$  to  $\bar{y} \in \mathbb{R}^{3 \times 256 \times \bar{W}}$  or  $3 \times \bar{H} \times 256$ , then apply the patch-based methods [61] on the reshaped images. Following that, the light masks  $\bar{\mathcal{M}}$  are interpolated to the original image size to obtain the  $\mathcal{M}$ , which can be regarded as the global light shift. After, the light factor  $f$  and the light mask  $\mathcal{M}$  will be fixed and utilized to generate the image patches of the original images, which will be finally recombined as the output images. In our experiments, low-light enhancement and HDR recovery problems can be tackled by this strategy.

## E. Further Ablation Study on the Guidance

To gain insight into the way of guidance, apart from GDP- $x_t$  and GDP- $x_0$ , two more variants GDP- $x_t$ -v1 and GDP- $x_0$ -v1 are devised for comparison.

The main difference among these four variants is the way of mean shift. The mean shift of four variants can be written

---

**Algorithm 7: GDP- $x_t$ -v1** with fixed degradation model: Conditioner guided diffusion sampling on  $x_t$ , given a diffusion model  $(\mu_\theta(x_t), \Sigma_\theta(x_t))$ , corrupted image conditioner  $y$ .

---

**Input:** Corrupted image  $y$ , gradient scale  $s$ , degradation model  $\mathcal{D}$ , distance measure  $\mathcal{L}$ , optional quality enhancement loss  $\mathcal{Q}$ , quality enhancement scale  $\lambda$ .

**Output:** Output image  $x_0$  conditioned on  $y$   
Sample  $x_T$  from  $\mathcal{N}(0, \mathbf{I})$

**for**  $t$  from  $T$  to  $1$  **do**

$\mu, \Sigma = \mu_\theta(x_t), \Sigma_\theta(x_t)$

$\tilde{x}_0 = \frac{x_t}{\sqrt{\alpha_t}} - \frac{\sqrt{1-\alpha_t}\epsilon_\theta(x_t, t)}{\sqrt{\alpha_t}}$

$\mathcal{L}_{x_t}^{total} = \mathcal{L}(y, \mathcal{D}(x_t)) + \mathcal{Q}(x_t)$

$x_t \leftarrow x_t - s \nabla_{x_t} \mathcal{L}(y, \mathcal{D}(x_t))$

Sample  $x_{t-1}$  by  $q(x_{t-1} | x_t, \tilde{x}_0) =$

$\mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, \tilde{x}_0), \tilde{\beta}_t \mathbf{I})$ ,

where

$\tilde{\mu}_t(x_t, \tilde{x}_0) = \frac{\sqrt{\alpha_t-1}\beta_t}{1-\alpha_t} \tilde{x}_0 + \frac{\sqrt{\alpha_t}(1-\alpha_t-1)}{1-\alpha_t} x_t$

and  $\tilde{\beta}_t = \frac{1-\alpha_t-1}{1-\alpha_t} \beta_t$

**end**

**return**  $x_0$

---

as follow:

GDP- $x_0$  :  $\tilde{\mu}_t(x_t, \tilde{x}_0) =$

$$\frac{\sqrt{\alpha_t-1}\beta_t}{1-\alpha_t} \tilde{x}_0 + \frac{\sqrt{\alpha_t}(1-\alpha_t-1)}{1-\alpha_t} x_t + s \nabla_{\tilde{x}_0} \mathcal{L}_{x_0}^{total}$$

GDP- $x_t$  :  $\tilde{\mu}_t(x_t, \tilde{x}_0) =$

$$\frac{\sqrt{\alpha_t-1}\beta_t}{1-\alpha_t} \tilde{x}_0 + \frac{\sqrt{\alpha_t}(1-\alpha_t-1)}{1-\alpha_t} x_t + s \nabla_{x_t} \mathcal{L}_{x_t}^{total}$$

GDP- $x_0$ -v1 :  $\tilde{\mu}_t(x_t, \tilde{x}_0) =$

$$\frac{\sqrt{\alpha_t-1}\beta_t}{1-\alpha_t} (\tilde{x}_0 + s \nabla_{\tilde{x}_0} \mathcal{L}_{x_0}^{total}) + \frac{\sqrt{\alpha_t}(1-\alpha_t-1)}{1-\alpha_t} x_t$$

GDP- $x_t$ -v1 :  $\tilde{\mu}_t(x_t, \tilde{x}_0) =$

$$\frac{\sqrt{\alpha_t-1}\beta_t}{1-\alpha_t} \tilde{x}_0 + \frac{\sqrt{\alpha_t}(1-\alpha_t-1)}{1-\alpha_t} (x_t + s \nabla_{x_t} \mathcal{L}_{x_t}^{total}).$$

(15)

where GDP- $x_0$  directly add the mean shift  $s \nabla_{x_0} \mathcal{L}_{x_0}^{total}$  into  $\tilde{\mu}_t(x_t, \tilde{x}_0)$  without the coefficient  $\frac{\sqrt{\alpha_t-1}\beta_t}{1-\alpha_t}$ , compared with GDP- $x_0$ -v1.

It is experimentally found that GDP- $x_0$  and GDP- $x_t$  fulfills better performance on four linear tasks than GDP- $x_0$ -v1 and GDP- $x_t$ -v1 in Table. 8.

Table 7. **The guidance scales and the number of optimization per time step on the various tasks.** Note that these parameters may not be optimal due to the infinite number of possible combinations.

Tasks	Dataset	Guidance scale	The number of optimization per time step
4× Super-resolution	ImageNet [62]	2E+03	6
Deblurring	ImageNet [62]	6E+03	6
25% Inpainting	ImageNet [62]	4E+03	6
Colorization	ImageNet [62]	6E+03	6
Low-light enhancement	LOL dataset [88]	1E+05	6
HDR recovery	NTIRE dataset [64]	1E+05	1

Table 8. The performance of ablation studies on the way of guidance. We compare four ways of guidance in terms of FID.

FID	4x super-resolution	Deblur	25% Inpainting	Colorization
GDP- $\mathbf{x}_t$ -v1	108.06	88.52	113.47	102.37
GDP- $\mathbf{x}_0$ -v1	44.16	10.35	37.32	41.53
GDP- $\mathbf{x}_t$	64.67	5.00	20.24	66.43
GDP- $\mathbf{x}_0$	38.24	2.44	16.58	37.60

**Algorithm 8: GDP- $\mathbf{x}_0$ -v1:** Conditioner guided diffusion sampling on  $\tilde{\mathbf{x}}_0$ , given a diffusion model  $(\mu_\theta(\mathbf{x}_t), \Sigma_\theta(\mathbf{x}_t))$ , corrupted image conditioner  $\mathbf{y}$ .

**Input:** Corrupted image  $\mathbf{y}$ , gradient scale  $s$ , degradation model  $\mathcal{D}$ , distance measure  $\mathcal{L}$ .

**Output:** Output image  $\mathbf{x}_0$  conditioned on  $\mathbf{y}$

Sample  $\mathbf{x}_T$  from  $\mathcal{N}(0, \mathbf{I})$

**for**  $t$  **from**  $T$  **to**  $1$  **do**

$\mu, \Sigma = \mu_\theta(\mathbf{x}_t), \Sigma_\theta(\mathbf{x}_t)$   
 $\tilde{\mathbf{x}}_0 = \frac{\mathbf{x}_t}{\sqrt{\alpha_t}} - \frac{\sqrt{1-\alpha_t}\epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{\alpha_t}}$   
 $\mathcal{L}_{\tilde{\mathbf{x}}_0}^{total} = \mathcal{L}(\mathbf{y}, \mathcal{D}(\tilde{\mathbf{x}}_0)) + \mathcal{Q}(\tilde{\mathbf{x}}_0)$

$\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_0 - s \nabla_{\tilde{\mathbf{x}}_0} \mathcal{L}_{\tilde{\mathbf{x}}_0}^{total}$   
 Sample  $\mathbf{x}_{t-1}$  by  $q(\mathbf{x}_{t-1} | \mathbf{x}_t, \tilde{\mathbf{x}}_0) =$

$\mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t(\mathbf{x}_t, \tilde{\mathbf{x}}_0), \tilde{\beta}_t \mathbf{I})$ ,

where

$\tilde{\mu}_t(\mathbf{x}_t, \tilde{\mathbf{x}}_0) = \frac{\sqrt{\alpha_{t-1}}\beta_t}{1-\alpha_t} \tilde{\mathbf{x}}_0 + \frac{\sqrt{\alpha_t}(1-\alpha_{t-1})}{1-\alpha_t} \mathbf{x}_t$   
 and  $\tilde{\beta}_t = \frac{1-\alpha_{t-1}}{1-\alpha_t} \beta_t$

**end**

**return**  $\mathbf{x}_0$

## F. The ELBO objective of GDP

GDP is a Markov chain conditioned on  $\mathbf{y}$ , resulting in the following ELBO objective [77]:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0), \mathbf{y} \sim q(\mathbf{y} | \mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0 | \mathbf{y})] \geq \\ & - \mathbb{E} \left[ \sum_{t=1}^{T-1} \text{KL} \left( q^t(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{x}_0, \mathbf{y}) \| p_\theta^t(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) \right) \right] \\ & + \mathbb{E} [\log p_\theta^0(\mathbf{x}_0 | \mathbf{x}_1, \mathbf{y})] \\ & - \mathbb{E} \left[ \text{KL} \left( q^T(\mathbf{x}_T | \mathbf{x}_0, \mathbf{y}) \| p_\theta^T(\mathbf{x}_T | \mathbf{y}) \right) \right] \end{aligned} \quad (16)$$

where  $q(\mathbf{x}_0)$  denotes the data distribution,  $q(\mathbf{y} | \mathbf{x}_0)$  in the main paper, the expectation on the right-hand side is given by sampling  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ ,  $\mathbf{y} \sim q(\mathbf{y} | \mathbf{x}_0)$ ,  $\mathbf{x}_T \sim q^T(\mathbf{x}_T | \mathbf{x}_0, \mathbf{y})$ , and  $\mathbf{x}_t \sim q^t(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{x}_0, \mathbf{y})$  for  $t \in [1, T-1]$ .

## G. Sampling with DDIM

To accelerate the sampling strategy, GDP follows [58] to use DDIM, which skipping steps in the reverse process to speed up the DDPM generating process. We apply this method to the ImageNet dataset on the four tasks. We set the  $T=20$  in the sampling process, while DDRM also utilizes the same time steps for a fair comparison. As shown in Table 9, our GDP- $\mathbf{x}_0$ -DDIM(20) outperforms DDRM(20) on consistency and FID across four tasks. Although DDRM(20) obtains better PSNR and SSIM than our GDP- $\mathbf{x}_0$ -DDIM(20), the qualitative results of DDRM(20) are still worse than our GDP- $\mathbf{x}_0$ -DDIM(20), which can be seen from Figs. 14 and 15. Previous work [11, 13, 17, 73]

demonstrated that these conventional automated evaluation measures (PSNR and SSIM) do not correlate well with human perception when the input resolution is low, and the magnification is large. This is not surprising since these metrics tend to penalize any synthesized high-frequency detail that is not perfectly aligned with the target image.

## H. Image Guidance

A conditioner  $p(\mathbf{y} | \mathbf{x})$  is exploited to improve a diffusion generator. Specifically, we can utilize a conditioner  $p_\phi(\mathbf{y} | \mathbf{x}_t, t)$  on input images, and then use gradients  $\nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t, t)$  to guide the diffusion sampling process towards a given the degraded images  $\mathbf{y}$ .

In this section, we will describe how to use such conditioners to improve the quality of sampled images. The notation is chosen as  $p_\phi(\mathbf{y} | \mathbf{x}_t, t) = p_\phi(\mathbf{y} | \mathbf{x}_t)$  and  $\epsilon_\theta(\mathbf{x}_t, t) = \epsilon_\theta(\mathbf{x}_t)$  for brevity. Note that they refer to separate functions for each time step  $t$ .

### H.1. Conditional Reverse Process

Assume a diffusion model with an unconditional reverse noising process  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$ . In image restoration and enhancement, the corrupted inputs can be regarded as conditions. Therefore, we regard  $\mathbf{y}$  as the input images, and  $\mathbf{x}_t$  as the generated images in time step  $t$ . Then, the conditioner is formulated as follows:

$$p_\phi(\mathbf{y} | \mathbf{x}_t) = \frac{1}{K} \exp(-\mathcal{L}(\mathbf{y}, \mathcal{D}(\mathbf{x}_t))), \quad (17)$$

where  $\mathcal{D}$  represents the degradation function,  $\mathcal{L}$  stands for Mean Square Error together with optional Quality Enhancement Loss, and  $K$  is an arbitrary constant. In order to condition this on the input corrupted image  $\mathbf{y}$ , it is sufficient to sample each transition based on the following:

$$p_{\theta, \phi}(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) = C p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) p_\phi(\mathbf{y} | \mathbf{x}_t) \quad (18)$$

where  $C$  denotes a normalizing constant. It is typically intractable to sample from this distribution exactly, but Sohl-Dickstein *et al.* [76] show that it can be approximated as a perturbed Gaussian distribution. Sampling accurately from this distribution is often tricky, but Sohl-Dickstein *et al.* [76] prove that it could be approximated as a perturbed Gaussian distribution. It is formulated that the diffusion model samples the previous time step  $\mathbf{x}_t$  from time step  $\mathbf{x}_{t+1}$  via a Gaussian distribution:

$$p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) = \mathcal{N}(\mu, \Sigma) \quad (19)$$

$$\log p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) = -\frac{1}{2} (\mathbf{x}_t - \mu)_T \Sigma^{-1} (\mathbf{x}_t - \mu) + Z \quad (20)$$

We can assume that  $\log_\phi p(\mathbf{y} | \mathbf{x}_t)$  owns low curvature when compared with  $\Sigma^{-1}$ . This assumption is reasonable

under the constraint that the infinite diffusion step, where  $\|\Sigma\| \rightarrow 0$ . Under the circumstances,  $\log p_\phi(\mathbf{y} | \mathbf{x}_t)$  can be approximated via a Taylor expansion around  $\mathbf{x}_t = \mu$  as:

$$\begin{aligned} \log p_\phi(\mathbf{y} | \mathbf{x}_t) &\approx \log p_\phi(\mathbf{y} | \mathbf{x}_t)|_{\mathbf{x}_t=\mu} \\ &\quad + (\mathbf{x}_t - \mu) \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t)|_{\mathbf{x}_t=\mu} \\ &= (\mathbf{x}_t - \mu) g + Z_1 \end{aligned} \quad (21)$$

Here,  $g = \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y} | \mathbf{x}_t)|_{\mathbf{x}_t=\mu}$ , and  $Z_1$  is a constant. We can replace the  $g$  with Eq. 17 as follows:

$$\log p(\mathbf{y} | \mathbf{x}_t) = -\mathcal{L}(\mathbf{y}, \mathcal{D}(\mathbf{x}_t)) - \log K \quad (22)$$

$$g = \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) = -\nabla_{\mathbf{x}_t} \mathcal{L}(\mathbf{y}, \mathcal{D}(\mathbf{x}_t)) \quad (23)$$

This gives:

$$\begin{aligned} &\log(p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) p_\phi(\mathbf{y} | \mathbf{x}_t)) \\ &\approx -\frac{1}{2} (\mathbf{x}_t - \mu)^T \Sigma^{-1} (\mathbf{x}_t - \mu) + (\mathbf{x}_t - \mu) g + Z_2 \\ &= -\frac{1}{2} (\mathbf{x}_t - \mu - \Sigma g)^T \Sigma^{-1} (\mathbf{x}_t - \mu - \Sigma g) + \frac{1}{2} g^T \Sigma g + Z_2 \\ &= -\frac{1}{2} (\mathbf{x}_t - \mu - \Sigma g)^T \Sigma^{-1} (\mathbf{x}_t - \mu - \Sigma g) + Z_3 \\ &= \log p(z) + Z_4, z \sim \mathcal{N}(\mu + \Sigma g, \Sigma) \end{aligned} \quad (24)$$

where the constant term  $C_4$  could be safely ignored because it is equivalent to the normalizing coefficient  $Z$  in Eq. 18. Thus, we find that the conditional transition operator can be approximated by a Gaussian similar to the unconditional transition operator, but with a mean shifted by  $\Sigma g$ . Moreover, an optional scaling factor  $s$  is included for gradients, which will be described in more detail in Sec. H.3. However, it is experimentally found that this guidance way might not be effective enough, where our GDP- $x_0$  is systematically studied.

### H.2. Conditional Diffusion Process

Here, we figure out that conditional sampling can be fulfilled with a transition operator proportional to  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1}) p_\phi(\mathbf{y} | \mathbf{x}_t)$ , where  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$  approximates  $q(\mathbf{x}_t | \mathbf{x}_{t+1})$  and  $p_\phi(\mathbf{y} | \mathbf{x}_t)$  approximates the distribution of the input for a noised sample  $\mathbf{x}_t$ .

A conditional Markovian noising process  $\hat{q}$  is similar to  $q$ . And  $\hat{q}(\mathbf{y} | \mathbf{x}_0)$  is assumed as a known and readily available degraded images distribution for each sample.

$$\hat{q}(\mathbf{x}_0) := q(\mathbf{x}_0) \quad (25)$$

$$\hat{q}(\mathbf{y} | \mathbf{x}_0) := \text{Corrupted input image per sample} \quad (26)$$

$$\hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{y}) := q(\mathbf{x}_{t+1} | \mathbf{x}_t) \quad (27)$$

$$\hat{q}(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{y}) := \prod_{t=1}^T \hat{q}(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}) \quad (28)$$

Table 9. **The performances of DDRM (20) and GDM- $x_0$ -DDIM(20) towards the four tasks on ImageNet 1k.** The DDIM sample steps are all set to 20 to make a fair comparison.

Task	4× super resolution				Deblur				25% Impainting				Colorization			
	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID
DDRM(20) [32]	<b>26.53</b>	<b>0.784</b>	19.39	40.75	<b>35.64</b>	<b>0.978</b>	50.24	4.78	<b>34.28</b>	<b>0.958</b>	<b>4.08</b>	24.09	<b>22.12</b>	<b>0.924</b>	38.66	47.05
GDP- $x_0$ -DDIM(20)	23.77	0.623	<b>9.24</b>	<b>39.46</b>	24.87	0.683	<b>44.39</b>	<b>3.66</b>	30.82	0.892	7.10	<b>19.70</b>	21.13	0.840	<b>37.33</b>	<b>41.38</b>

Table 10. The time comparison of GDP- $x_0$ -DDIM(20) and GDP- $x_0$  on 4x super-resolution. These experiments are compared on Tesla A100.

	Guidance scale	Total steps	Guidance times per steps	Generation time per image
GDP- $x_0$	w.o. DDIM	2e3	1000	69.55
GDP- $x_0$ -DDIM(20)	w. DDIM	22e5	20	1.74

Assuming that the noising process  $\hat{q}$  is conditioned on  $\mathbf{y}$ , we can reveal that  $\hat{q}$  behaves exactly like  $q$  when not conditioned on  $\mathbf{y}$ . According to this idea, we first derive the unconditional noising operator  $\hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)$ :

$$\hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) = \int_{\mathbf{y}} \hat{q}(\mathbf{x}_{t+1}, \mathbf{y} | \mathbf{x}_t) d\mathbf{y} \quad (29)$$

$$= \int_{\mathbf{y}} \hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{y}) \hat{q}(\mathbf{y} | \mathbf{x}_t) d\mathbf{y} \quad (30)$$

$$= \int_{\mathbf{y}} q(\mathbf{x}_{t+1} | \mathbf{x}_t) \hat{q}(\mathbf{y} | \mathbf{x}_t) d\mathbf{y} \quad (31)$$

$$= q(\mathbf{x}_{t+1} | \mathbf{x}_t) \int_{\mathbf{y}} \hat{q}(\mathbf{y} | \mathbf{x}_t) d\mathbf{y} \quad (32)$$

$$= q(\mathbf{x}_{t+1} | \mathbf{x}_t) \quad (33)$$

$$= \hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{y}) \quad (34)$$

Similarly, the joint distribution  $\hat{q}(\mathbf{x}_{1:T} | \mathbf{x}_0)$  can be written as:

$$\hat{q}(\mathbf{x}_{1:T} | \mathbf{x}_0) = \int_{\mathbf{y}} \hat{q}(\mathbf{x}_{1:T}, \mathbf{y} | \mathbf{x}_0) d\mathbf{y} \quad (35)$$

$$= \int_{\mathbf{y}} \hat{q}(\mathbf{y} | \mathbf{x}_0) \hat{q}(\mathbf{x}_{1:T} | \mathbf{x}_0, \mathbf{y}) d\mathbf{y} \quad (36)$$

$$= \int_{\mathbf{y}} \hat{q}(\mathbf{y} | \mathbf{x}_0) \prod_{t=1}^T \hat{q}(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{y}) d\mathbf{y} \quad (37)$$

$$= \int_{\mathbf{y}} \hat{q}(\mathbf{y} | \mathbf{x}_0) \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) d\mathbf{y} \quad (38)$$

$$= \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \int_{\mathbf{y}} \hat{q}(\mathbf{y} | \mathbf{x}_0) d\mathbf{y} \quad (39)$$

$$= \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (40)$$

$$= q(\mathbf{x}_{1:T} | \mathbf{x}_0) \quad (41)$$

$\hat{q}(\mathbf{x}_t)$  can be derived by using Eq. 41 as follows:

$$\hat{q}(\mathbf{x}_t) = \int_{\mathbf{x}_{0:t-1}} \hat{q}(\mathbf{x}_0, \dots, \mathbf{x}_t) d\mathbf{x}_{0:t-1} \quad (42)$$

$$= \int_{\mathbf{x}_{0:t-1}} \hat{q}(\mathbf{x}_0) \hat{q}(\mathbf{x}_1, \dots, \mathbf{x}_t | \mathbf{x}_0) d\mathbf{x}_{0:t-1} \quad (43)$$

$$= \int_{\mathbf{x}_{0:t-1}} q(\mathbf{x}_0) q(\mathbf{x}_1, \dots, \mathbf{x}_t | \mathbf{x}_0) d\mathbf{x}_{0:t-1} \quad (44)$$

$$= \int_{\mathbf{x}_{0:t-1}} q(\mathbf{x}_0, \dots, \mathbf{x}_t) d\mathbf{x}_{0:t-1} \quad (45)$$

$$= q(\mathbf{x}_t) \quad (46)$$

It is proved by Bayes rule that the unconditional reverse process  $\hat{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) = q(\mathbf{x}_t | \mathbf{x}_{t+1})$  when using the identities  $\hat{q}(\mathbf{x}_t) = q(\mathbf{x}_t)$  and  $\hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) = q(\mathbf{x}_{t+1} | \mathbf{x}_t)$ .

Note that  $\hat{q}$  is able to produce an input function  $\hat{q}(\mathbf{y} | \mathbf{x}_t)$ . It is shown that this distribution of the input does not depend on  $\mathbf{x}_{t+1}$  (the noisy version of  $\mathbf{x}_t$ ), we will discuss this fact later by exploiting:

$$\hat{q}(\mathbf{y} | \mathbf{x}_t, \mathbf{x}_{t+1}) = \hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{y}) \frac{\hat{q}(\mathbf{y} | \mathbf{x}_t)}{\hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} \quad (47)$$

$$= \hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t) \frac{\hat{q}(\mathbf{y} | \mathbf{x}_t)}{\hat{q}(\mathbf{x}_{t+1} | \mathbf{x}_t)} \quad (48)$$

$$= \hat{q}(\mathbf{y} | \mathbf{x}_t) \quad (49)$$

In this way, the conditional reverse process can be derived

Table 11. The quantitative comparison of performance on CelebA.

CelebA	4x SR				Deblur				25% Inpainting			
	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID	PSNR	SSIM	Consistency	FID
DDRM	29.50	0.863	6.82	87.71	<b>36.51</b>	<b>0.98</b>	35.91	14.30	31.99	0.918	<b>0.47</b>	69.46
GDP- $x_t$	29.19	0.847	14.11	94.98	27.35	0.81	34.87	9.97	36.19	0.963	1.94	22.53
GDP- $x_0$	<b>30.26</b>	<b>0.868</b>	<b>5.33</b>	<b>46.64</b>	28.66	0.83	<b>32.66</b>	<b>4.50</b>	<b>37.70</b>	<b>0.972</b>	0.51	<b>11.62</b>

Table 12. The quantitative comparison of results on LSUN bedroom.

LSUN Bedroom	4x SR		Deblur		25% Inpainting		Colorization	
	Consistency	FID	Consistency	FID	Consistency	FID	Consistency	FID
DDRM	20.33	40.12	43.78	10.16	<b>5.33</b>	22.49	35.16	45.22
GDP- $x_t$	70.46	58.62	46.90	12.50	9.33	20.63	66.88	57.13
GDP- $x_0$	<b>7.66</b>	<b>36.94</b>	<b>42.28</b>	<b>9.51</b>	6.77	<b>18.34</b>	<b>33.51</b>	<b>34.59</b>

as:

$$\hat{q}(\mathbf{x}_t | \mathbf{x}_{t+1}, \mathbf{y}) = \frac{\hat{q}(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{y})}{\hat{q}(\mathbf{x}_{t+1}, \mathbf{y})} \quad (50)$$

$$= \frac{\hat{q}(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{y})}{\hat{q}(\mathbf{y} | \mathbf{x}_{t+1}) \hat{q}(\mathbf{x}_{t+1})} \quad (51)$$

$$= \frac{\hat{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) \hat{q}(\mathbf{y} | \mathbf{x}_t, \mathbf{x}_{t+1}) \hat{q}(\mathbf{x}_{t+1})}{\hat{q}(\mathbf{y} | \mathbf{x}_{t+1}) \hat{q}(\mathbf{x}_{t+1})} \quad (52)$$

$$= \frac{\hat{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) \hat{q}(\mathbf{y} | \mathbf{x}_t, \mathbf{x}_{t+1})}{\hat{q}(\mathbf{y} | \mathbf{x}_{t+1})} \quad (53)$$

$$= \frac{\hat{q}(\mathbf{x}_t | \mathbf{x}_{t+1}) \hat{q}(\mathbf{y} | \mathbf{x}_t)}{\hat{q}(\mathbf{y} | \mathbf{x}_{t+1})} \quad (54)$$

$$= \frac{q(\mathbf{x}_t | \mathbf{x}_{t+1}) \hat{q}(\mathbf{y} | \mathbf{x}_t)}{\hat{q}(\mathbf{y} | \mathbf{x}_{t+1})} \quad (55)$$

where the  $\hat{q}(\mathbf{y} | \mathbf{x}_{t+1})$  can be treated as a constant because it does not depend on  $\mathbf{x}_{t+1}$ . Therefore, we want to sample from the distribution  $Cq(\mathbf{x}_t | \mathbf{x}_{t+1}) \hat{q}(\mathbf{y} | \mathbf{x}_t)$  where  $C$  denotes the normalization constant. We already have a neural network approximation of  $q(\mathbf{x}_t | \mathbf{x}_{t+1})$  called  $p_\theta(\mathbf{x}_t | \mathbf{x}_{t+1})$ , so the rest is  $\hat{q}(\mathbf{y} | \mathbf{x}_t)$  that can be obtained by computing a conditioner  $p_\phi(\mathbf{y} | \mathbf{x}_t)$  on noised images  $\mathbf{x}_t$  derived by sampling from  $q(\mathbf{x}_t)$ .

### H.3. Scaling Conditioner Gradients

The conditioner is incorporated into the sampling process of the diffusion model using Eq. 24. To unveil the effect of scaling conditioner gradients, note that  $s \cdot \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x}) = \nabla_{\mathbf{x}} \log \frac{1}{K} p(\mathbf{y} | \mathbf{x})^s$ , where  $K$  is an arbitrary constant. Thus, the conditioning process is still theoretically based on the re-normalized distribution of the input proportional to  $p(\mathbf{y} | \mathbf{x})^s$ . If  $s > 1$ , this distribution becomes sharper than  $p(\mathbf{y} | \mathbf{x})$  because larger values are

exponentially magnified. Therefore, using a larger gradient scale to focus more on the modes of the conditioner may be beneficial in producing higher fidelity (but less diverse) samples. In this paper, due to the observation that  $\Sigma$  might exert a negative influence on the quality of images. Therefore, with the absence of the  $\Sigma$ , the guidance scale can be a variable scale  $\hat{s}$ , where  $s = \Sigma \hat{s}$ . Thanks to this variable scale  $\hat{s}$ , the quality of images can be promoted

## I. Additional Results on Linear inverse problems

We provide additional figures below showing GDP’s versatility across different datasets and linear inverse problems (Figures 16, 17, 18, 19), and 21). We present more uncurated samples from the ImageNet experiments in Figures 20, 22, 23, 24, 25, and 26. Moreover, our GDP is also able to recover the corrupted images that undergo multi-linear degradations, as shown in Fig. 27

## J. Additional Results on Low-light Enhancement

In addition to the linear inverse problems, we further show more samples on the blind and non-linear task of low-light enhancement. As shown in 28, 31, and 33, our GDP performs well under the three datasets, including LOL, VE-LOL-L, and LoLi-phone, indicating the effectiveness of GDP under the different distributions of the images. Moreover, we also compare the GDP with other methods on the three datasets. As seen in 30, and 32, GDP- $x_0$  is able to generate more satisfactory images than other supervised learning, unsupervised learning, self-supervised, and zero-shot learning methods. Note that GDP- $x_t$  tends to yield images lighter than the ones generated by GDP- $x_0$ . Furthermore, GDP can adjust the brightness of generated images by the

Table 13. The weight of reconstruction loss and quality enhancement loss.

	MSE loss	Exposure Control Loss	Color Constancy Loss	Illumination Smoothness Loss
Colorization	1	0	500	0
Low-light Enhancement	1	1/100	1/200	1
HDR recovery	1	1/100	1/200	1

Exposure Control Loss. As shown in 29, users can change the gray level  $E$  in the RGB color space to obtain the target images with specific brightness.

## K. Additional Results on HDR Recovery

As shown in 35, our HDR-GDP- $x_0$  is capable of adjusting the over-exposed and under-exposed areas of the picture in various scenes. It is noted that since the model used by GDP is pre-trained on ImageNet, the tone of the generated picture will be slightly different from ground truth images. Moreover, we also show more samples compared with the state-of-the-art methods, including AHDRNet [91], HDR-GAN [59], DeepHDR [90] and deep-high-dynamic-range [30]. As seen in Fig. 36, our HDR-GDP- $x_0$  can recover more realistic images with more details.

## L. Additional Results on Ablation Study

The visualization comparisons of the ablation study on the trainable degradation and the patch-based tactic are shown in Figs. 37 and 38. It is shown that Model A fails to generate high-quality images due to the interpolation operation, while Model B generates images with more artifacts because of the naive restoration. Model C predicts the outputs in an uncontrollable way thanks to the randomly initiated and fixed parameters.



Figure 14. **More samples from the  $4 \times$  super-resolution task of GDP- $x_0$ -DDIM (20) compare with DDRM (20) on  $256 \times 256$  ImageNet 1K.** The generated images by the DDRM (20) are still blurred, while our proposed GDP- $x_0$  with 20 steps of DDIM sampling can restore more details.

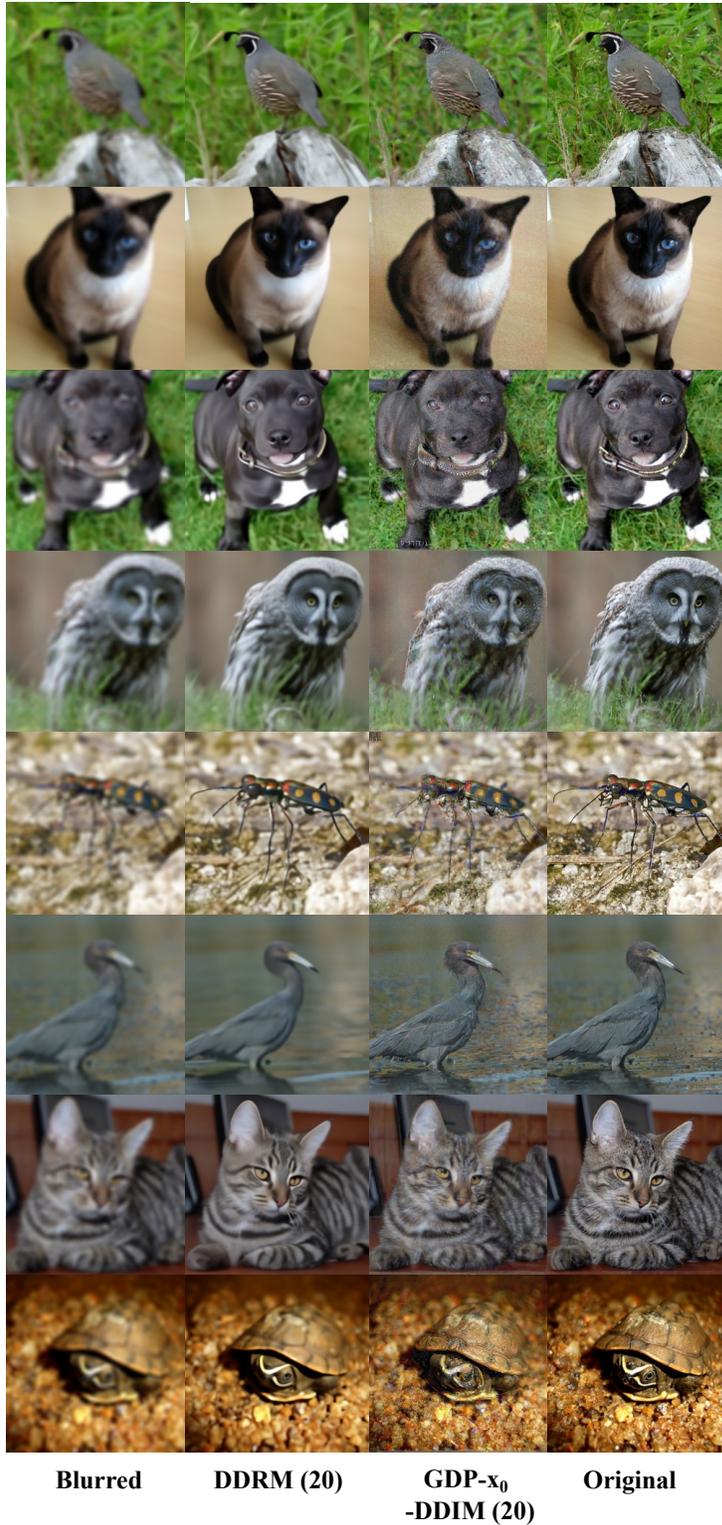


Figure 15. More samples from the deblurring task of GDP- $x_0$ -DDIM (20) compare with DDRM (20) on  $256 \times 256$  ImageNet 1K. Our GDP- $x_0$ -DDIM (20) can recover more details than DDRM (20) under the same DDIM steps.

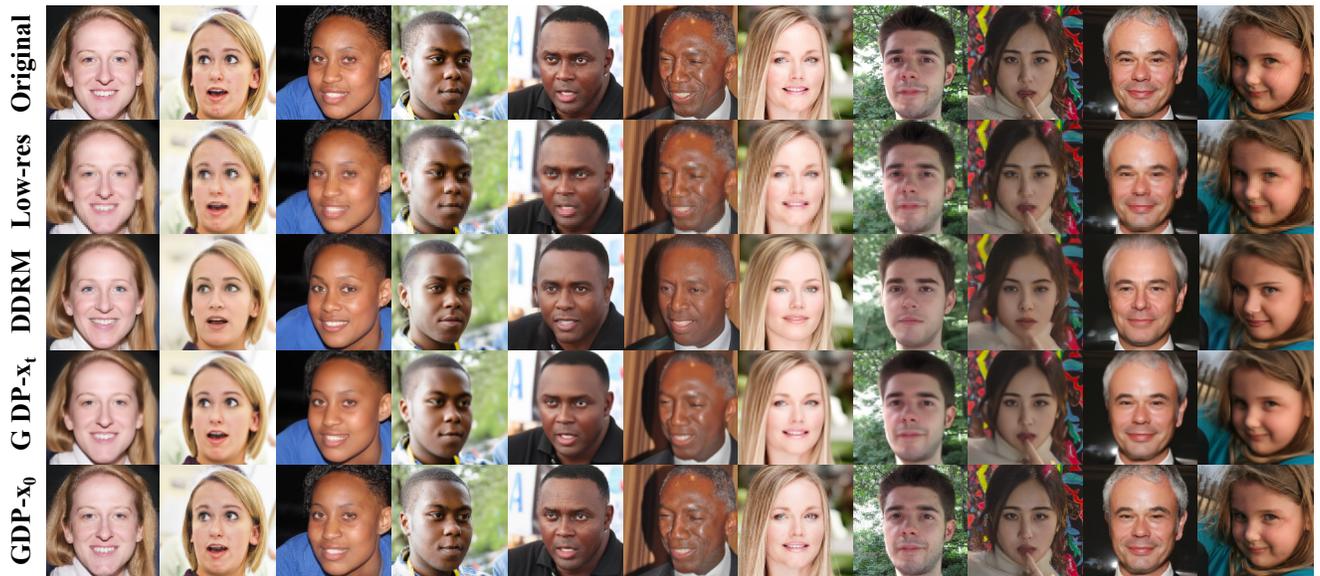


Figure 16.  $4 \times$  super-resolution results of DDRM, GDP- $x_t$ , and GDP- $x_0$  on CelebA face images. Compared with GDP- $x_t$  and DDRM, GDP- $x_0$  can restore more realistic faces, such as the wrinkles on the faces, systematically demonstrating the superiority of the guidance on  $x_0$  protocol.

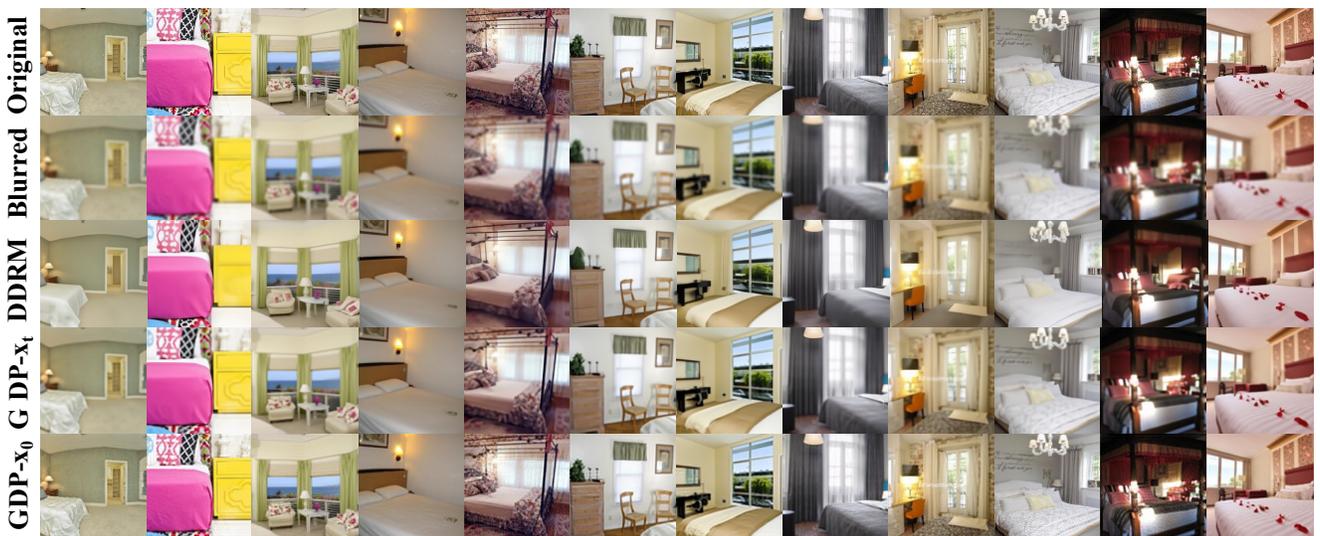


Figure 17. Deblurring results of DDRM, GDP- $x_t$ , and GDP- $x_0$  on LSUN bedroom images.

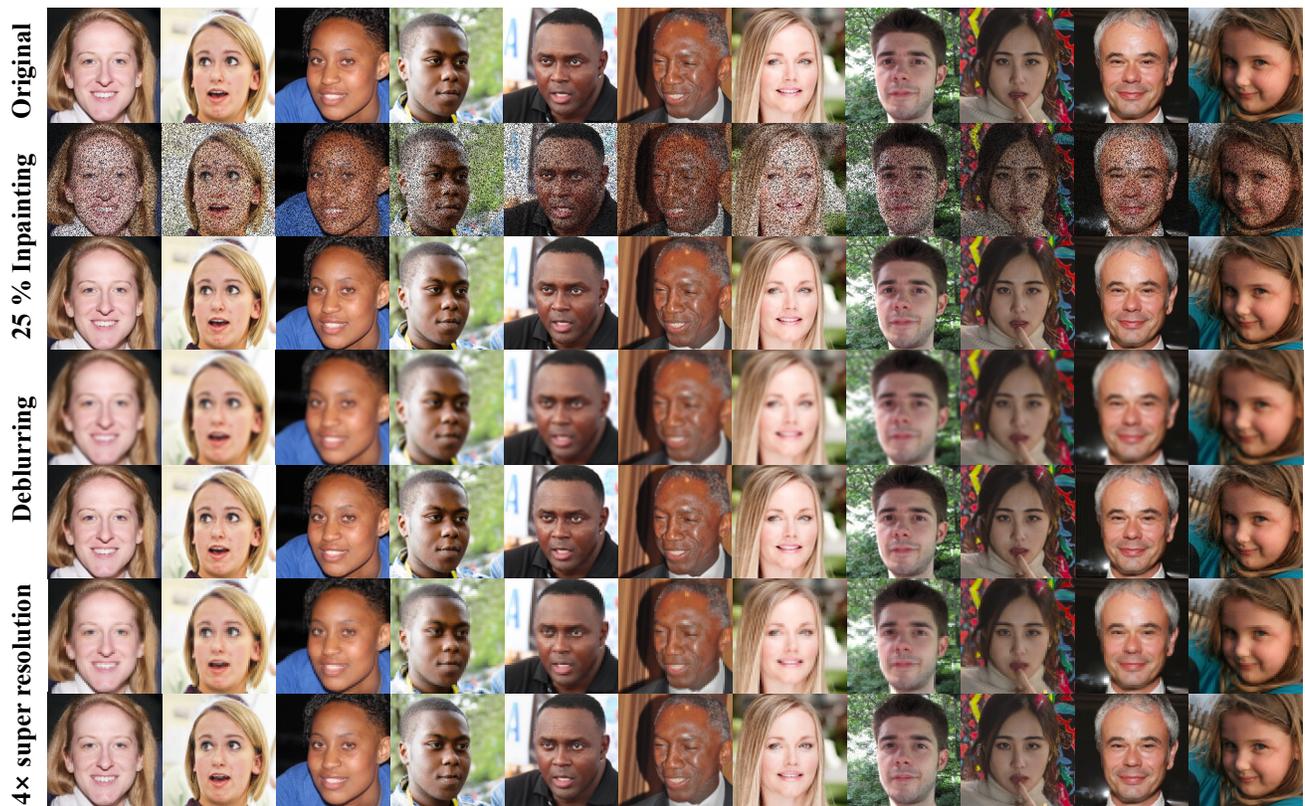


Figure 18. Pairs of degraded and recovered  $256 \times 256$  CelebA face images with a  $GDP-x_0$ . Three tasks including 25% inpainting, deblurring and  $4 \times$  super-resolution are vividly depicted.



Figure 19. Pairs of degraded and recovered  $256 \times 256$  LSUN bedroom images with a  $\text{GDP-}x_0$ . We show more samples under the 25% inpainting, colorization, deblurring and  $4 \times$  super-resolution.

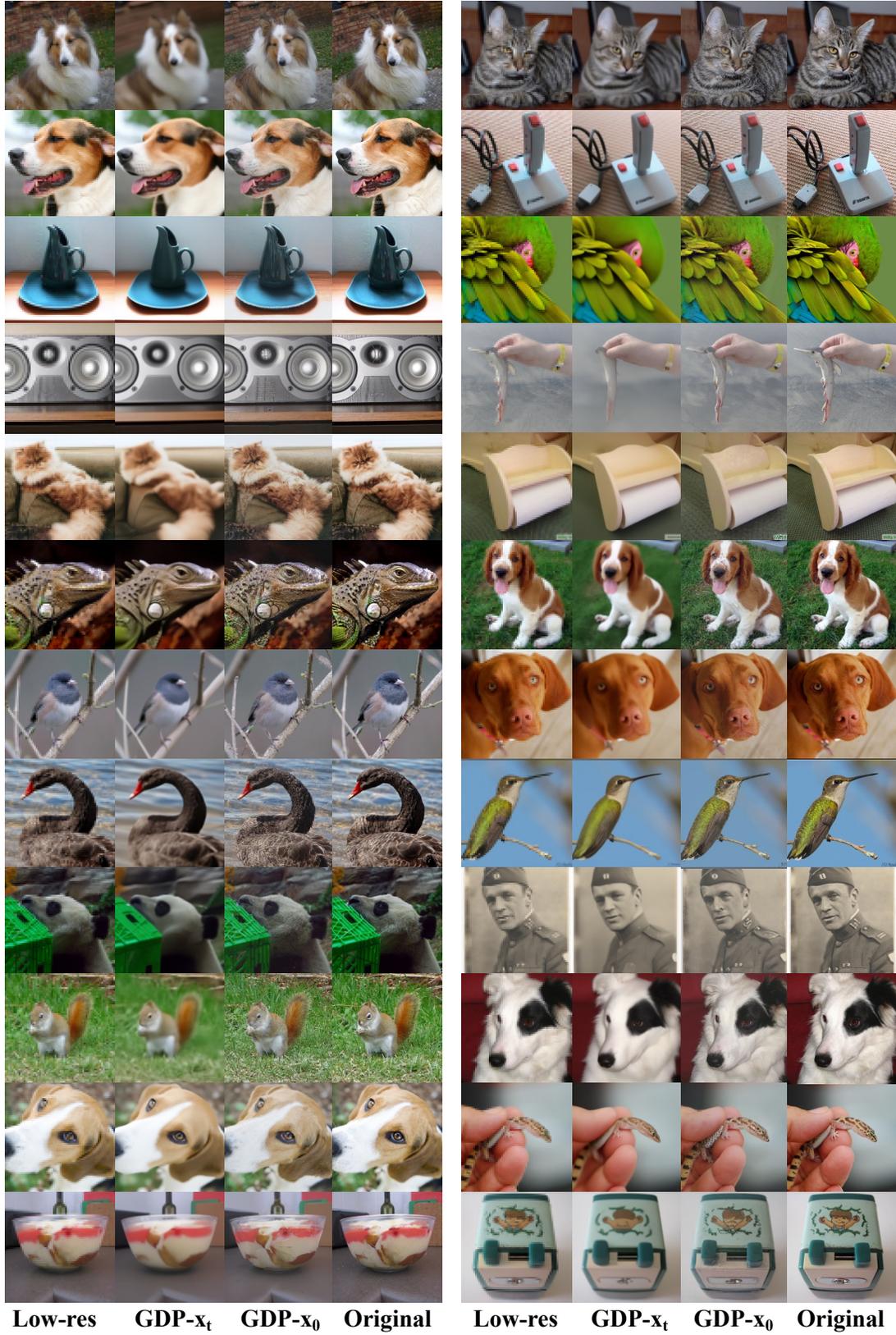
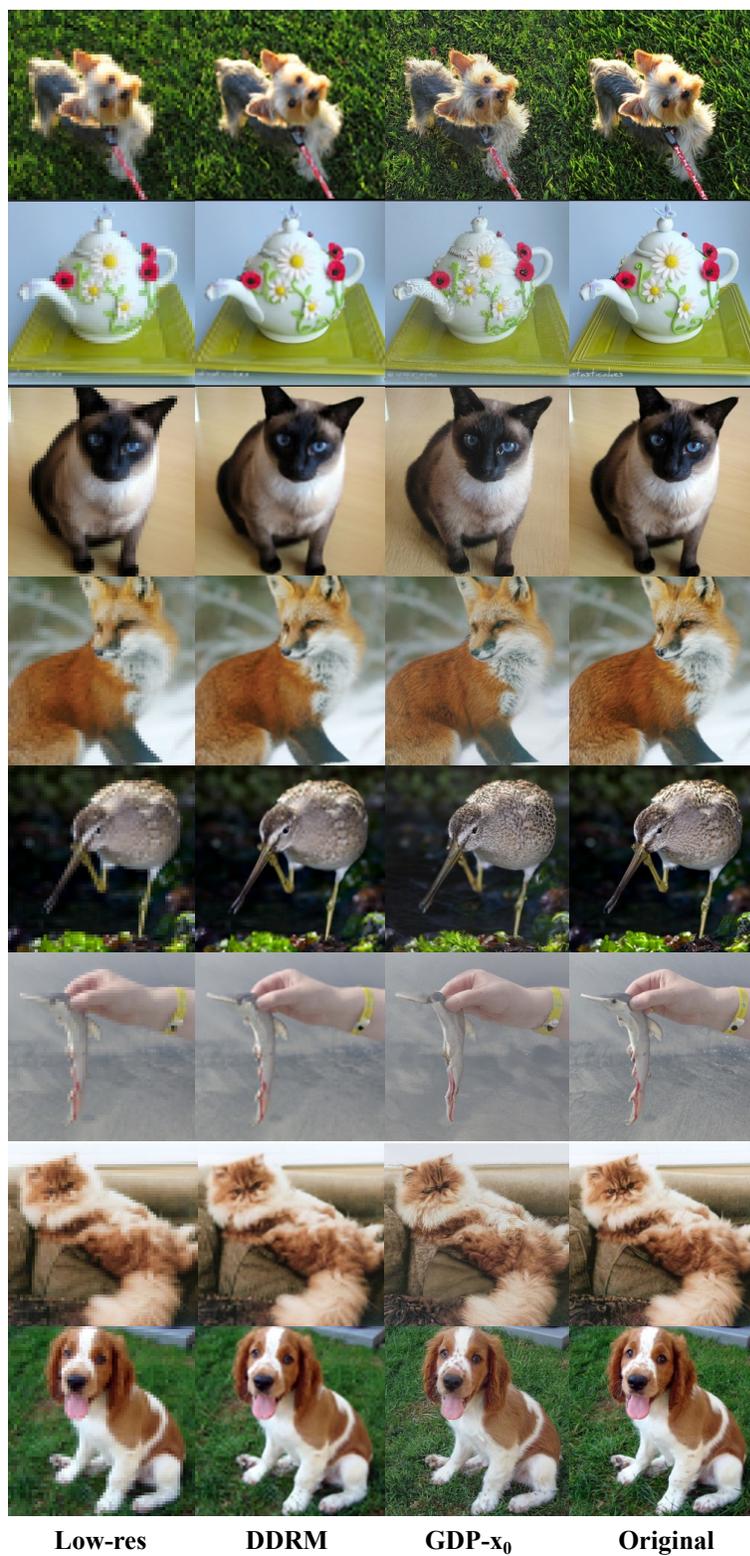


Figure 20. Uncurated samples from the  $4 \times$  super-resolution task on  $256 \times 256$  ImageNet 1K.



Low-res          DDRM          GDP- $x_0$           Original

Figure 21. More samples from the  $4 \times$  super-resolution task compare with DDRM on  $256 \times 256$  ImageNet 1K. As we mentioned above, DDRM adds guidance on the  $x_t$ , leading to the less satisfactory results than our GDP- $x_0$ .

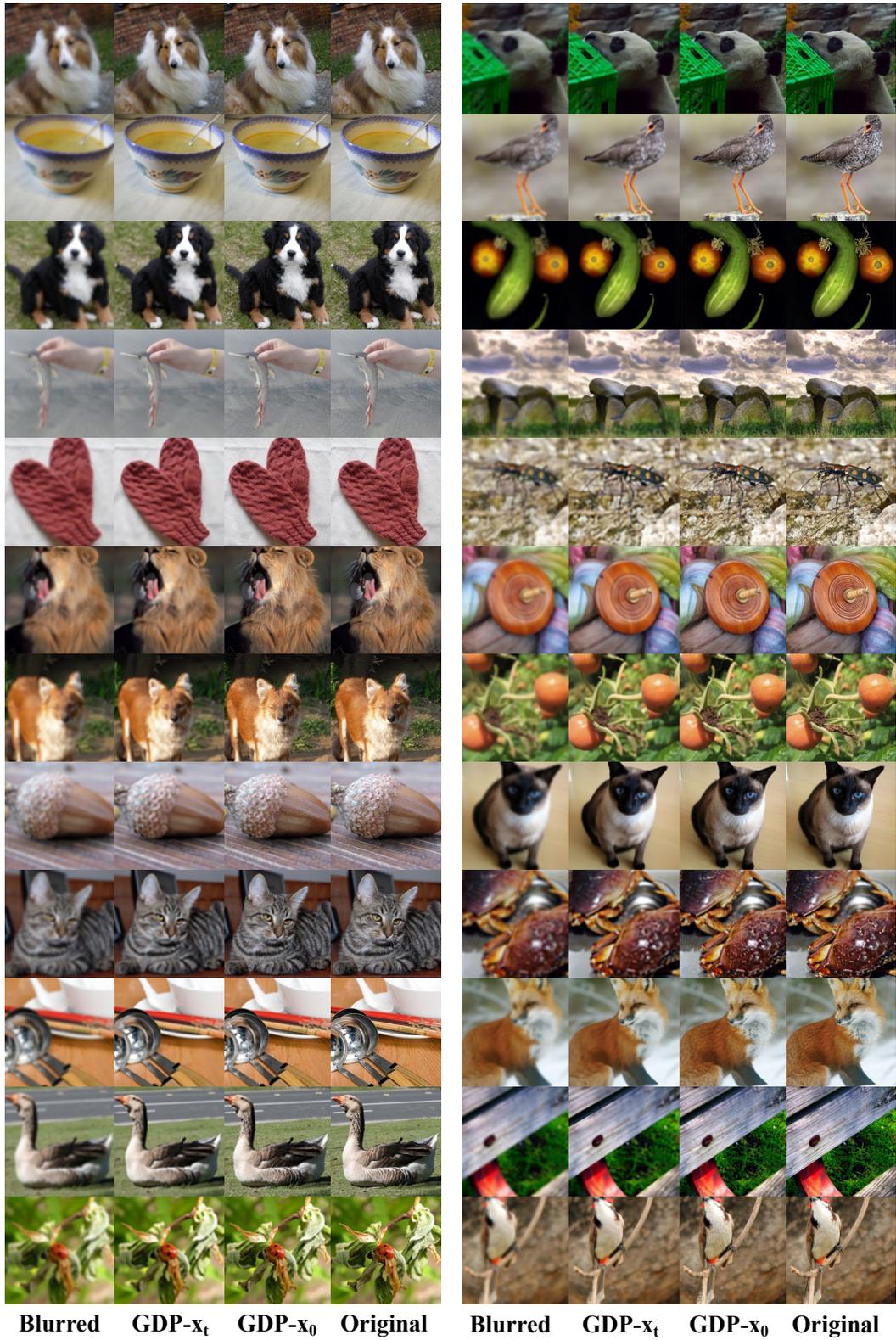


Figure 22. Uncurated samples from the deblurring task on  $256 \times 256$  ImageNet 1K.

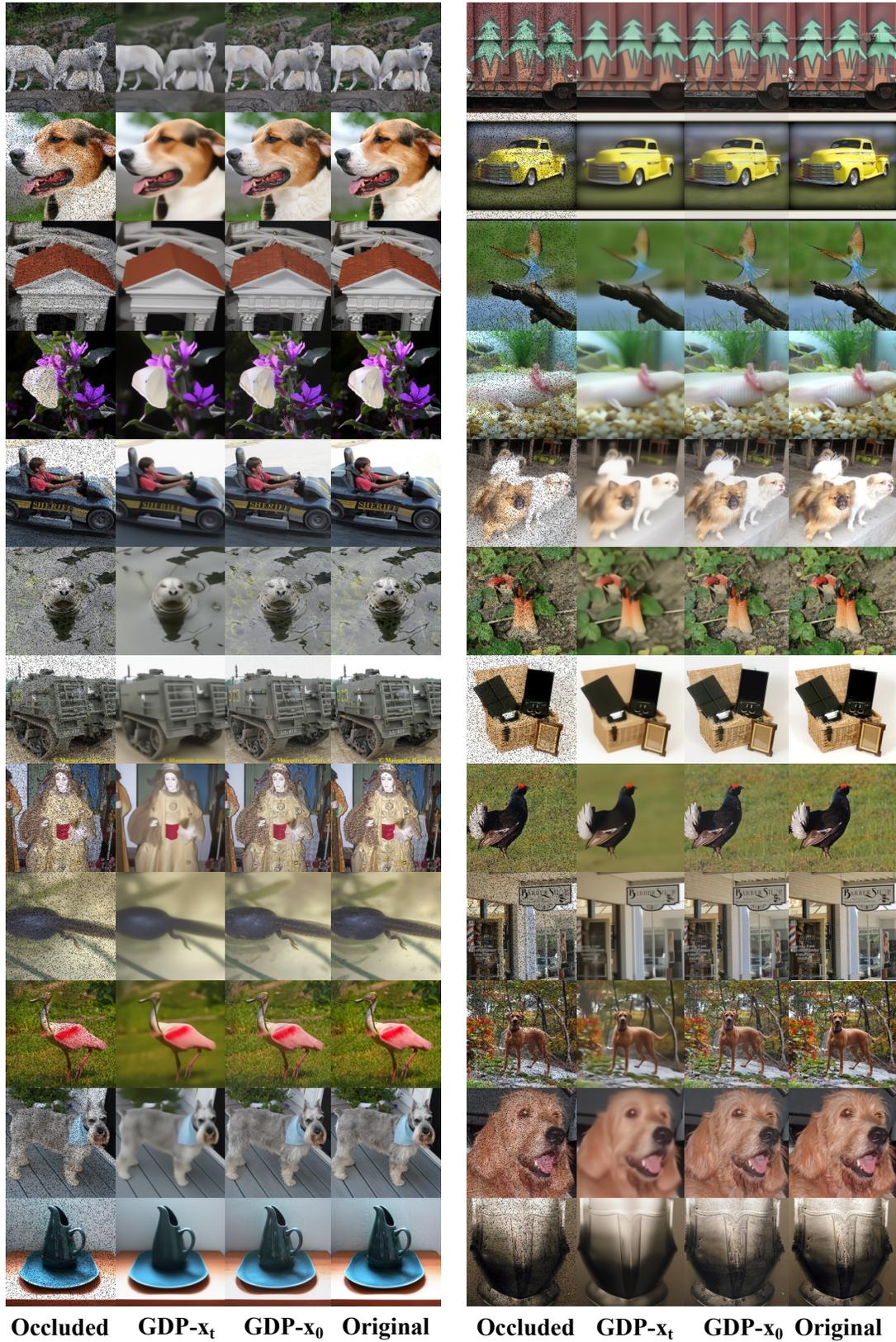


Figure 23. Uncurated samples from the 10% inpainting task on  $256 \times 256$  ImageNet 1K.

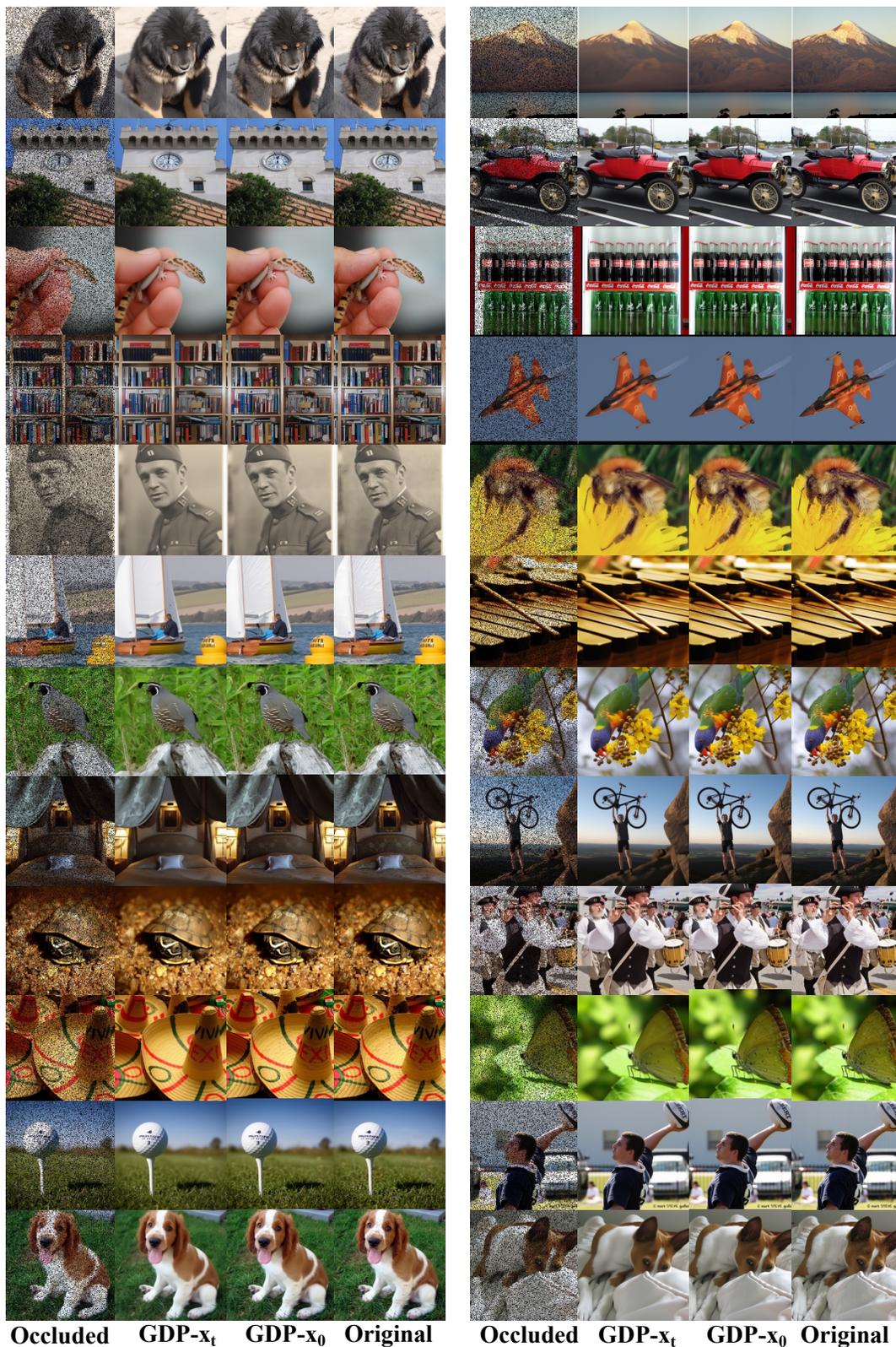


Figure 24. Uncurated samples from the 25% inpainting task on  $256 \times 256$  ImageNet 1K.



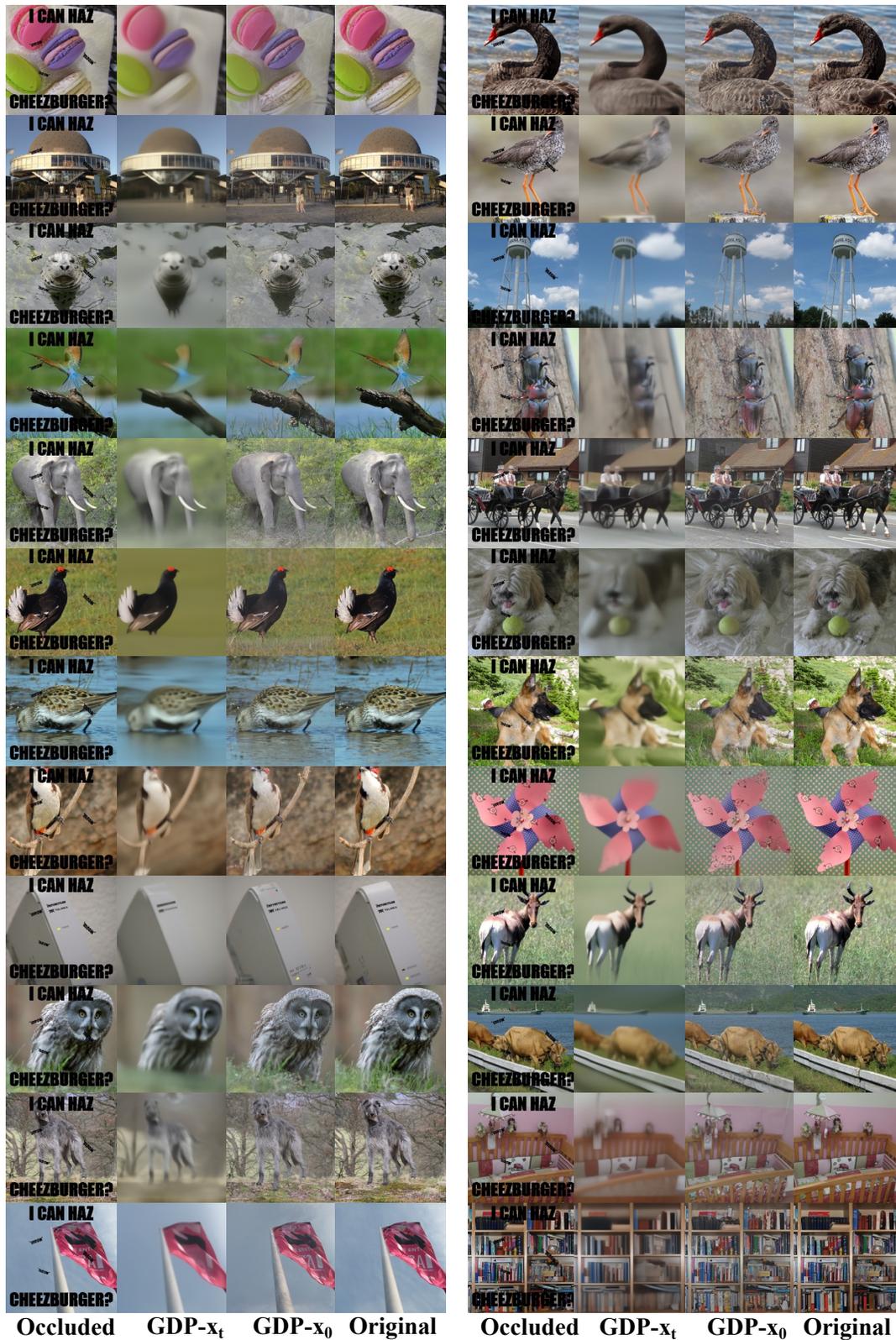


Figure 26. Uncurated samples from the inpainting task on  $256 \times 256$  ImageNet 1K.

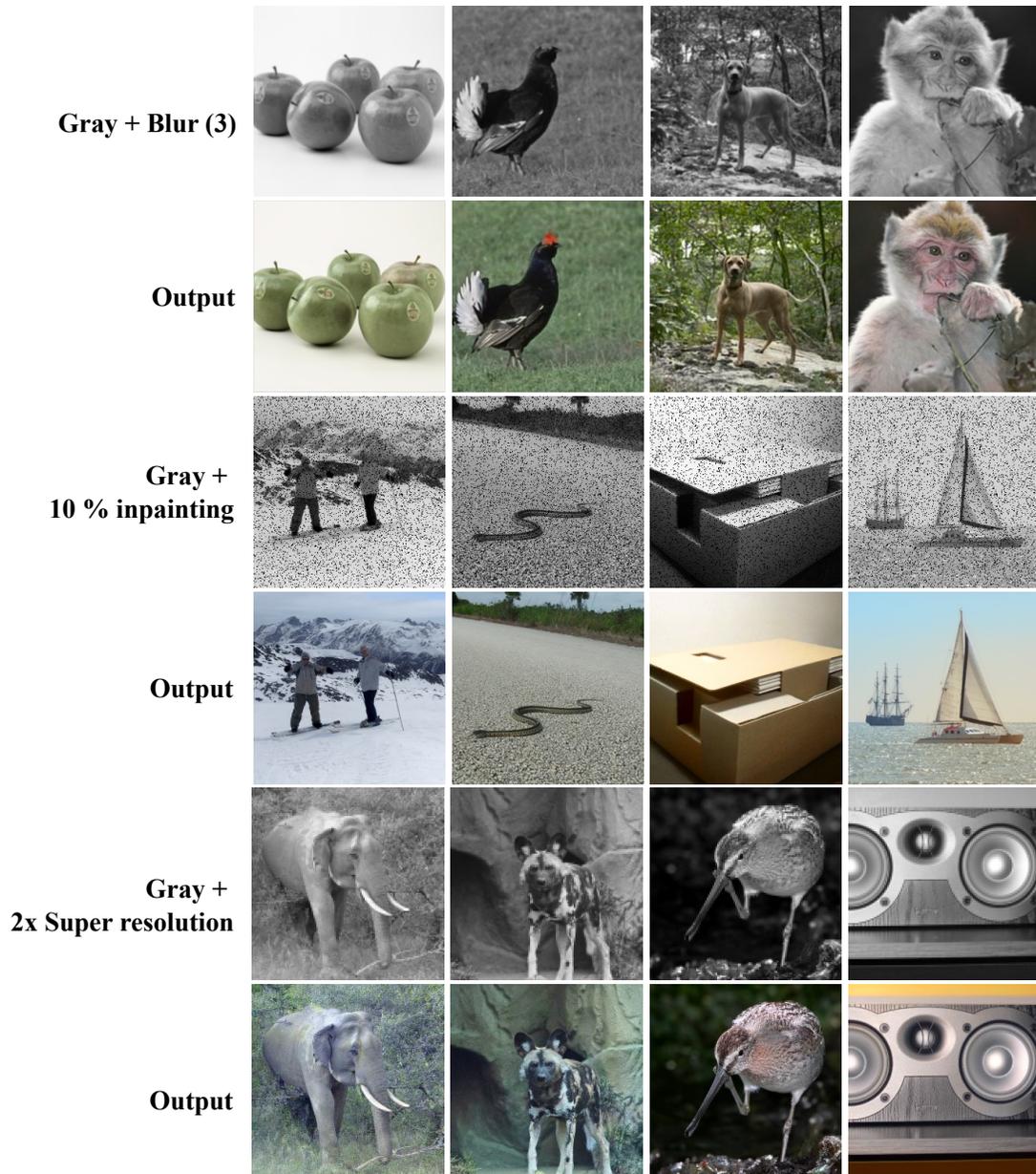
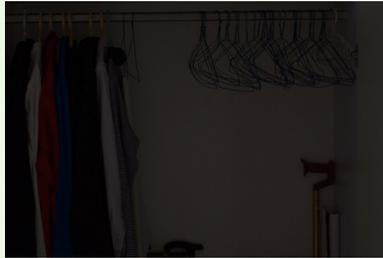


Figure 27. **Samples from the multi-degradation tasks on  $256 \times 256$  ImageNet 1K.** It is shown that GDP can recover the corrupted images undergoing multiple degradations, such as gray + blur, gray + inpainting, and gray + down-sampling. It is noted that multi-linear degradation should be only one degradation model that will damage the contents of the images. In other words, the restoration will be more difficult if two content-damaged degradations occur at the same time, such as down-sampling + mask.

## LOL Dataset

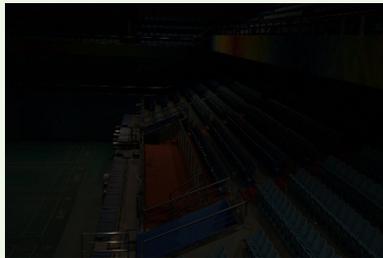
Original



GDP-x<sub>0</sub>



Original



GDP-x<sub>0</sub>



Figure 28. Results of low-light image enhancement on LOL dataset.

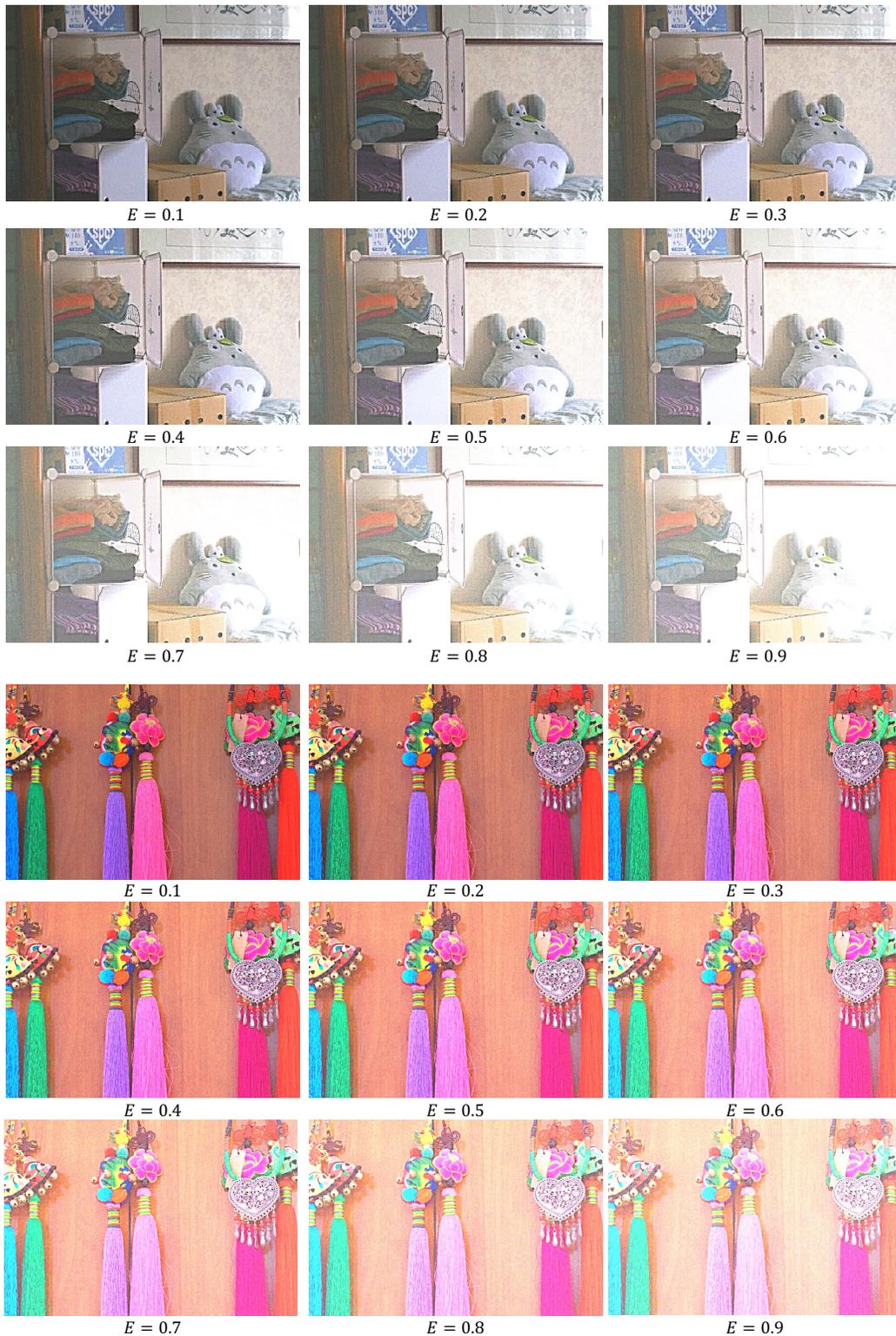


Figure 29. **Results of light control on LOL dataset.** We can adjust the brightness of the generated images with the help of Exposure Control Loss. Users can adjust the gray level  $E$  in the RGB color space to obtain the images according to their needs.

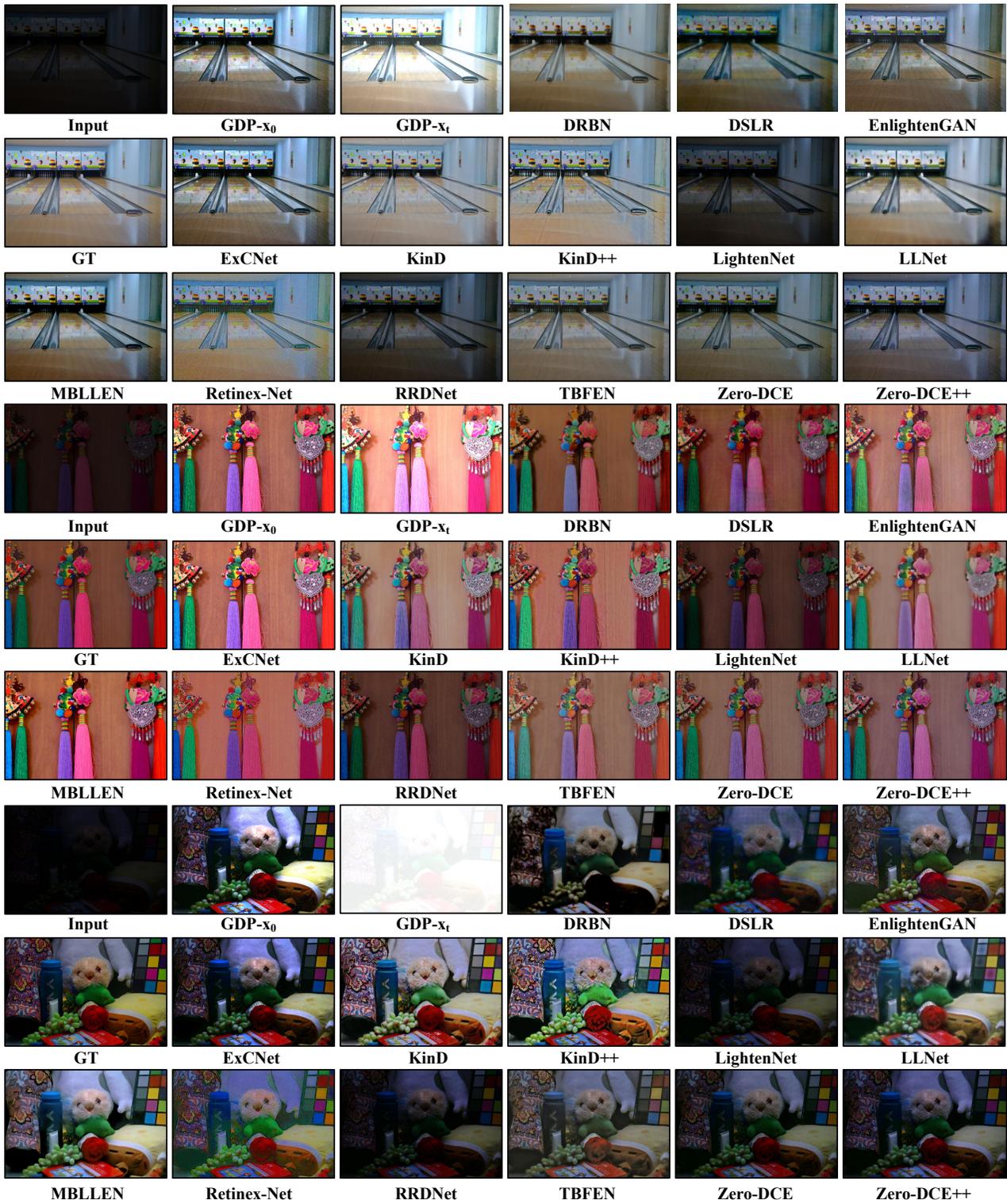


Figure 30. The comparison of our GDP and other methods on the LOL datasets towards low-light enhancement.



Figure 31. Results of low-light image enhancement on VE-LOL-L dataset.

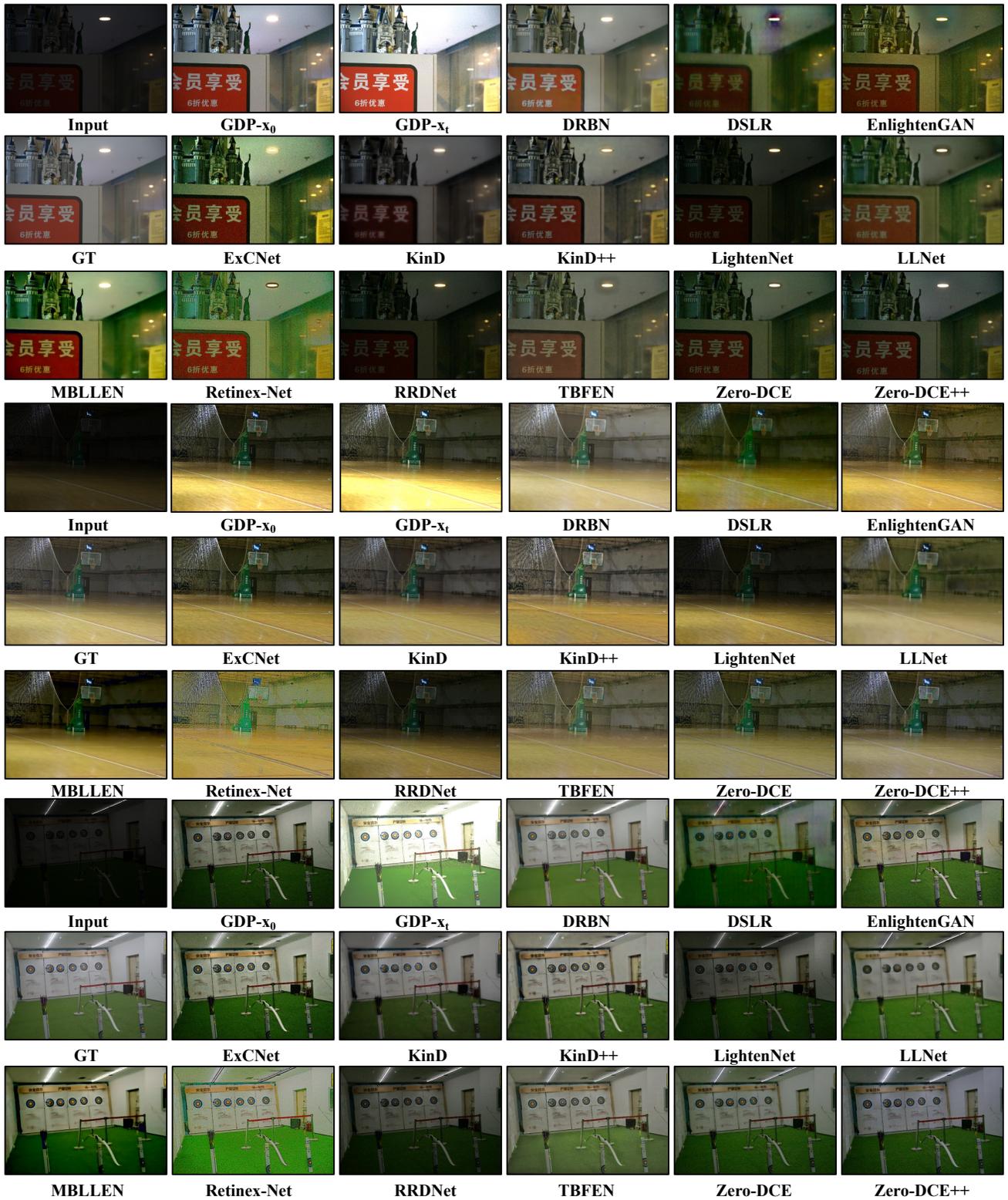


Figure 32. The comparison of our GDP and other methods on the VE-LOL-L datasets towards low-light enhancement.

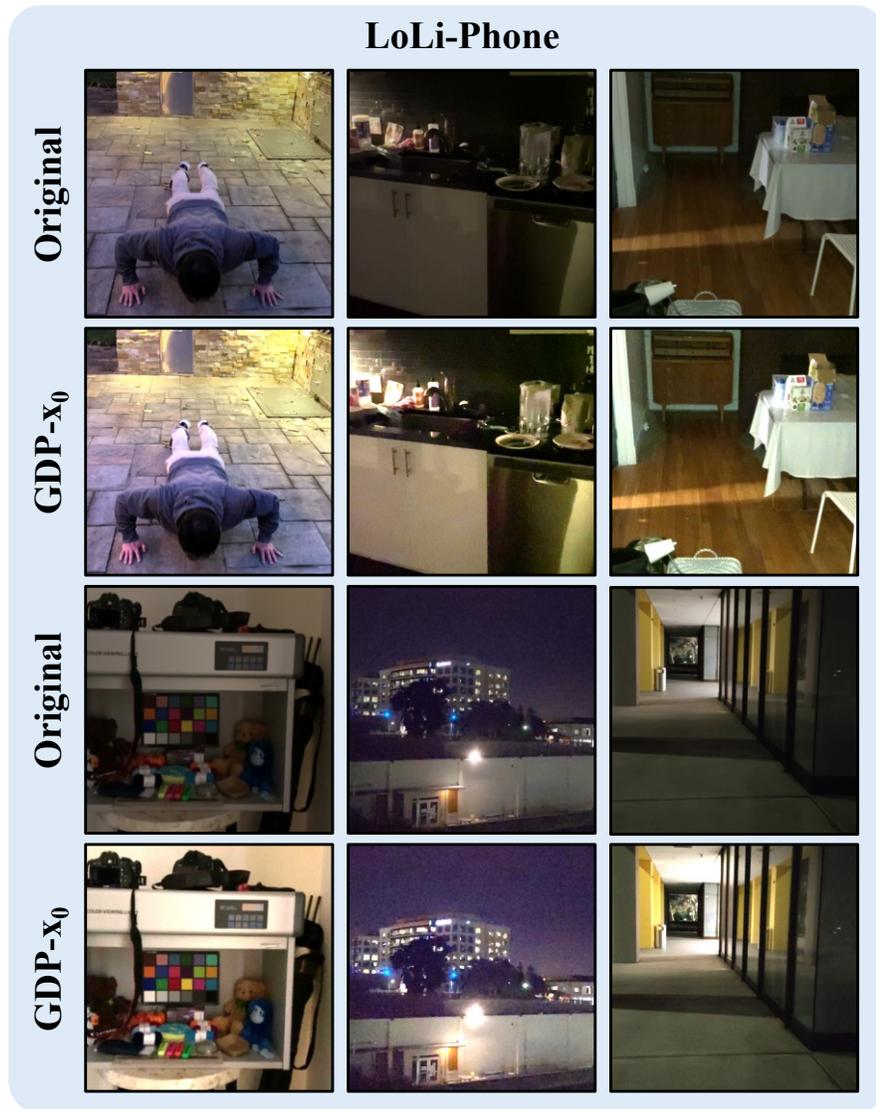


Figure 33. Results of low-light image enhancement on LoLi-Phone dataset.

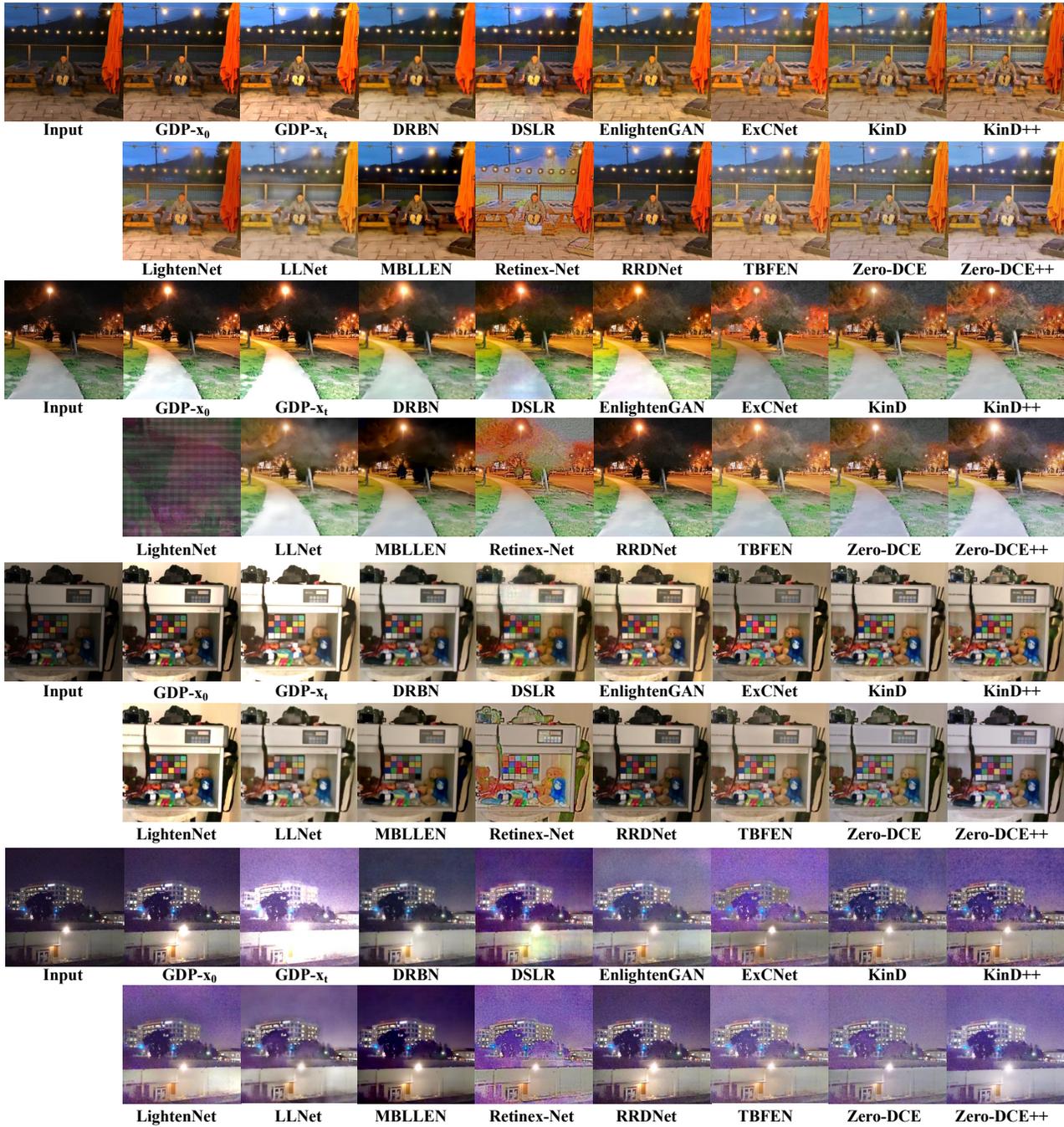
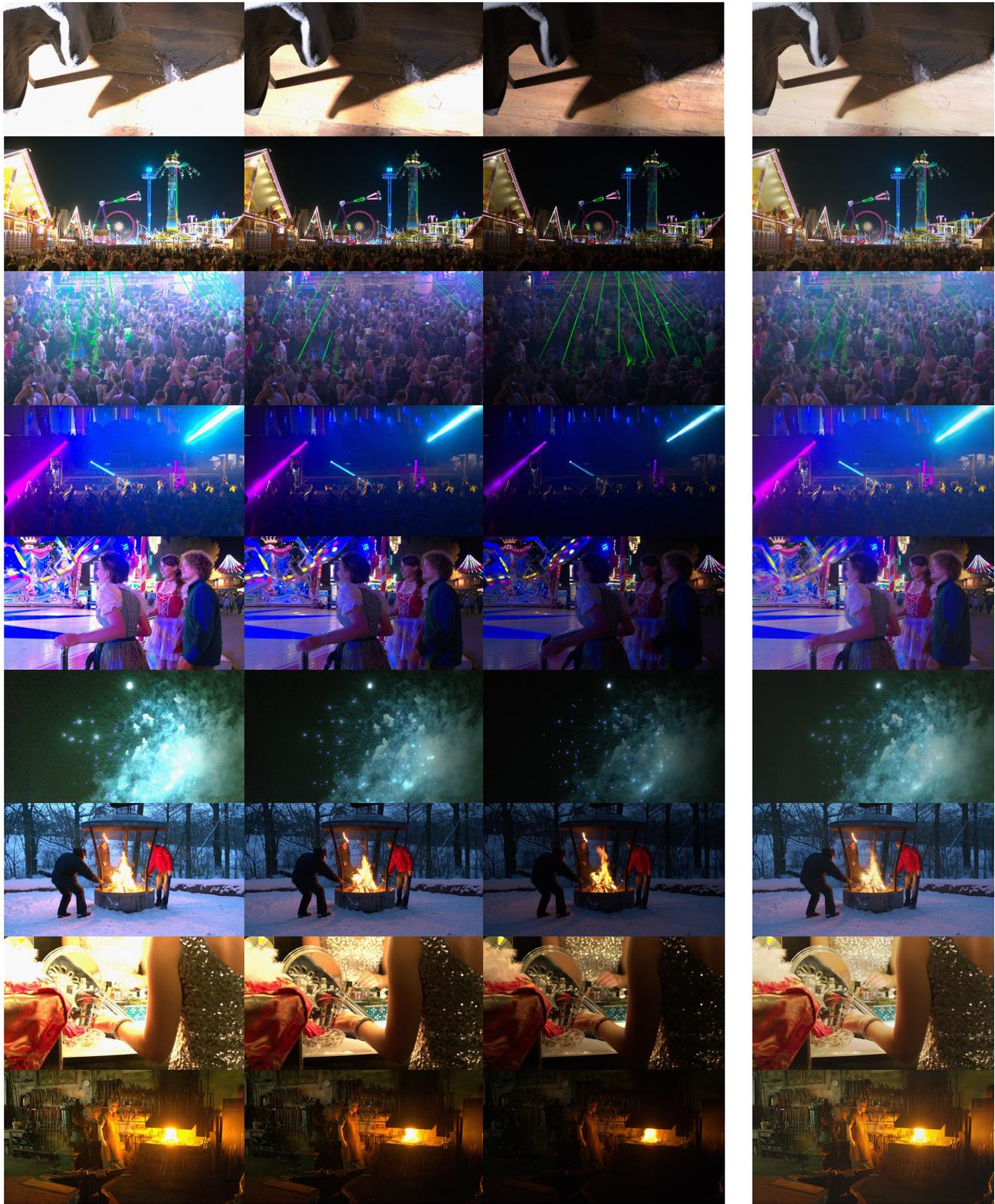


Figure 34. The comparison of our GDP and other methods on the LoLi-phone datasets towards low-light enhancement.



Long

Medium

Short

HDR-GDP- $x_0$

Figure 35. Results of HDR image recovery on NTIRE2021 dataset.



Figure 36. The comparison of HDR image recovery on NTIRE2021 dataset.



Figure 37. **Qualitative comparison of ablation study on LOL dataset.** Model A recovers the images in  $256 \times N$  or  $256 \times N$  sizes and is interpolated by the nearest neighbor to the original size. Model B is devised to naively restore the images from patches and patches where the parameters are not related. Model C is designed with fixed parameters for all patches in the images.



**HDR-GDP- $x_0$**



**Model A  
Naïve restoration**



**Model B  
Fixed parameters**

Figure 38. Qualitative comparison of ablation study on NTIRE2021 dataset.