

ADIR: Adaptive Diffusion for Image Reconstruction

Shady Abu-Hussein
Tel Aviv University
Tel Aviv, Israel
shady.abh@gmail.com

Tom Tirer
Bar-Ilan University
Ramat Gan, Israel
tirer.tom@gmail.com

Raja Giryes
Tel Aviv University
Tel Aviv, Israel
raja@tauex.tau.ac.il

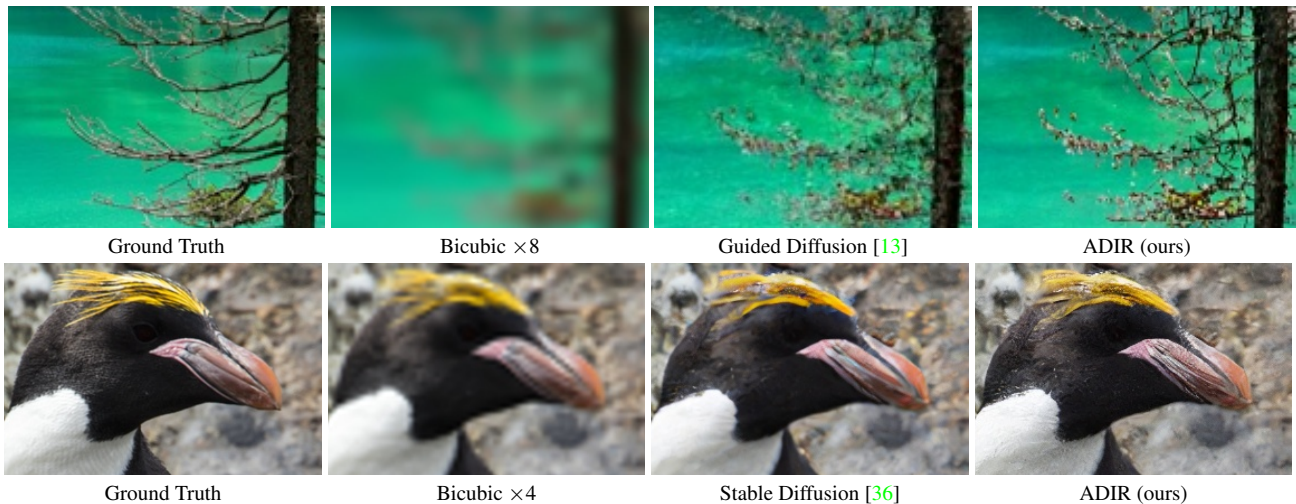


Figure 1. Super-resolution with scale factors 4 and 8, using Stable Diffusion and Guided Diffusion, and our method ADIR. The adaptability of ADIR allows reconstructing finer details.

Abstract

In recent years, denoising diffusion models have demonstrated outstanding image generation performance. The information on natural images captured by these models is useful for many image reconstruction applications, where the task is to restore a clean image from its degraded observations. In this work, we propose a conditional sampling scheme that exploits the prior learned by diffusion models while retaining agreement with the observations. We then combine it with a novel approach for adapting pre-trained diffusion denoising networks to their input. We examine two adaption strategies: the first uses only the degraded image, while the second, which we advocate, is performed using images that are “nearest neighbors” of the degraded image, retrieved from a diverse dataset using an off-the-shelf visual-language model. To evaluate our method, we test it on two state-of-the-art publicly available diffusion models, Stable Diffusion and Guided Diffusion. We show that our proposed ‘adaptive diffusion for image reconstruction’ (ADIR) approach achieves a significant improvement in the

super-resolution, deblurring, and text-based editing tasks. Our code and additional results are available online in the [project web page](#).

1. Introduction

Image reconstruction problems appear in a wide range of applications, where one would like to reconstruct an unknown clean image $\mathbf{x} \in \mathbb{R}^n$ from its degraded version $\mathbf{y} \in \mathbb{R}^m$, which can be noisy, blurry, low-resolution, etc. The acquisition (forward) model of \mathbf{y} in many important degradation settings can be formulated using the following linear model

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}, \quad (1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the measurement operator (blurring, masking, sub-sampling, etc.) and $\mathbf{e} \in \mathbb{R}^m \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_m)$ is the measurement noise. Typically, just fitting the observation model is not sufficient for recovering \mathbf{x} successfully. Thus, prior knowledge on the characteristics of \mathbf{x} is needed.

Over the past decade, many works suggested solving the

inverse problem in Eq. (1) using a single execution of a deep neural network, which has been trained on pairs of clean $\{\mathbf{x}_i\}$ images and their degraded versions $\{\mathbf{y}_i\}$ obtained by applying the forward model (1) on $\{\mathbf{x}_i\}$ [14, 28, 45, 51]. Yet, these approaches tend to overfit the observation model and perform poorly on setups that have not been considered in training [20, 41, 47]. Tackling this limitation with dedicated training for each application is not only computationally inefficient but also often impractical. This is because the exact observation model may not be known before inference time.

Several alternative approaches such as Deep Image Prior [48], zero-shot-super-resolution [41] or GSURE-based test-time optimization [1] rely solely on the observation image \mathbf{y} . They utilize the implicit bias of deep neural networks and gradient-based optimizers, as well as the self-recurrence of patterns in natural images when training a neural model directly on the observation and in this way reconstruct the original image. Although these methods are not limited to a family of observation models, they usually perform worse than data-driven methods, since they do not exploit the robust prior information that the unknown image \mathbf{x} share with external data that may contain images of the same kind. Other popular approaches, which do exploit external data while remaining flexible to the observation model, use deep models for imposing only the prior. Such methods typically use pretrained deep denoisers [3, 46, 50, 52] or generative adversarial networks (GANs) [8, 12, 21] within their optimization schemes, where consistency of the reconstruction with the observation \mathbf{y} is maintained by minimizing a data-fidelity term.

Recently, diffusion models [13, 18, 31, 42] have shown remarkable capabilities in generating high-fidelity images. These models are based on constructing a Markov chain diffusion process of length T for each training sample. In addition, they learn the reverse process, namely, the denoising operation between each two points in the chain. Sampling images via pretrained diffusion models is performed by starting with a pure white Gaussian noise image. This is followed by progressively sampling a less noisy image, given the previous one, until reaching a clean image after T iterations. Since diffusion models capture prior knowledge of the data, one may utilize them as deep priors/regularization for inverse problems of the form (1) [5, 9, 22, 29, 36, 44].

In this work, we propose an Adaptive Diffusion framework for Image Reconstruction (ADIR). First, we devise a diffusion guidance sampling scheme that solves (1) while restricting the reconstruction of \mathbf{x} to the range of a pretrained diffusion model. Our scheme is based on novel modifications to the guidance used in [13] (see Section 3.2 for details). Then, we propose two techniques that use the observations \mathbf{y} to adapt the diffusion network to patterns beneficial for recovering the unknown \mathbf{x} . Adapting the model’s parameters is based either directly on \mathbf{y} or on K ex-

ternal images similar to \mathbf{y} in some neural embedding space that is not sensitive to the degradation of \mathbf{y} . These images may be retrieved from a diverse dataset and the embedding can be calculated using an off-the-shelf encoder model for images such as CLIP [34]. As far as we know, our ADIR framework is the first *adaptive* diffusion approach for inverse problems. We evaluate it with two state-of-the-art diffusion models: Stable Diffusion [36] and Guided Diffusion [13], and show that it outperforms existing methods in the super-resolution and deblurring tasks. Finally, we demonstrate that our proposed adaptation strategy can be employed also for text-guided image editing.

2. Related Work

Diffusion models In recent years, many works utilized diffusion models for image manipulation and reconstruction tasks [22, 23, 36, 39, 49], where a denoising network is trained to learn the prior distribution of the data. At test time, some conditioning mechanism is combined with the learned prior to solve very challenging imaging tasks [4, 5, 10].

In [39, 49] the problems of deblurring and super-resolution were considered. Then, a diffusion model has been trained to perform this task where instead of adding noise at each of its steps, a blur or downsampling is performed. In this way, the model learns to carry out the deblurring or super-resolution task directly. Notice that these models are trained for one specific task and cannot be used for the other as is.

The closest works to us are [16, 23, 40]. These very recent concurrent works consider the task of image editing and perform an adaptation of the used diffusion model using the provided input and external data. Yet, notice that neither of these works consider the task of image reconstruction as we do here or apply our proposed sampling scheme for this task.

Image-Adaptive Reconstruction Adaptation of pretrained deep models, which serve as priors in inverse problems, to the unknown true \mathbf{x} through its observations at hand was proposed in [21, 47]. These works improve the reconstruction performance by fine-tuning the parameters of pretrained deep denoisers [47] and GANs [21] via the observed image \mathbf{y} instead of keeping them fixed during inference time. The image-adaptive GAN (IAGAN) approach [21] has led to many follow up works with different applications, e.g., [6, 32, 33, 35]. Recently, it has been shown that one may even fine-tune a masked-autoencoder to the input data at test-time for improving the adaptivity of classification neural networks to new domains [15].

In this paper we consider test-time adaptation of diffusion models for inverse problems. As far as we know, adaptation of diffusion models has not been proposed. Further-

more, while existing works fine-tune the deep priors directly using \mathbf{y} , we propose an improved strategy where the tuning is based on K external images similar to \mathbf{y} that are automatically retrieved from an external dataset.

3. Method

We now turn to present our proposed approach. Yet, before doing that we first provide a short introduction to regular denoising diffusion models. After that we describe our proposed strategy for modifying the sampling scheme of diffusion models for the task of image reconstruction. Finally, we present our suggested adaptation scheme.

3.1. Denoising Diffusion Models

Denoising diffusion models [18, 42] are latent variable generative models, with latent variables $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T \in \mathbb{R}^n$ (the same dimensionality as the data $\mathbf{x} \sim q_{\mathbf{x}}$). Given a training sample $\mathbf{x}_0 \sim q_{\mathbf{x}}$, these models are based on constructing a diffusion process (forward process) of the variables $\mathbf{x}_{1:T} := \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$ as a Markov chain from \mathbf{x}_0 to \mathbf{x}_T of the form

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) := \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}), \quad (2)$$

where $q(\mathbf{x}_t|\mathbf{x}_{t-1}) := \mathcal{N}(\sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}_n)$, and $0 < \beta_1 < \dots < \beta_T = 1$ is the diffusion variance schedule (hyperparameters of the model). Note that sampling $\mathbf{x}_t|\mathbf{x}_0$ can be done via a simplified way using the parametrization [18]:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad (3)$$

where $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$. The goal of these models is to learn the distribution of the reverse chain from \mathbf{x}_T to \mathbf{x}_0 , which is parameterized as the Markov chain

$$p_{\theta}(\mathbf{x}_{0:T}) := p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t), \quad (4)$$

where $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) := \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$,

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) := \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}}\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)\right), \quad (5)$$

and θ denotes all the learnable parameters. Essentially, $\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t)$ is an estimator for the noise in \mathbf{x}_t (up to scaling).

The parameters θ of the diffusion model $(\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$ are optimized by minimizing evidence lower bound [42], a simplified score-matching loss [18, 43], or a combination of both [13, 31]. For example, the simplified loss involves the minimization of

$$\ell_{\text{simple}} = \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}, t)\|_2^2 \quad (6)$$

in each training iteration, where \mathbf{x}_0 is drawn from the training data, t uniformly drawn from $\{1, \dots, T\}$ and the noise $\boldsymbol{\epsilon}$ is drawn from $\mathcal{N}(\mathbf{0}, \mathbf{I}_n)$.

Given a trained diffusion model $(\boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$, one may generate a sample \mathbf{x}_0 from the learned data distribution p_{θ} by initializing $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ and running the reverse diffusion process by sampling

$$\mathbf{x}_{t-1} \sim \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)), \quad (7)$$

where $0 < t \leq T$ and $\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)$ is defined in (5).

The class-guided sampling method that has been proposed in [13] modifies the sampling procedure in (7) by adding to the mean of the Gaussian a term that depends on the gradient of an offline-trained classifier, which has been trained using noisy images $\{\mathbf{x}_t\}$ for each t , and approximates the likelihood $p_{c|\mathbf{x}_t}$, where c is the desired class. This procedure has been shown to improve the quality of the samples generated for the learned classes.

3.2. Diffusion based Image Reconstruction

We turn now to extend the guidance method of [13] to image reconstruction. First, we conceptually generalize their framework to inverse problems of the form (1). Namely, given the observed image \mathbf{y} , we modify the guided reverse diffusion process to generate possible reconstructions of \mathbf{x} that are associated with \mathbf{y} rather than arbitrary samples of a certain class. Similarly to [13], ideally, the guiding direction at the t iteration should follow from (the gradient of) the likelihood function $p_{\mathbf{y}|\mathbf{x}_t}$.

The key difference between our framework and [13] is that we need to base our method on the specific degraded image \mathbf{y} rather than on a classifier that has been trained for each level of noise of $\{\mathbf{x}_t\}$. However, only the likelihood function $p_{\mathbf{y}|\mathbf{x}_0}$ is known, i.e., of the clean image \mathbf{x}_0 that is available only at the end of the procedure, and not for every $1 \leq t \leq T$. To overcome this issue, we propose a surrogate for the intermediate likelihood functions $p_{\mathbf{y}|\mathbf{x}_t}$. Our relaxation resembles the one in a recent concurrent work [11]. Yet, their sampling scheme is significantly different and has no adaptation ingredient.

Similar to [13], we guide the diffusion progression using the log-likelihood gradient. Formally, we are interested in sampling from the posterior

$$p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) \propto p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1})p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t), \quad (8)$$

where $p_{\mathbf{y}|\mathbf{x}_t}(\cdot|\mathbf{x}_t)$ is the distribution of \mathbf{y} conditioned on \mathbf{x}_t , and $p_{\theta}(\mathbf{x}_t|\mathbf{x}_{t+1}) = \mathcal{N}(\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t+1), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t+1))$ is the learned diffusion prior. For brevity, we omit the arguments of $\boldsymbol{\mu}_{\theta}$ and $\boldsymbol{\Sigma}_{\theta}$ in the rest of this subsection.

Under the assumption that the likelihood $\log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\cdot)$ has low curvature compared to $\boldsymbol{\Sigma}_{\theta}^{-1}$ [13], the following approximation using Taylor expansion around $\mathbf{x}_t = \boldsymbol{\mu}_{\theta}$ is

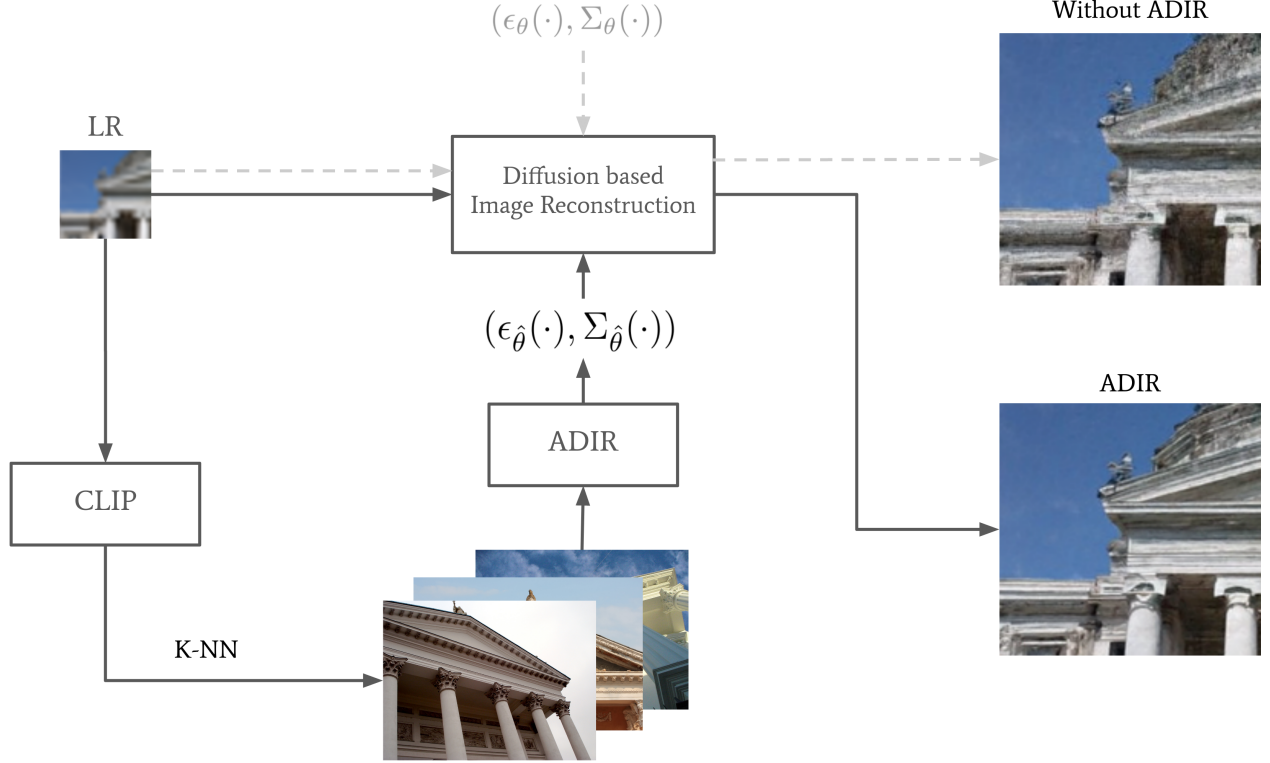


Figure 2. Diagram of our proposed method ADIR (Adaptive Diffusion for Image Reconstruction) applied to the super resolution task. Given a pretrained diffusion model $(\epsilon_\theta(\cdot), \Sigma_\theta(\cdot))$ and a Low Resolution (LR) image, we look for the K nearest neighbor images to the LR image, then using ADIR we adapt the diffusion model and use it for reconstruction.

valid

$$\begin{aligned} \log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t) &\approx \log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t)|_{\mathbf{x}_t=\boldsymbol{\mu}_\theta} \\ &+ (\mathbf{x}_t - \boldsymbol{\mu}_\theta)^\top \nabla_{\mathbf{x}_t} \log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t)|_{\mathbf{x}_t=\boldsymbol{\mu}_\theta} \\ &= (\mathbf{x}_t - \boldsymbol{\mu}_\theta)^\top \mathbf{g} + C_1, \end{aligned} \quad (9)$$

where $\mathbf{g} = \nabla_{\mathbf{x}_t} \log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t)|_{\mathbf{x}_t=\boldsymbol{\mu}_\theta}$, and C_1 is a constant that does not depend on \mathbf{x}_t . Then, similar to the computation in [13], we can use (9) to express the posterior (8) in the following form

$$\log(p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1})p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t)) \approx C_2 + \log p(\mathbf{z}), \quad (10)$$

where $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_\theta + \Sigma_\theta \mathbf{g}, \Sigma_\theta)$, and C_2 is some constant that does not depend on \mathbf{x}_t . Therefore, for conditioning the diffusion reverse process on \mathbf{y} , one needs to evaluate the derivative \mathbf{g} from a (different) log-likelihood function $\log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\cdot)$ at each iteration t .

Observe that we know the exact log-likelihood function for $t = 0$. Since the noise \mathbf{e} in (1) is white Gaussian with variance σ^2 , we therefore have following distribution

$$p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{A}\mathbf{x}, \sigma^2 \mathbf{I}_m) \propto e^{-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}. \quad (11)$$

In the denoising diffusion setup, \mathbf{y} is related to \mathbf{x}_0 using the observation model (1). Therefore,

$$\log p_{\mathbf{y}|\mathbf{x}_0}(\mathbf{y}|\mathbf{x}_0) \propto -\|\mathbf{A}\mathbf{x}_0 - \mathbf{y}\|_2^2. \quad (12)$$

However, we do not have tractable expressions for the likelihood functions $\{p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\cdot)\}_{t=1}^T$. Therefore, motivated by the expression above, we propose the following approximation

$$\log p_{\mathbf{y}|\mathbf{x}_t}(\mathbf{y}|\mathbf{x}_t) \approx \log p_{\mathbf{y}|\mathbf{x}_0}(\mathbf{y}|\hat{\mathbf{x}}_0(\mathbf{x}_t)), \quad (13)$$

where

$$\hat{\mathbf{x}}_0(\mathbf{x}_t) := (\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)) / \sqrt{\bar{\alpha}_t} \quad (14)$$

is an estimation of \mathbf{x}_0 from \mathbf{x}_t , which is based on the (stochastic) relation of \mathbf{x}_t and \mathbf{x}_0 in (3) and the random noise $\boldsymbol{\epsilon}$ is replaced by its estimation $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$.

From (11) and (13) it follows that \mathbf{g} in (9) can be approximated at each iteration t by evaluating (e.g., via automatic-differentiation)

$$\mathbf{g} \approx -\nabla_{\mathbf{x}_t} \|\mathbf{A}\hat{\mathbf{x}}_0(\mathbf{x}_t) - \mathbf{y}\|_2^2|_{\mathbf{x}_t=\boldsymbol{\mu}_\theta}. \quad (15)$$

Note that existing methods [11,22,44] either use a term that

Algorithm 1 Proposed guided diffusion sampling for image reconstruction given a diffusion model $(\epsilon_\theta(\cdot), \Sigma_\theta(\cdot))$, and a guidance scale s

Require: $(\epsilon_\theta(\cdot), \Sigma_\theta(\cdot)), \mathbf{y}, s$

- 1: $\mathbf{x}_T \leftarrow$ sample from $\mathcal{N}(\mathbf{0}, \mathbf{I}_n)$
- 2: **for all** t from T to 1 **do**
- 3: $\hat{\epsilon}, \hat{\Sigma} \leftarrow \epsilon_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)$
- 4: $\hat{\boldsymbol{\mu}} \leftarrow \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}}\hat{\epsilon})$
- 5: $\mathbf{y}_t \leftarrow \sqrt{\alpha_t}\mathbf{y} + \sqrt{1-\alpha_t}\mathbf{A}\hat{\epsilon}$
- 6: $\mathbf{g} \leftarrow -2\mathbf{A}^T(\mathbf{A}\hat{\boldsymbol{\mu}} - \mathbf{y}_t)$
- 7: $\mathbf{x}_{t-1} \leftarrow$ sample from $\mathcal{N}(\hat{\boldsymbol{\mu}} + s\hat{\Sigma}\mathbf{g}, \hat{\Sigma})$
- 8: **end for**
- 9: **return** \mathbf{x}_0

resembles (15) with the naive approximation $\hat{\mathbf{x}}_0(\mathbf{x}_t) = \mathbf{x}_t$ [22,44], or significantly modify (15) before computing it via automatic derivation framework [11] (we observed that trying to compute the exact (15) is unstable due to numerical issues). For example, in the official implementation of [11], which uses automatic derivation, the squaring of the norm in (15) is dropped even though this is not stated in their paper (otherwise, the reconstruction suffers from significant artifacts). In our case, we use the following relaxation to overcome the stability issue of using (15) directly. For a pretrained denoiser predicting ϵ_θ from \mathbf{x}_t and $0 < t \leq T$ we have

$$\begin{aligned} \|\mathbf{A}\hat{\mathbf{x}}_0(\mathbf{x}_t) - \mathbf{y}\|_2^2 &= \|\mathbf{A}(\mathbf{x}_t - \sqrt{1-\alpha_t}\epsilon_\theta)/\sqrt{\alpha_t} - \mathbf{y}\|_2^2 \\ &\propto \|\mathbf{A}\mathbf{x}_t - \sqrt{1-\alpha_t}\mathbf{A}\epsilon_\theta - \sqrt{\alpha_t}\mathbf{y}\|_2^2 \\ &= \|\mathbf{A}\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{y} - \sqrt{1-\alpha_t}\mathbf{A}\epsilon_\theta\|_2^2 \\ &= \|\mathbf{A}\mathbf{x}_t - \mathbf{y}_t\|_2^2, \end{aligned} \quad (16)$$

where $\mathbf{y}_t := \sqrt{\alpha_t}\mathbf{y} + \sqrt{1-\alpha_t}\mathbf{A}\epsilon_\theta$. Consequently, we propose to replace the expression for \mathbf{g} (the guiding likelihood direction at each iteration t) that is given in (15) with a surrogate obtained by evaluating the derivative of (16) w.r.t. \mathbf{x}_t , which is given by

$$\mathbf{g} \approx -2\mathbf{A}^T(\mathbf{A}\mathbf{x}_t - \mathbf{y}_t)|_{\mathbf{x}_t=\boldsymbol{\mu}_\theta} \quad (17)$$

that can be used for sampling the posterior distribution as detailed in Algorithm 1.

3.3. Adaptive Diffusion

Having defined the guided inverse diffusion flow for image reconstruction, we turn to discuss how one may adapt a given diffusion model to a given degraded image \mathbf{y} as defined in (1). Assume we have a pretrained diffusion model $(\epsilon_\theta(\cdot), \Sigma_\theta(\cdot))$, then the adaptation scheme is defined by the

following minimization problem

$$\hat{\theta} = \arg \min_{\theta} \sum_{t=1}^T \ell_{\text{simple}}(\mathbf{y}, \epsilon_\theta) \quad (18)$$

which can be solved efficiently using stochastic gradient descent, where at each iteration the gradient step is performed on a single term of the sum above, for $0 < t \leq T$ chosen randomly.

Adapting the denoising network to the measurement image \mathbf{y} , allows it to learn cross-scale features recurring in the image. Such an approach has been proven to be very helpful in reconstruction-based algorithms as shown in [21,47].

However, in some cases where the image does not satisfy the assumption of recurring patterns across scales, this approach can lose some of the sharpness captured in training. Therefore, in this work we extend the approach to few-shot fine-tuning adaptation, where instead of solving (18) w.r.t. \mathbf{y} , we propose an algorithm for retrieving K images similar to \mathbf{x} from a large dataset of diverse images, using off-the-shelf embedding distance.

Let $(\xi_v(\cdot), \xi_\ell(\cdot))$ be some off-the-shelf multi-modal encoder trained on visual-language modalities. Examples include CLIP [34], BLIP [27], and CyCLIP [17]). Let $\xi_v(\cdot)$ and $\xi_\ell(\cdot)$ be the visual and language encoders respectively. Then, given a large diverse dataset of natural images, we propose to retrieve K images, denoted by $\{\mathbf{z}_k\}_{k=1}^K$, with minimal embedding distance from \mathbf{y} . Formally, let \mathcal{D}_{IA} be an arbitrary external dataset, then

$$\begin{aligned} \{\mathbf{z}_k\}_{k=1}^K &= \{\mathbf{z}_1, \dots, \mathbf{z}_K \mid \phi_\xi(\mathbf{z}_1, \mathbf{y}) \leq \dots \leq \phi_\xi(\mathbf{z}_K, \mathbf{y}) \\ &\leq \phi_\xi(\mathbf{z}, \mathbf{y}), \forall \mathbf{z} \in \mathcal{D}_{\text{IA}} \setminus \{\mathbf{z}_1, \dots, \mathbf{z}_K\}\}, \end{aligned} \quad (19)$$

where $\phi_\xi(\mathbf{a}, \mathbf{b}) = 2 \arcsin(0.5\|\xi(\mathbf{a}) - \xi(\mathbf{b})\|_2)$ is the spherical distance and ξ can be either the visual or language encoder depending on the provided conditioning of the application.

After retrieving K -NN images $\{\mathbf{z}_k\}_{k=1}^K$ from \mathcal{D}_{IA} , we fine-tune the diffusion model on them, which adapts the denoising network to the context of \mathbf{y} , where we use the loss in (6) to modify the denoiser parameters θ . We refer to this K-NN based adaptation technique as ADIR (Adaptive Diffusion for Image Reconstruction), which is described schematically in Figure 2.

4. Experiments

We evaluate our method on two state-of-the-art diffusion models, Guided Diffusion (GD) [13] and Stable Diffusion (SD) [36], showing results on super-resolution and deblurring. In addition, we show how the adaptive diffusion can be used for the task of text-based editing using stable diffusion. In the [project web page](#) we provide more comparisons and setups.



Figure 3. Image deblurring using Guided Diffusion approach from section 3.2 and ADIR, using the unconditional model from [13]. The degradation is performed using 5×5 uniform blur filter with 10 levels of additive Gaussian noise. Note the better quality of our method.

Guided diffusion [13] provides several models with a conditioning mechanism built-in the denoiser. However, in our case, we perform the conditioning using the log-likelihood term. Therefore, we used the unconditional model that was trained on ImageNet [37] and produces images of size 256×256 . In the original work, the conditioning for generating an image from an arbitrary class was performed using a classifier trained to classify the noisy sample \mathbf{x}_t directly, where the log-likelihood derivative can be obtained by deriving the corresponding logits w.r.t. \mathbf{x}_t directly. In our setup, the conditioning is performed using \mathbf{g} in (17), where \mathbf{A} is defined by the reconstruction task, which we specify in the sequel.

In addition to GD, we demonstrate the improvement that can be achieved using stable diffusion [36], where we use the publicly available super-resolution and text-based editing models for it. Instead of training the denoiser on the natural images domain directly, they suggest to use a Variational Auto Encoder (VAE) and train the denoiser using a latent representation of the data. Note that the lower dimensionality of the latent enables the network to be trained on higher resolutions.

In all cases, we adapt the diffusion models in the image adaptive scheme presented in section 3.3, using the Google Open Dataset [26] as the external dataset \mathcal{D}_{IA} , from which we retrieve K images, where $K = 20$ for GD and $K = 50$ for SD (several examples of retrieved images are shown in Figure 7). For optimizing the network parameters we use Adam [25] and stabilize the training using Exponential

Moving Average (EMA). The specific implementation configurations are detailed in Table 4. We run all of our experiments on a NVIDIA RTX A6000 48GB card, which allows us to fine-tune the models by randomly sampling a batch of 6 images from $\{\mathbf{z}_k\}_{k=1}^K$, where in each iterations we use the same $0 < t \leq T$ for images in the batch.

	Bicubic	Guided Diffusion	IA	ADIR (GD)
SRx4	29.847 / 3.617	64.630 / 4.876	63.973/4.807	68.756 / 5.163
SRx8	17.915 / 3.523	54.406 / 4.433	-	58.580 / 4.476

Table 1. Super resolution with a scale factor of 4 ($128^2 \rightarrow 512^2$) results for the unconditional guided diffusion model [13]. The results are averaged on the first 50 images of the DIV2K validation set [2]. We compare our method to the baseline approach presented in Section 3.2, Bicubic upsampling, and the Image Adaptive (IA) method relying solely on \mathbf{y} with no external data. We use both the AVA-MUSIQ and KonIQ-MUSIQ [24] for evaluation.

	Bicubic	Stable Diffusion	ADIR (SD)
SRx4	36.953 / 4.159	69.184 / 5.073	72.703 / 5.545

Table 2. Super resolution with a scale factor of 4 ($256^2 \rightarrow 1024^2$) using the Stable Diffusion SR model [36]. Similar to Table 1, the results are averaged on the first 50 images of the DIV2K validation set [2]. We compare our method to the baseline approach presented by Stable Diffusion and Bicubic upsampling. We use both the AVA-MUSIQ and KonIQ-MUSIQ [24] for evaluation.

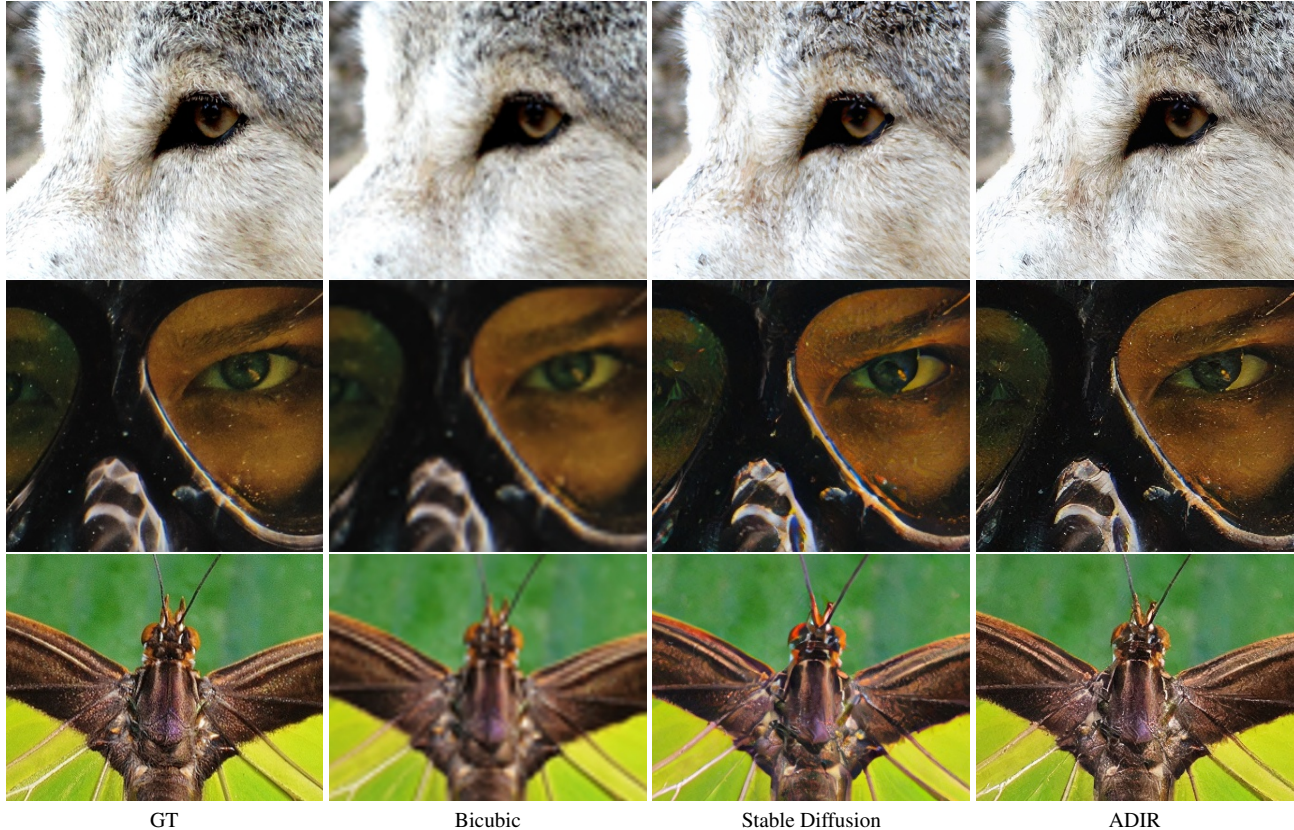


Figure 4. Comparison of super resolution ($256^2 \rightarrow 1024^2$) results of Stable Diffusion model [36] and our method (ADIR). As can be seen from the images, our method outperforms Stable Diffusion in both sharpness and reconstructing details.

	Guided Diffusion	ADIR (GD)
Uniform Deblur (256)	49.195 / 4.196	56.874 / 4.403
Uniform Deblur (512)	58.665 / 4.812	61.655 / 4.880
Gaussian deblur (256)	48.114 / 4.013	51.804 / 4.194

Table 3. Deblurring with 5×5 box filter and 10 noise levels results for the unconditional guided diffusion model [13]. Similar to SR in Table 1, the results are averaged on the first 50 images of the DIV2K validation set [2]. We compare our method to the baseline approach presented in Section 3.2. We use both the AVA-MUSIQ and KonIQ-MUSIQ [24] for evaluation.

4.1. Super Resolution

In the Super-Resolution (SR) task one would like to reconstruct a high resolution image x from its low resolution image y , where in this case A represents an anti-aliasing filter followed by sub-sampling with stride γ , which we refer to as the scaling factor. In our use-case we employ a bicubic anti-aliasing filter and assume $e = 0$, similarly to most SR works.

Here we apply our approach on two different diffusion based SR methods, Stable Diffusion [36], and section 3.2

approach combined with the unconditional diffusion model from [13]. In Stable Diffusion, the low-resolution image y is upsampled from 256×256 to 1024×1024 , while in Guided Diffusion we use the unconditional model trained on 256×256 images, to upscale y from 128×128 to 512×512 resolution. When adapting Stable diffusion, we downsample random crops of the K -NN images using A , which we encode using the VAE and plug into the network conditioning mechanism. We fine-tune both models using random crops of the K -NN images, to which we then add noise using the scheduler provided by each model.

The perception preference of generative models-based image reconstruction has been seen in many works [7,8,21]. Therefore, we chose a perception-based measure to evaluate the performance of our method. Specifically, we use the state-of-the-art AVA-MUSIQ and KonIQ-MUSIQ perceptual quality assessment measures [24], which are state-of-the-art image quality assessment measures. We report our results using the two measures averaged on the first 50 validation images of the DIV2K [2] dataset. As can be seen in Tables 1, 2, our method significantly outperforms both Stable Diffusion and GD-based reconstruction approaches. We compare our SR results to Stable Diffusion SR and Guided

Diffusion without adaptation, as well as using Image Adaptation (IA) performed on \mathbf{y} with no external data. The latter is done only for guided diffusion and show inferior performance compared to using external data. Therefore, in the other experiments we use only the external data for improving the optimization. A clear dominance of our method can be seen in the tables.

We also show qualitative results in Figures 1 and 4. As can be seen, our method achieves superior restoration quality, where in some cases it restores fine details that were blurred out in the acquisition of the ground-truth image.

4.2. Deblurring

In deblurring, \mathbf{y} is obtained by applying a uniform blur filter of size 5×5 on \mathbf{x} , followed by adding measurement noise $\mathbf{e} \sim \mathcal{N}(0, \sigma^2 I_n)$, which in our setting is $\sigma = 10$.

We apply the approach in section 3.2 on Guided Diffusion unconditional model [13] to solve the task. In this case, \mathbf{A} can be implemented by applying the blur kernel on a given image.

We use the unconditional diffusion model provided by GD [13], which was trained on 256×256 size images. Yet, in our tests, we solve the deblurring task on images of size 256×256 as well as 512×512 , which emphasizes the remarkable benefit of the adaptation, as it allows the model to generalize to resolutions not seen in training.

Similar to SR, in Table 3 we report the KonIQ-MUSIQ and AVA-MUSIQ [24] measures, averaged on the first 50 DIV2K validation images [2], where we compare our approach to the guided diffusion reconstruction without image adaptation. A visual comparisons are also available in Figure 3, where a very significant improvement can be seen in both robustness to noise and reconstructing details.

4.3. Text-Guided Editing

Text-guided image editing is the task of completing a masked region of \mathbf{x} according to a prompt provided by the user. In this case, the diffusion model needs predict objects and textures correspondent to the provided prompt, therefore we chose to adapt the network on $\{\mathbf{z}_k\}_{k=1}^K$ retrieved using the text encoder, i.e. by solving (19) using ξ_ℓ . For evaluating our method on this application, we use the state-of-the-art inpainting model of Stable Diffusion [36]. Where \mathbf{y} encoded and concatenated with the mask resized to latent dimension, which are then plugged to the denoising network. When adapting the network, we follow the training scheme of Stable Diffusion, where we use random masks and the classifier-free conditioning approach [19] used for training Stable Diffusion, where the text embedding is randomly chosen to either be the encoded prompt or the embedding of an empty prompt. Notice that we cannot compare to [16, 23, 40] as there is no code available for them. For some of them, we do not even have access to the diffusion

model that they adapt [38]

Figure 6 presents the editing results and compares to both stable diffusion and GLIDE. GLIDE is the basis of the popular DALL-E-2 model. The images of GLIDE are taken from the paper. We use ADIR with stable diffusion and optimize them using the same seed.

Since Stable Diffusion was trained using a lossy latent representation with smaller dimensionality than the data, it is clear that GLIDE can achieve better results. However, because our method adapts the network to a specific scenario, it enables the model to produce cleaner and more accurate generations, as can be seen in Figure 6. In the first image we see that Stable Diffusion adds an object that does not blend well and has artifacts, while when combined with our approach the quality improves significantly. Similarly, in the second image we see that Stable Diffusion produces an inaccurate edit, where it adds a brown hair instead of red hair. This is again improved by our adaptation method.

5. Conclusion

We have presented the Adaptive Diffusion Image Reconstruction (ADIR) method, in which we improve the reconstruction results in several imaging tasks using off-the-shelf diffusion models. We have demonstrated how our adaptation can significantly improve existing state-of-the-art methods, e.g. Stable Diffusion for super resolution, where the exploitation of external data with the same context as \mathbf{y} , combined with our adaptation scheme leads to a significant improvement. Specifically, the produced images are sharper and have more details than the original ground truth image.

One limitation of our approach is that as is the case with all diffusion models, there is randomness in the generation results. Therefore, they quality of the output may depend on the seed used. Yet, we still find that when we compare our approach and the baseline with the same seed, we get an improvement in the vast majority of cases. We provide additional examples in the [project web page](#) of such randomly generated pairs with different random seeds to further show the consistent improvement that can be achieved by our proposed strategy.

References

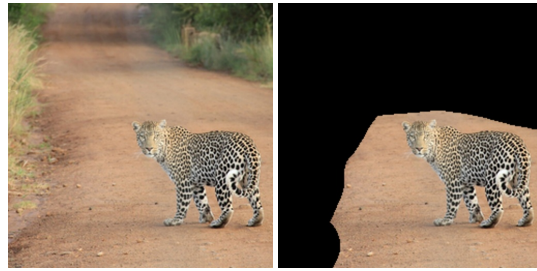
- [1] Shady Abu-Hussein, Tom Tirer, Se Young Chun, Yonina C Eldar, and Raja Giryes. Image restoration by deep projected gsure. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3602–3611, 2022. 2
- [2] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 126–135, 2017. 6, 7, 8

- [3] Siavash Arjomand Bigdeli, Matthias Zwicker, Paolo Favaro, and Meiguang Jin. Deep mean-shift priors for image restoration. *Advances in Neural Information Processing Systems*, 30, 2017. 2
- [4] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *arXiv preprint arXiv:2206.02779*, 2022. 2
- [5] 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. 2
- [6] Sayantan Bhadra, Weimin Zhou, and Mark A Anastasio. Medical image reconstruction with image-adaptive priors learned by use of generative adversarial networks. In *Medical Imaging 2020: Physics of Medical Imaging*, volume 11312, pages 206–213. SPIE, 2020. 2
- [7] Yochai Blau and Tomer Michaeli. The perception-distortion tradeoff. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6228–6237, 2018. 7
- [8] Ashish Bora, Ajil Jalal, Eric Price, and Alexandros G Dimakis. Compressed sensing using generative models. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 537–546. JMLR. org, 2017. 2, 7
- [9] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14347–14356. IEEE, 2021. 2
- [10] Hyungjin Chung, Eun Sun Lee, and Jong Chul Ye. Mr image denoising and super-resolution using regularized reverse diffusion. *arXiv preprint arXiv:2203.12621*, 2022. 2
- [11] Hyungjin Chung, Byeongsu Sim, Dohoon Ryu, and Jong Chul Ye. Improving diffusion models for inverse problems using manifold constraints. *arXiv preprint arXiv:2206.00941*, 2022. 3, 4, 5
- [12] Manik Dhar, Aditya Grover, and Stefano Ermon. Modeling sparse deviations for compressed sensing using generative models. In *International Conference on Machine Learning*, pages 1214–1223. PMLR, 2018. 2
- [13] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 1, 2, 3, 4, 5, 6, 7, 8
- [14] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2015. 2
- [15] Yossi Gandelsman, Yu Sun, Xinlei Chen, and Alexei A Efros. Test-time training with masked autoencoders. *arXiv preprint arXiv:2209.07522*, 2022. 2
- [16] Giorgio Giannone, Didrik Nielsen, and Ole Winther. Few-shot diffusion models. *10.48550/ARXIV.2205.15463*, 2022. 2, 8
- [17] Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan A Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *arXiv preprint arXiv:2205.14459*, 2022. 5
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 2, 3
- [19] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 8
- [20] Shady Abu Hussein, Tom Tirer, and Raja Giryes. Correction filter for single image super-resolution: Robustifying off-the-shelf deep super-resolvers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1428–1437, 2020. 2
- [21] Shady Abu Hussein, Tom Tirer, and Raja Giryes. Image-adaptive gan based reconstruction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3121–3129, 2020. 2, 5, 7
- [22] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *arXiv preprint arXiv:2201.11793*, 2022. 2, 4, 5
- [23] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. 2, 8
- [24] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5148–5157, 2021. 6, 7, 8
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [26] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020. 6, 14
- [27] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. *arXiv preprint arXiv:2201.12086*, 2022. 5
- [28] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 136–144, 2017. 2
- [29] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471, 2022. 2
- [30] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 13
- [31] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 2, 3

- [32] Yotam Nitzan, Kfir Aberman, Qiurui He, Orly Liba, Michal Yarom, Yossi Gandelsman, Inbar Mosseri, Yael Pritch, and Daniel Cohen-Or. Mystyle: A personalized generative prior. *arXiv preprint arXiv:2203.17272*, 2022. [2](#)
- [33] 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](#)
- [34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. [2](#), [5](#), [14](#)
- [35] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022. [2](#)
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#), [12](#), [13](#)
- [37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. [6](#)
- [38] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv:2205.11487*, 2022. [8](#)
- [39] 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. [2](#)
- [40] Shelly Sheynin, Oron Ashual, Adam Polyak, Uriel Singer, Oran Gafni, Eliya Nachmani, and Yaniv Taigman. Knn-diffusion: Image generation via large-scale retrieval. *arXiv:2204.02849*, 2022. [2](#), [8](#)
- [41] 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 (CVPR)*, pages 3118–3126, 2018. [2](#)
- [42] 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](#), [3](#)
- [43] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019. [3](#)
- [44] Yang Song, Liyue Shen, Lei Xing, and Stefano Ermon. Solving inverse problems in medical imaging with score-based generative models. *arXiv preprint arXiv:2111.08005*, 2021. [2](#), [4](#), [5](#)
- [45] Jian Sun, Wenfei Cao, Zongben Xu, and Jean Ponce. Learning a convolutional neural network for non-uniform motion blur removal. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 769–777, 2015. [2](#)
- [46] Tom Tirer and Raja Giryes. Image restoration by iterative denoising and backward projections. *IEEE Transactions on Image Processing*, 28(3):1220–1234, 2018. [2](#)
- [47] Tom Tirer and Raja Giryes. Super-resolution via image-adapted denoising cnns: Incorporating external and internal learning. *IEEE Signal Processing Letters*, 26(7):1080–1084, 2019. [2](#), [5](#)
- [48] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9446–9454, 2018. [2](#)
- [49] Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. Deblurring via stochastic refinement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16293–16303, 2022. [2](#)
- [50] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [2](#)
- [51] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a Gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE Transactions on Image Processing*, 26(7):3142–3155, 2017. [2](#)
- [52] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. Learning deep cnn denoiser prior for image restoration. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3929–3938, 2017. [2](#)

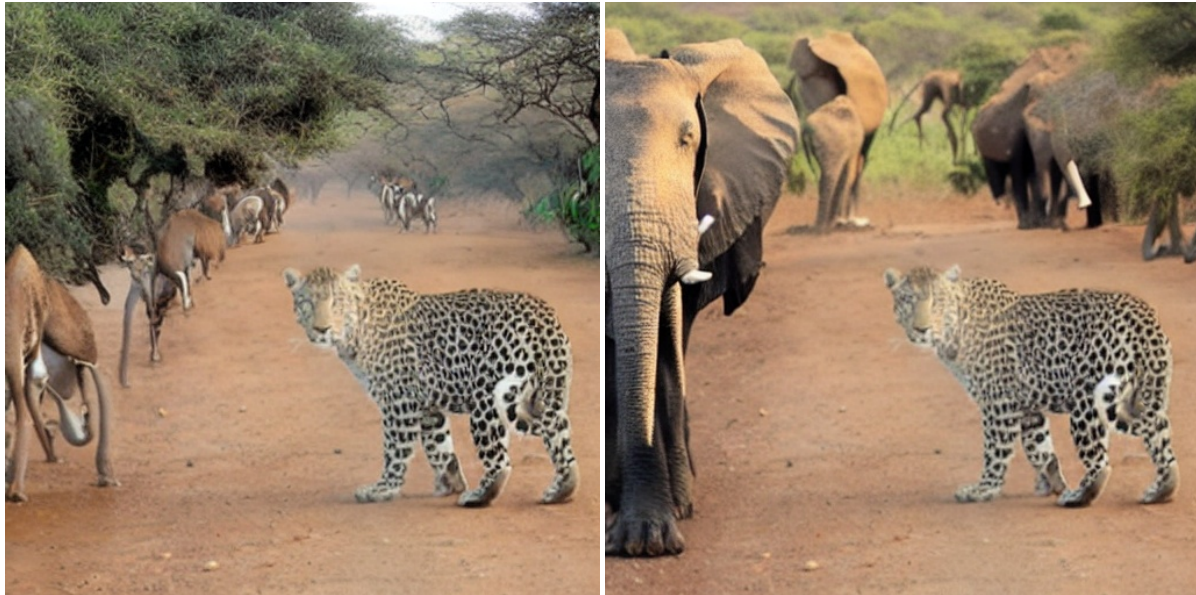
Configuration Method	IA iterations	learning rate	EMA	NN images	s	diffusion steps
Super Resolution x4						
GD	0	-	-	-	10	1000
IA-GD	100	10^{-4}	0.95	-	10	1000
ADIR-GD	400	10^{-4}	0.8	20	10	1000
SD	-	-	-	-	-	500
ADIR-SD	400	10^{-4}	0.95	50	-	500
Super Resolution x8						
GD	0	-	-	-	20	1000
ADIR-GD	400	10^{-4}	0.8	20	20	1000
Deblurring						
GD	0	-	-	-	10	1000
ADIR-GD	400	10^{-4}	0.8	20	10	1000
Text based editing						
SD	-	-	-	-	-	500
ADIR-SD	1000	10^{-6}	0	50	-	500

Table 4. Configurations used for ADIR.

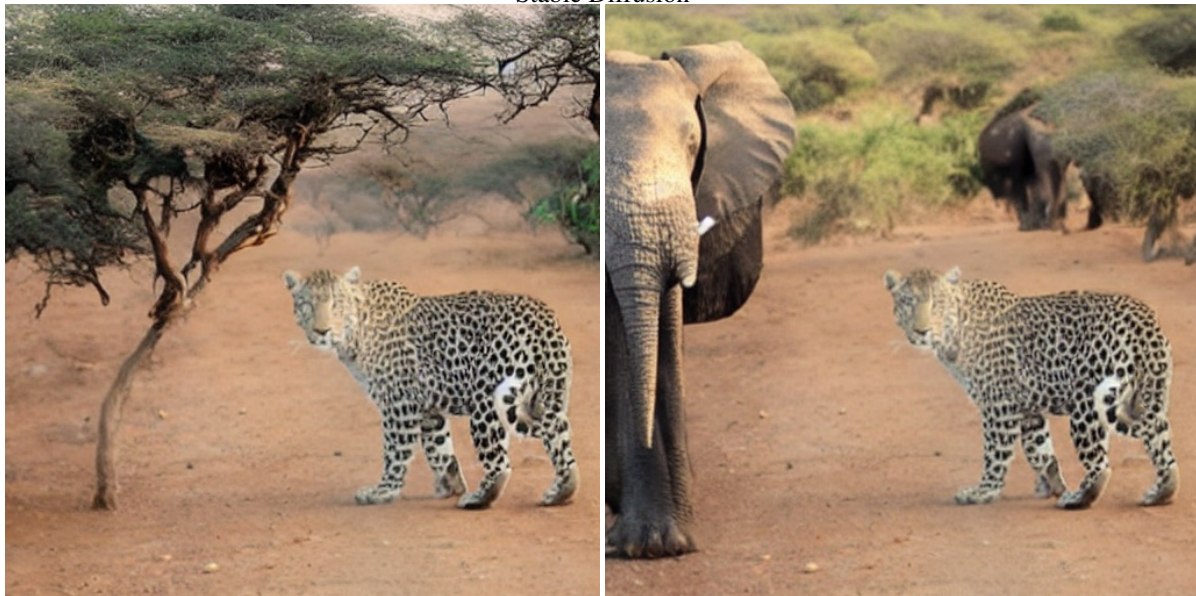


Original Image

Masked



Stable Diffusion



ADIR

Figure 5. Text-based image editing comparison between Stable Diffusion [36] and ADIR, using the prompt “Africa” for two different seeds.

“A vase of flowers on the table of a living room”



“A man with red hair”



“An old car in a snowy forest”



input

GLIDE

Stable Diffusion

ADIR

Figure 6. Text-based image editing comparison between GLIDE (full) [30], Stable Diffusion [36] and ADIR applied to the Stable Diffusion model. The images are taken from [30], since their official high-resolution model was not publicly released. As can be seen, our method produces more realistic images in cases where Stable Diffusion either was not accurate (brown hair instead of red) or in terms of artifacts.



Figure 7. Examples of images retrieved from Google Open Dataset [26] using CLIP [34] for super resolution with scale factor of 8 ($64^2 \rightarrow 512^2$).