

# Exploiting Deep Generative Prior for Versatile Image Restoration and Manipulation

Xingang Pan<sup>1</sup>, Xiaohang Zhan<sup>1</sup>, Bo Dai<sup>1</sup>,  
Dahua Lin<sup>1</sup>, Chen Change Loy<sup>2</sup>, and Ping Luo<sup>3</sup>

<sup>1</sup> The Chinese University of Hong Kong  
{px117,zx017,bdai,dhlin}@ie.cuhk.edu.hk

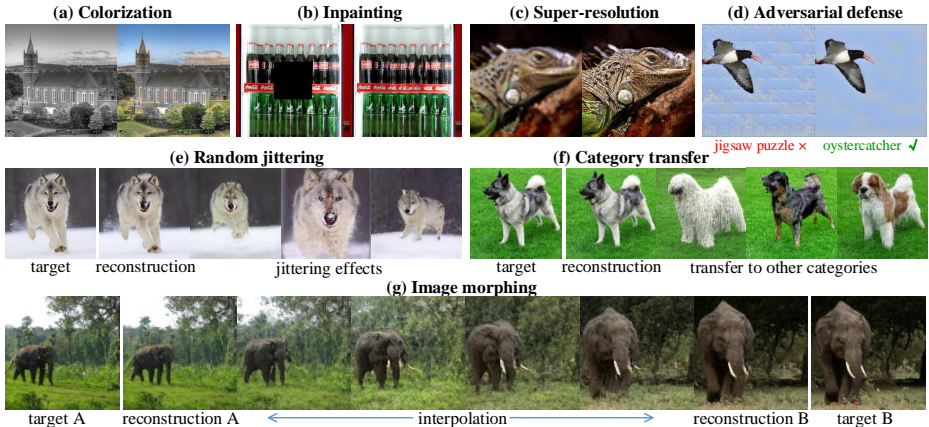
<sup>2</sup> Nanyang Technological University <sup>3</sup> The University of Hong Kong  
ccloy@ntu.edu.sg pluo@cs.hku.hk

**Abstract.** Learning a good image prior is a long-term goal for image restoration and manipulation. While existing methods like deep image prior (DIP) capture low-level image statistics, there are still gaps toward an image prior that captures rich image semantics including color, spatial coherence, textures, and high-level concepts. This work presents an effective way to exploit the image prior captured by a generative adversarial network (GAN) trained on large-scale natural images. As shown in Fig. 1, the deep generative prior (DGP) provides compelling results to restore missing semantics, *e.g.*, color, patch, resolution, of various degraded images. It also enables diverse image manipulation including random jittering, image morphing, and category transfer. Such highly flexible restoration and manipulation are made possible through relaxing the assumption of existing GAN-inversion methods, which tend to fix the generator. Notably, we allow the generator to be fine-tuned on-the-fly in a progressive manner regularized by feature distance obtained by the discriminator in GAN. We show that these easy-to-implement and practical changes help preserve the reconstruction to remain in the manifold of nature image, and thus lead to more precise and faithful reconstruction for real images. Code is available at <https://github.com/XingangPan/deep-generative-prior>.

## 1 Introduction

Learning image prior models is important to solve various tasks of image restoration and manipulation, such as *image colorization* [25,43], *image inpainting* [41], *super-resolution* [14,26], and *adversarial defense* [33]. In the past decades, many image priors [30,47,15,18,31] have been proposed to capture certain statistics of natural images. Despite their successes, these priors often serve a dedicated purpose. For instance, markov random field [30,47,15] is often used to model the correlation among neighboring pixels, while dark channel prior [18] and total variation [31] are developed for dehazing and denoising respectively.

There is a surge of interest to seek for more general priors that capture richer statistics of images through deep learning models. For instance, the seminal work on deep image prior (DIP) [36] showed that the structure of a randomly initialized Convolutional Neural Network (CNN) implicitly captures texture-level

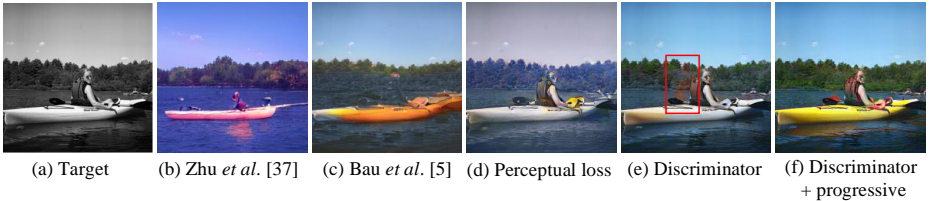


**Fig. 1.** These image restoration(a)(b)(c)(d) and manipulation(e)(f)(g) effects are achieved by leveraging the rich generative prior of a GAN. The GAN does not see these images during training

image prior, thus can be used for restoration by fine-tuning it to reconstruct a corrupted image. SinGAN [34] further shows that a randomly-initialized generative adversarial network (GAN) model is able to capture rich patch statistics after training from a single image. These priors have shown impressive results on some low-level image restoration and manipulation tasks like super-resolution and harmonizing. In both the representative works, the CNN and GAN are trained from a single image of interest from scratch.

In this study, we are interested to go one step further, examining how we could leverage a GAN [16] trained on large-scale natural images for richer priors beyond a single image. GAN is a good approximator for natural image manifold. By learning from large image datasets, it captures rich knowledge on natural images including color, spatial coherence, textures, and high-level concepts, which are useful for broader image restoration and manipulation effects. Specifically, we take a collapsed image (*e.g.*, gray-scale image) as a partial observation of the original natural image, and reconstruct it in the observation space (*e.g.*, gray-scale space) with the GAN, the image prior of the GAN would tend to restore the missing semantics (*e.g.*, color) in a faithful way to match natural images. Despite its enormous potentials, it remains a challenging task to exploit a GAN as a prior for general image restoration and manipulation. The key challenge lies in the needs in coping with arbitrary images from different tasks with distinctly different natures. The reconstruction also needs to produce sharp and faithful images obeying the natural image manifold.

An appealing option for our problem is GAN-inversion [45,10,2,5]. Existing GAN-inversion methods typically reconstruct a target image by optimizing over the latent vector, *i.e.*,  $\mathbf{z}^* = \arg \min_{\mathbf{z} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, G(\mathbf{z}; \theta))$ , where  $\mathbf{x}$  is the target image,  $G$  is a fixed generator,  $\mathbf{z}$  and  $\theta$  are the latent vector and generator parameters, respectively. In practice, we found that this strategy fails in dealing



**Fig. 2.** Comparison of various methods in reconstructing a gray image under the gray-scale observation space using a GAN. Conventional GAN-inversion strategies like (b)[45] and (c)[5] produce imprecise reconstruction for the existing semantics. In this work, we relax the generator so that it can be fine-tuned on-the-fly, achieving more accurate reconstruction as in (d)(e)(f), of which optimization is based on (d) VGG perceptual loss, (e) discriminator feature matching loss, and (f) combined with progressive reconstruction, respectively. We highlight that discriminator is important to preserve the generative prior so as to achieve better restoration for the missing information (*i.e.*, color). The proposed progressive strategy eliminates the ‘information lingering’ artifacts as in the red box in (e)

with complex real-world images. In particular, it often results in mismatched reconstructions, whose details (*e.g.*, objects, texture, and background) appear inconsistent with the original images, as Fig. 2 (b)(c) show. On one hand, existing GAN-inversion methods still suffer from the issues of mode collapse and limited generator capacity, affecting their capability in capturing the desired data manifold. On the other hand, perhaps a more crucial limitation is that when a generator is fixed, the GAN is inevitably limited by the training distribution and its inversion cannot faithfully reconstruct unseen and complex images. It is infeasible to carry such assumptions while using a GAN as prior for general image restoration and manipulation.

Despite the gap between the approximated manifold and the real one, the GAN generator still captures rich statistics of natural images. In order to make use of these statistics while avoiding the aforementioned limitation, in this paper we present a relaxed and more practical reconstruction formulation for mining the priors in GAN. Our first reformulation is to allow the generator parameters to be fine-tuned on the target image on-the-fly, *i.e.*,  $\theta^*, \mathbf{z}^* = \arg \min_{\theta, \mathbf{z}} \mathcal{L}(\mathbf{x}, G(\mathbf{z}; \theta))$ . This lifts the constraint of confining the reconstruction within the training distribution. Relaxing the assumption with fine-tuning, however, is still not sufficient to ensure good reconstruction quality for arbitrary target images. We found that fine-tuning using a standard loss such as perceptual loss [22] or mean squared error (MSE) in DIP could risk wiping out the originally rich priors. Consequently, the reconstruction may become increasingly unnatural during the reconstruction of a degraded image. Fig. 2(d) shows an example, suggesting that a new loss and reconstruction strategy is needed.

Thus, in our second reformulation, we devise an effective reconstruction strategy that consists of two components:

1) *Feature matching loss from the coupled discriminator* - we make full use of the discriminator of a trained GAN to regularize the reconstruction. Note that

during training, the generator is optimized to mimic massive natural images via gradients provided by the discriminator. It is reasonable to still adopt the discriminator in guiding the generator to match a single image as the discriminator preserves the original parameter structure of the generator better than other distance metrics. Thus deriving a feature matching loss from the discriminator can help maintain the reconstruction to remain in the natural image space. Although the feature matching loss is not new in the literature [37], its significance to GAN reconstruction has not been investigated before.

2) *Progressive reconstruction* - we observe that a joint fine-tuning of all parameters of the generator could lead to ‘*information lingering*’, where missing semantics (*e.g.*, color) do not naturally change along with the content when reconstructing a degraded image. This is because the deep layers of the generator start to match the low-level textures before the high-level configurations are aligned. To address this issue, we propose a progressive reconstruction strategy that fine-tunes the generator gradually from the shallowest layers to the deepest layers. This allows the reconstruction to start with matching high-level configurations and gradually shift its focus on low-level details.

Thanks to the proposed techniques that enable faithful reconstruction while maintaining the generator prior, our approach, which we name as Deep Generative Prior (DGP), generalizes well to various kinds of image restoration and manipulation tasks, despite that our method is not specially designed for each task. When reconstructing a corrupted image in a task-dependent observation space, DGP tends to restore the missing information, while keeping existing semantic information unchanged. As shown in Fig. 1 (a)(b)(c), color, missing patches, and details of the given images are well restored, respectively. As illustrated in Fig. 1 (e)(f), we can manipulate the content of an image by tweaking the latent vector or category condition of the generator. Fig. 1 (g) shows that image morphing is possible by interpolating between the parameters of two fine-tuned generators and the corresponding latent vectors of these images. To our knowledge, it is the first time these jittering and morphing effects are achieved on a dataset with complex images like ImageNet [12]. We show more interesting examples in the experiments and Appendix.

## 2 Related Work

**Image Prior.** Image priors that describe various statistics of natural images have been widely adopted in computer vision, including markov random fields [30,47,15], dark channel prior [18], and total variation regularizer [31]. Recently, the work of deep image prior (DIP) [36] shows that image statistics are implicitly captured by the structure of CNN, which is also a kind of prior, and could be used to restore corrupted images. SinGAN [34] fine-tunes a randomly initialized GAN on patches of a single image, achieving various image editing or restoration effects. As DIP and SinGAN are trained from scratch, they have limited access to image statistics beyond the input image, which restrains their applicability in tasks such as image colorization. There are also other deep pri-

ors developed for low-level restoration tasks like deep denoiser prior [42,6] and TNRD [8], but competing with them is not our goal. Instead, our goal is to study and exploit the prior that is captured in GAN for versatile restoration as well as manipulation tasks. Existing attempts that use a pre-trained GAN as a source of image statistics include [4] and [20], which respectively applies to image manipulation, *e.g.*, editing partial areas of an image, and image restoration, *e.g.*, compressed sensing and super-resolution for human faces. As we will show in our experiments, by using a discriminator based distance metric and a progressive fine-tuning strategy, DGP can better preserve image statistics learned by the GAN and thus allows richer restoration and manipulation effects.

Recently, a concurrent work of multi-code GAN prior [17] also conducts image processing by solving the GAN-inversion problem. It uses multiple latent vectors to reconstruct the target image and keeps the generator fixed, while our method makes the generator image-adaptive by allowing it to be fine-tuned on-the-fly.

**Image Restoration and Manipulation.** In this paper we demonstrate the effect of applying DGP to multiple tasks of image processing, including image colorization [25], image inpainting [41], super-resolution [14,26], adversarial defence [33], and semantic manipulation [45,46,9]. While many task-specific models and loss functions have been proposed to pursue a better performance on a specific restoration task [25,43,41,14,26,33], there are also works that apply GAN and design task-specific pipelines to achieve various image manipulation effects [46,9,37,4,35,40], such as CycleGAN [46] and StarGAN [9]. In this work we are more interested in uncovering the potential of exploiting the GAN prior as a task-agnostic solution, where we propose several techniques to achieve this goal. Moreover, as shown in Fig. 1(e)(g), with an improved reconstruction process we successfully achieve image jittering and morphing on ImageNet, while previous methods are insufficient to handle these effects on such complex data.

**GAN-Inversion.** As mentioned in Sec.3, a straightforward way to utilize generative prior is conducting image reconstruction based on GAN-inversion. GAN-inversion aims at finding a vector in the latent space that best reconstructs a given image, where the GAN generator is fixed. Previous attempts either optimize the latent vector directly via gradient back-propagation [10,2] or leverage an additional encoder mapping images to latent vectors [45,13]. A more recent approach [5] proposes to add small perturbations to shallow blocks of the generator to ease the inversion task. While these methods could handle datasets with limited complexities or synthetic images sampled by the GAN itself, we empirically found in our experiments they may produce imprecise reconstructions for complex real scenes, *e.g.*, images in the ImageNet [12]. Recently, the work of StyleGAN [23] enables a new way for GAN-inversion by operating in intermediate latent spaces [1], but noticeable mismatches are still observed and the inversion for vanilla GAN (*e.g.*, BigGAN [7]) is still challenging. In this paper, instead of directly applying standard GAN-inversion, we devise a more practical way to reconstruct a given image using the generative prior, which is shown to achieve better reconstruction results.

### 3 Method

We first provide some preliminaries on DIP and GAN before discussing how we exploit DGP for image restoration and manipulation.

**Deep Image Prior.** Ulyanov *et al* [36] show that image statistics are implicitly captured by the structure of CNN. These statistics can be seen as a kind of image prior, which can be exploited in various image restoration tasks by tuning a randomly initialized CNN on the degraded image:  $\theta^* = \arg \min_{\theta} E(\hat{\mathbf{x}}, f(\mathbf{z}; \theta))$ ,  $\mathbf{x}^* = f(\mathbf{z}; \theta^*)$ , where  $E$  is a task-dependent distance metric,  $\mathbf{z}$  is a randomly chosen latent vector, and  $f$  is a CNN with  $\theta$  being its parameters.  $\hat{\mathbf{x}}$  and  $\mathbf{x}^*$  are the degraded image and restored image respectively. One limitation of DIP is that the restoration process mainly resorts to existing statistics in the input image, it is thus infeasible to apply DIP on tasks that require more general statistics, such as image colorization [25] and manipulation [45].

**Generative Adversarial Networks (GANs).** GANs are widely used for modeling complex data such as natural images [16,39,11,23]. In GAN, the underlying manifold of natural images is approximated by the combination of a parametric generator  $G$  and a prior latent space  $\mathcal{Z}$ , so that an image can be generated by sampling a latent vector  $\mathbf{z}$  from  $\mathcal{Z}$  and applying  $G$  as  $G(\mathbf{z})$ . GAN jointly trains  $G$  with a parametric discriminator  $D$  in an adversarial manner, where  $D$  is supposed to distinguish generated images from real ones. Although extensive efforts have been made to improve the power of GAN, there inevitably exists a gap between GAN’s approximated manifold and the actual one, due to issues such as insufficient capacity and mode collapse.

#### 3.1 Deep Generative Prior

Suppose  $\hat{\mathbf{x}}$  is obtained via  $\hat{\mathbf{x}} = \phi(\mathbf{x})$ , where  $\mathbf{x}$  is the original natural image and  $\phi$  is a degradation transform. *e.g.*,  $\phi$  could be a graying transform that turns  $\mathbf{x}$  into a grayscale image. Many tasks of image restoration can be regarded as recovering  $\mathbf{x}$  given  $\hat{\mathbf{x}}$ . A common practice is learning a mapping from  $\hat{\mathbf{x}}$  to  $\mathbf{x}$ , which often requires task-specific training for different  $\phi$ s. Alternatively, we can also 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  $\hat{\mathbf{x}}$ , viewing  $\hat{\mathbf{x}}$  as partial observations of  $\mathbf{x}$ .

While various priors have been proposed [30,36,34] in the second line of research, in this paper we are interested in studying a more generic image prior, *i.e.*, a GAN generator trained on large-scale natural images for image synthesis. Specifically, a straightforward realization is a reconstruction process based on GAN-inversion, which optimizes the following objective:

$$\begin{aligned} \mathbf{z}^* &= \arg \min_{\mathbf{z} \in \mathbb{R}^d} E(\hat{\mathbf{x}}, G(\mathbf{z}; \theta)), & \mathbf{x}^* &= G(\mathbf{z}^*; \theta), \\ &= \arg \min_{\mathbf{z} \in \mathbb{R}^d} \mathcal{L}(\hat{\mathbf{x}}, \phi(G(\mathbf{z}; \theta))), \end{aligned} \quad (1)$$

where  $\mathcal{L}$  is a distance metric such as the L2 distance,  $G$  is a GAN generator parameterized by  $\theta$  and trained on natural images. Ideally, if  $G$  is sufficiently

powerful that the data manifold of natural images is well captured in  $G$ , the above objective will drag  $\mathbf{z}$  in the latent space and locate the optimal natural image  $\mathbf{x}^* = G(\mathbf{z}^*; \boldsymbol{\theta})$ , which contains the missing semantics of  $\hat{\mathbf{x}}$  and matches  $\hat{\mathbf{x}}$  under  $\phi$ . For example, if  $\phi$  is a graying transform,  $\mathbf{x}^*$  will be an image with a natural color configuration subject to  $\phi(\mathbf{x}^*) = \hat{\mathbf{x}}$ . However, in practice it is not always the case.

As the GAN generator is fixed in Eq.(1) and its improved versions, *e.g.*, adding an extra encoder [45,13], these reconstruction methods based on the standard GAN-inversion suffer from an intrinsic limitation, *i.e.*, there is a gap between the approximated manifold of natural images and the actual one. On one hand, due to issues including mode collapse and insufficient capacity, the GAN generator cannot perfectly grasp the training manifold represented by a dataset of natural images. On the other hand, the training manifold itself is also an approximation of the actual one. Such two levels of approximations inevitably lead to a gap. Consequently, a sub-optimal  $\mathbf{x}^*$  is often retrieved, which often contains significant mismatches to  $\hat{\mathbf{x}}$ , especially when the original image  $\mathbf{x}$  is a complex image, *e.g.*, ImageNet [12] images, or an image located outside the training manifold. See Fig. 2 and existing literature [5,13] for an illustration.

**A Relaxed GAN Reconstruction Formulation.** Despite the gap between the approximated manifold and the real one, a well trained GAN generator still covers rich statistics of natural images. In order to make use of these statistics while avoiding the aforementioned limitation, we propose a relaxed GAN reconstruction formulation by allowing parameters  $\boldsymbol{\theta}$  of the generator to be moderately fine-tuned along with the latent vector  $\mathbf{z}$ . Such a relaxation on  $\boldsymbol{\theta}$  gives rise to an updated objective:

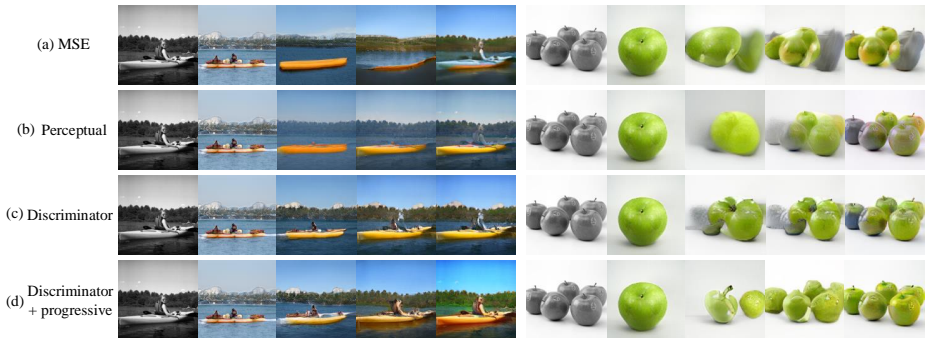
$$\boldsymbol{\theta}^*, \mathbf{z}^* = \arg \min_{\boldsymbol{\theta}, \mathbf{z}} \mathcal{L}(\hat{\mathbf{x}}, \phi(G(\mathbf{z}; \boldsymbol{\theta}))), \quad \mathbf{x}^* = G(\mathbf{z}^*; \boldsymbol{\theta}^*). \quad (2)$$

We refer to this updated objective as Deep Generative Prior (DGP). With this relaxation, DGP significantly improves the chance of locating an optimal  $\mathbf{x}^*$  for  $\hat{\mathbf{x}}$ , as fitting the generator to a single image is much more achievable than fully capturing a data manifold. Note that the generative prior buried in  $G$ , *e.g.*, its ability to output faithful natural images, might be deteriorated during the fine-tuning process. The key to preserve the generative prior lies in the design of a good distance metric  $\mathcal{L}$  and a proper optimization strategy.

### 3.2 Discriminator Guided Progressive Reconstruction

To fit the GAN generator to the input image  $\hat{\mathbf{x}}$  while retaining a natural output, in this section we introduce a discriminator based distance metric, and a progressive fine-tuning strategy.

**Discriminator Matters.** Given an input image  $\hat{\mathbf{x}}$ , DGP will start with an initial latent vector  $\mathbf{z}_0$ . In practice, we obtain  $\mathbf{z}_0$  by randomly sampling a few hundreds of candidates from the latent space  $\mathcal{Z}$  and selecting the one that its corresponding image  $G(\mathbf{z}; \boldsymbol{\theta})$  best resembles  $\hat{\mathbf{x}}$  under the metric  $\mathcal{L}$  we used in



**Fig. 3.** Comparison of different loss types when fine-tuning the generator to reconstruct the image

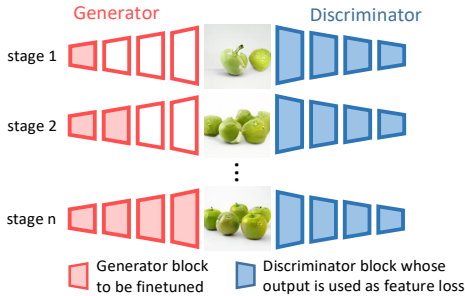
Eq.(2). As shown in Fig. 3, the choice of  $\mathcal{L}$  significantly affects the optimization of Eq.(2). Existing literature often adopts the Mean-Squared-Error (MSE) [36] or the AlexNet/VGGNet based Perceptual loss [22,45] as  $\mathcal{L}$ , which respectively emphasize the pixel-wise appearance and the low-level/mid-level texture. However, we empirically found using these metrics in Eq.(2) often cause unfaithful outputs at the beginning of optimization, leading to sub-optimal results at the end. We thus propose to replace them with a discriminator-based distance metric, which measures the L1 distance in the *discriminator feature space*:

$$\mathcal{L}(\mathbf{x}_1, \mathbf{x}_2) = \sum_{i \in \mathcal{I}} \|D(\mathbf{x}_1, i), D(\mathbf{x}_2, i)\|_1, \quad (3)$$

where  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are two images, corresponding to  $\hat{\mathbf{x}}$  and  $\phi(G(\mathbf{z}; \boldsymbol{\theta}))$  in Eq.1 and Eq.2, and  $D$  is the discriminator that is coupled with the generator.  $D(\mathbf{x}, i)$  returns the feature of  $\mathbf{x}$  at  $i$ -block of  $D$ , and  $\mathcal{I}$  is the index set of used blocks. Compared to the AlexNet/VGGNet based perceptual loss, the discriminator  $D$  is trained along with  $G$ , instead of being trained for a separate task.  $D$ , being a distance metric, thus is less likely to break the parameter structure of  $G$ , as they are well aligned during the pre-training. Moreover, we found the optimization of DGP using such a distance metric visually works like an image morphing process. *e.g.*, as shown in Fig. 3, the person on the boat is preserved and all intermediate outputs are all vivid natural images. It is worth pointing out again while the feature matching loss is not new, this is the first time it serves as a regularizer during GAN reconstruction.

**Progressive Reconstruction.** Typically, we will fine-tune all parameters of  $\boldsymbol{\theta}$  simultaneously during the optimization of Eq.(2). However, we observe an adverse effect of ‘*information lingering*’, where missing semantics (*e.g.* color) do not shift along with existing context. Taking Fig. 3 (c) as an example, the leftmost apple fails to inherit the green color of the initial apple when it emerges. One possible reason is deep blocks of the generator  $G$  start to match low-level textures before high-level configurations are completely aligned. To overcome this problem, we propose a progressive reconstruction strategy for some restoration tasks.





**Fig. 4.** Progressive reconstruction of the generator can better preserve the consistency between missing and existing semantics in comparison to simultaneous fine-tuning on all the parameters at once. Here the list of images shown in the middle are the outputs of the generator in different fine-tuning stages.

Specifically, as illustrated in Fig. 4, we first fine-tune the shallowest block of the generator, and gradually continue with blocks at deeper depths, so that DGP can control the global configuration at the beginning and gradually shift its attention to details at lower levels. A demonstration of the proposed strategy is included in Fig. 3 (d), where DGP splits the apple from one to two at first, then increases the number to five, and finally refines the details of apples. Compared to the non-progressive counterpart, such a progressive strategy better preserves the consistency between missing and existing semantics.

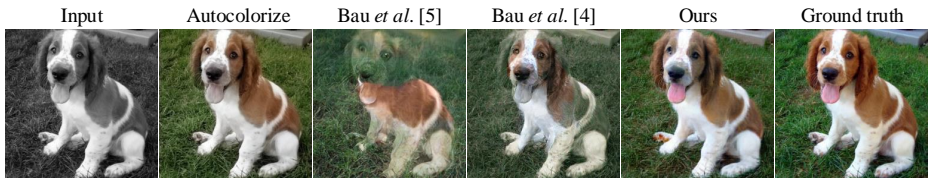
## 4 Applications

We first compare our method with other GAN inversion methods for reconstruction, and then show the application of DGP in a number of image restoration and image manipulation tasks. We adopt a BigGAN [7] to progressively reconstruct given images based on discriminator feature loss. The BigGAN is pre-trained on the ImageNet training set for conditional image synthesis. BigGAN is selected due to its excellent performance in image generation. Other GANs are possible. For dataset, we use the ImageNet [12] validation set that has not been observed by BigGAN. To quantitatively evaluate our method on image restoration tasks, we test on 1k images from the ImageNet validation set, where the first image for each class is collected to form the test set. We recommend readers to refer to the Appendix for implementation details and more qualitative results.

**Comparison with other GAN-inversion methods.** To begin with, we compare with other GAN-inversion methods [10,2,45,5] for image reconstruction. As shown in Table 1, our method achieves a very high PSNR and SSIM scores, outperforming other GAN-inversion methods by a large margin. It can be seen from Fig. 2 that conventional GAN-inversion methods like [45,5] suffer from obvious mismatches between reconstructed images and the target one, where the details or even contents are not well aligned. In contrast, the reconstruction error of DGP is almost visually imperceptible. More qualitative examples are provided in the Appendix. In the following sections we show that our method also well exploits the generative prior in various applications.

**Table 1.** Comparison with other GAN-inversion methods, including (a) optimizing latent vector [10,2], (b) learning an encoder [45], (c) a combination of (a)(b) [45], and (d) adding small perturbations to early stages based on (c) [5]. We reported PSNR, SSIM, and MSE of image reconstruction. The results are evaluated on the 1k ImageNet validation set

	(a)	(b)	(c)	(d)	Ours
PSNR $\uparrow$	15.97	11.39	16.46	22.49	<b>32.89</b>
SSIM $\uparrow$	46.84	32.08	47.78	73.17	<b>95.95</b>
MSE $\downarrow$ ( $\times e-3$ )	29.61	85.04	28.32	6.91	<b>1.26</b>



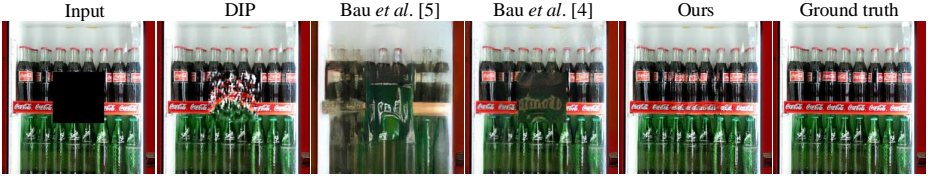
**Fig. 5. Colorization.** Qualitative comparison of Autocolorize [25], other GAN-inversion methods [5][4], and our DGP

#### 4.1 Image Restoration

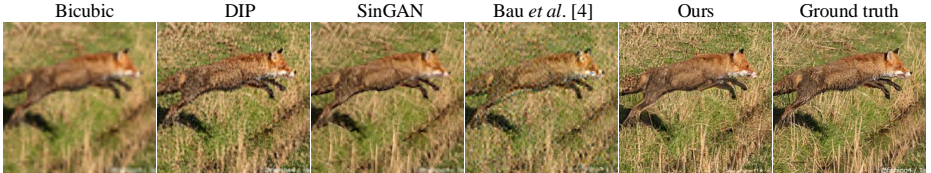
**Colorization.** Image colorization aims at restoring a gray-scale image  $\hat{\mathbf{x}} \in \mathbb{R}^{H \times W}$  to a colorful image with RGB channels  $\mathbf{x} \in \mathbb{R}^{3 \times H \times W}$ . To obtain  $\hat{\mathbf{x}}$  from the colorful image  $\mathbf{x}$ , the degradation transform  $\phi$  is a graying transform that only preserves the brightness of  $\mathbf{x}$ . By taking this degradation transform to Eq.(2), the goal becomes finding the colorful image  $\mathbf{x}^*$  whose gray-scale image is the same as  $\hat{\mathbf{x}}$ . We optimize Eq.(2) using back-propagation and the progressive discriminator based reconstruction technique in Section 3.2. Fig. 3(d) shows the reconstruction process. Note that the colorization task only requires to predict the “ab” dimensions of the Lab color space. Therefore, we transform  $\mathbf{x}^*$  to the Lab space, and adopt its “ab” dimensions as well as the given brightness dimension  $\hat{\mathbf{x}}$  to produce the final colorful image.

Fig. 5 presents the qualitative comparisons with the Autocolorize [25] method. Note that Autocolorize is directly optimized to predict color from gray-scale images while our method does not adopt such task-specific training. Despite so, our method is visually better or comparable to Autocolorize. To evaluate the colorization quality, we report the classification accuracy of a ResNet50 [19] model on the colorized images. The ResNet50 accuracy for Autocolorize [25], Bau *et al* [5], Bau *et al* [4], and ours are 51.5%, 56.2%, 56.0%, and 62.8% respectively, showing that DGP outperforms other baselines on this perceptual metric.

**Inpainting.** 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  $\mathbf{m}$ :  $\phi(\mathbf{x}) = \mathbf{x} \odot \mathbf{m}$ , where  $\odot$  is Hadamard’s product. As before, we put this degradation transform to Eq.(2), and reconstruct target images with missing boxes. Thanks to the generative image prior of the generator, the missing part tends to be recovered in harmony with the context, as illustrated in Fig. 6. In contrast, the absence of a learned image prior would



**Fig. 6. Inpainting.** Compared with DIP and [5][4], the proposed DGP could preserve the spatial coherence in image inpainting with large missing regions



**Fig. 7. Super-resolution ( $\times 4$ )** on  $64 \times 64$  size images. The comparisons of our method with DIP, SinGAN, and [4] are shown, where DGP produces sharper super-resolution results

**Table 2.** Inpainting evaluation. We reported PSNR and SSIM of the inpainted area. The results are evaluated on the 1k ImageNet validation set

	DIP	Zhu <i>et al</i> [45]	Bau <i>et al</i> [5]	Bau <i>et al</i> [4]	Ours
PSNR $\uparrow$	14.58	13.70	15.01	14.33	<b>16.97</b>
SSIM $\uparrow$	29.37	33.09	33.95	30.60	<b>45.89</b>

**Table 3.** Super-resolution ( $\times 4$ ) evaluation. We reported widely used NIQE [27], PSNR, and RMSE scores. The results are evaluated on the 1k ImageNet validation set. (MSE) and (D) indicate which kind of loss DGP is biased to use

	DIP	SinGAN	Bau <i>et al</i> [4]	Ours (MSE)	Ours (D)
NIQE $\downarrow$	6.03	6.28	5.05	5.30	<b>4.90</b>
PSNR $\uparrow$	23.02	20.80	19.89	<b>23.30</b>	22.00
RMSE $\downarrow$	17.84	19.78	25.42	<b>17.40</b>	20.09

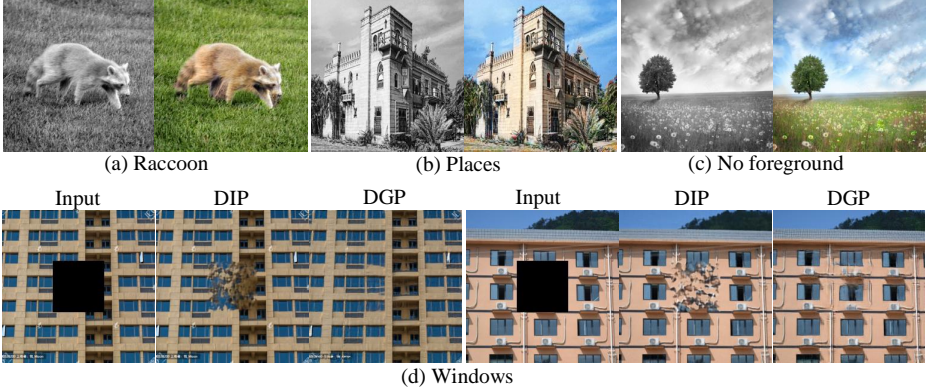
result in messy inpainting results, as in DIP. Quantitative results indicate that DGP outperforms DIP and other GAN-inversion methods by a large margin, as Table 2 shows.

**Super-Resolution.** In this task, one is given with a low-resolution image  $\hat{\mathbf{x}} \in \mathbb{R}^{3 \times H \times W}$ , and the purpose is to generate the corresponding high-resolution image  $\mathbf{x} \in \mathbb{R}^{3 \times fH \times fW}$ , where  $f$  is the upsampling factor. In this case, the degradation transform  $\phi$  is to downsample the input image by a factor  $f$ . Following DIP [36], we adopt the Lanczos downsampling operator in this work.

Fig. 7 and Table 3 show the comparison of DGP with DIP, SinGAN, and Bau *et al* [4]. Our method achieves sharper and more faithful super-resolution results than its counterparts. For quantitative results, we could trade off between perceptual quality like NIQE and commonly used PSNR score by using different combination ratios of discriminator loss and MSE loss at the final fine-tuning stage. For instance, when using higher MSE loss, DGP has excellent PSNR and RMSE performance, and outperforms other counterparts in all the metrics



**Fig. 8.** (a) Colorizing an image under different class conditions. (b) Simultaneously conduct colorization, inpainting, and super-resolution ( $\times 2$ )



**Fig. 9.** Evaluation of DGP on non-ImageNet images, including (a) ‘Raccoon’, a category not belonging to ImageNet categories, (b) image from Places dataset [44], (c) image without foreground object, and (d) windows. (a)(c)(d) are scratched from Internet

involved. And the perceptual quality NIQE could be further improved by biasing towards discriminator loss.

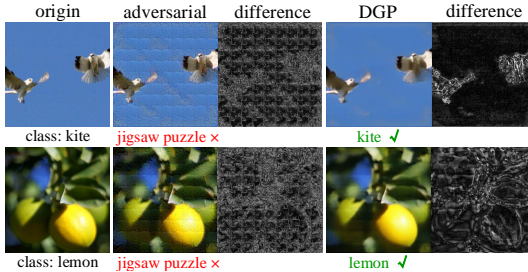
**Flexibility of DGP.** The generic paradigm of DGP provides more flexibility in restoration tasks. For example, an image of gray-scale bird may have many possibilities when restored in the color space. Since the BigGAN used in our method is a conditional GAN, we could achieve diversity in colorization by using different class conditions when restoring the image, as Fig. 8 (a) shows. Furthermore, our method allows hybrid restoration, *i.e.*, jointly conducting colorization, inpainting, and super-resolution. This could naturally be achieved by using a composite of degrade transform  $\phi(\mathbf{x}) = \phi_a(\phi_b(\phi_c(\mathbf{x})))$ , as shown in Fig. 8 (b).

**Generalization of DGP.** We also test our method on images not belonging to ImageNet. As Fig.9 shows, DGP restores the color and missed patches of these images reasonably well. Particularly, compared with DIP, DGP fills the missed patches to be well aligned with the context. This indicates that DGP does capture the ‘*spatial coherence*’ prior of natural images, instead of memorizing the ImageNet dataset. We scratch a small dataset with 18 images of windows, stones, and libraries to test our method, where DGP achieves 15.34 for PSNR and 41.53 for SSIM, while DIP has only 12.60 for PSNR and 21.12 for SSIM.

**Ablation Study.** To validate the effectiveness of the proposed discriminator guided progressive reconstruction method, we compare different fine-tuning

**Table 4.** Comparison of different loss type and fine-tuning strategy

Task	Metric	MSE	Perceptual	Discriminator	Discriminator +Progressive
Colorization	ResNet50 $\uparrow$	49.1	53.9	56.8	<b>62.8</b>
SR	NIQE $\downarrow$	6.54	6.27	6.06	<b>4.90</b>
	PSNR $\uparrow$	21.24	20.30	21.58	<b>22.00</b>

**Fig. 10. Adversarial defense.** DGP is capable of filtering out unnatural perturbations in the adversarial samples by reconstructing them**Table 5.** Adversarial defense evaluation. We reported the classification accuracy of a ResNet50. The results are evaluated on the 1k ImageNet validation set

method	clean image	adversarial	DefenceGAN	DIP	Ours
top1 acc. (%)	74.9	1.4	0.2	37.5	41.3
top5 acc. (%)	92.7	12.0	1.4	61.2	65.9

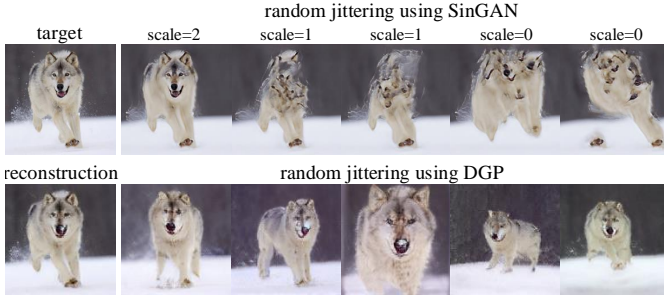
strategies in Table 4. There is a clear improvement of discriminator feature matching loss over MSE and perceptual loss, and the combination of the progressive reconstruction further boosts the performance. Fig. 2, Fig. 3, and Appendix provide qualitative comparisons. The results show that the progressive strategy effectively eliminates the ‘information lingering’ artifacts.

**Adversarial Defense.** Adversarial attack methods aim at fooling a CNN classifier by adding a certain perturbation  $\Delta\mathbf{x}$  to a target image  $\mathbf{x}$  [28]. In contrast, adversarial defense aims at preventing the model from being fooled by attackers. Specifically, the work of DefenceGAN [33] proposed to restore a perturbed image to a natural image by reconstructing it with a GAN. It works well for simple data like MNIST, but would fail for complex data like ImageNet due to poor reconstruction. Here we show the potential of DGP in adversarial defense under a black-box attack setting [3], where the attacker does not have access to the classifier and defender.

For adversarial attack, the degradation transform is  $\phi(\mathbf{x}) = \mathbf{x} + \Delta\mathbf{x}$ , where  $\Delta\mathbf{x}$  is the perturbation generated by the attacker. Since calculating  $\phi(\mathbf{x})$  is generally not differentiable, here we adopt DGP to directly reconstruct the adversarial image  $\hat{\mathbf{x}}$ . To prevent  $\mathbf{x}^*$  from overfitting to  $\hat{\mathbf{x}}$ , we stop the reconstruction when the MSE loss reaches  $5e-3$ . We adopt the adversarial transformation networks attacker [3] to produce the adversarial samples<sup>1</sup>.

<sup>1</sup> We use the code at <https://github.com/pfnet-research/nips17-adversarial-attack>





**Fig. 11.** Comparison of **random jittering** using SinGAN (above) and DGP (below)

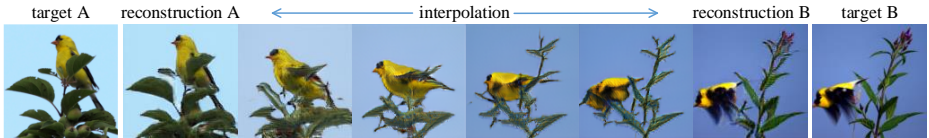
As Fig. 10 shows, the generated adversarial image contains unnatural perturbations, leading to misclassification for a ResNet50 [19]. After reconstructing the adversarial samples using DGP, the perturbations are largely alleviated, and the samples are thus correctly classified. The comparisons of our method with DefenseGAN and DIP are shown in Table 5. DefenseGAN yields poor defense performance due to inaccurate reconstruction. And DGP outperforms DIP, thanks to the learned image prior that produces more natural restored images.

## 4.2 Image Manipulation

Since DGP enables precise GAN reconstruction while preserving the generative property, it becomes straightforward to apply the fascinating capabilities of GAN to real images like random jittering, image morphing, and category transfer. In this section, we show the application of our method in these image manipulation tasks.

**Random Jittering.** We show the random jittering effects of DGP, and compare it with SinGAN. Specifically, after reconstructing a target image using DGP, we add Gaussian noise to the latent vector  $\mathbf{z}^*$  and see how the output changes. As shown in Fig. 11, the dog in the image changes in pose, action, and size, where each variant looks like a natural shift of the original image. For SinGAN, however, the jittering effects seem to preserve some texture, but losing the concept of ‘dog’. This is because it cannot learn a valid representation of dog by looking at only one dog. In contrast, in DGP the generator is fine-tuned in a moderate way such that the structure of image manifold captured by the generator is well preserved. Therefore, perturbing  $\mathbf{z}^*$  corresponds to shifting the image in the natural image manifold.

**Image Morphing.** The purpose of image morphing is to achieve a visually sound transition from one image to another. Given a GAN generator  $G$  and two latent vectors  $\mathbf{z}_A$  and  $\mathbf{z}_B$ , morphing between  $G(\mathbf{z}_A)$  and  $G(\mathbf{z}_B)$  could naturally be done by interpolating between  $\mathbf{z}_A$  and  $\mathbf{z}_B$ . In the case of DGP, however, reconstructing two target images  $\mathbf{x}_A$  and  $\mathbf{x}_B$  would result in two generators  $G_{\theta_A}$  and  $G_{\theta_B}$ , and the corresponding latent vectors  $\mathbf{z}_A$  and  $\mathbf{z}_B$ . Inspired by [38], to morph between  $\mathbf{x}_A$  and  $\mathbf{x}_B$ , we apply linear interpolation to both the latent



**Fig. 12. Image morphing.** Our method achieves visually realistic image morphing effects



**Fig. 13. Category transfer.** DGP enables the editing of semantics of objects in images

vectors and the generator parameters:  $\mathbf{z} = \lambda\mathbf{z}_A + (1 - \lambda)\mathbf{z}_B$ ,  $\boldsymbol{\theta} = \lambda\boldsymbol{\theta}_A + (1 - \lambda)\boldsymbol{\theta}_B$ ,  $\lambda \in (0, 1)$ , and generate images with the new  $\mathbf{z}$  and  $\boldsymbol{\theta}$ .

As Fig. 12 shows, our method enables highly photo-realistic image morphing effects. Despite the existence of complex backgrounds, the imagery contents shift in a natural way. To quantitatively evaluate image morphing quality, we apply image morphing to every consecutive image pairs for each class in the ImageNet validation set, and collect the intermediate images where  $\lambda = 0.5$ . For 50k images with 1k classes, this would create 49k generated images. We evaluate the image quality using Inception Score (IS) [32], and compare DGP with DIP, which adopts a similar network interpolation strategy. Finally, DGP achieves a satisfactory IS, 59.9, while DIP fails to create valid morphing results, leading to only 3.1 of IS.

**Category Transfer.** In conditional GAN, the class condition controls the content to be generated. So after reconstructing a given image via DGP, we can manipulate its content by tweaking the class condition. Fig. 1 (f) and Fig. 13 present examples of transferring the object category of given images. Our method can transfer the dog and bird to various other categories without changing the pose, size, and image configurations.

## 5 Conclusion

To summarise, we have shown that a GAN generator trained on massive natural images could be used as a generic image prior, namely deep generative prior (DGP). Embedded with rich knowledge on natural images, DGP could be used to restore the missing information of a degraded image by progressively reconstructing it under the discriminator metric. Meanwhile, such reconstruction strategy addresses the challenge of GAN-inversion, achieving multiple visually realistic image manipulation effects. Our results uncover the potential of a universal image prior captured by a GAN in image restoration and manipulation.

**Acknowledgment** We would like to thank Xintao Wang for helpful discussions.

## References

1. Abdal, R., Qin, Y., Wonka, P.: Image2stylegan: How to embed images into the stylegan latent space? In: ICCV. pp. 4432–4441 (2019)
2. Albright, M., McCloskey, S.: Source generator attribution via inversion. In: CVPR Workshops (2019)
3. Baluja, S., Fischer, I.: Adversarial transformation networks: Learning to generate adversarial examples. arXiv preprint arXiv:1703.09387 (2017)
4. Bau, D., Strobel, H., Peebles, W., Wulff, J., Zhou, B., Zhu, J.Y., Torralba, A.: Semantic photo manipulation with a generative image prior. *ACM Transactions on Graphics (TOG)* **38**(4), 59 (2019)
5. Bau, D., Zhu, J.Y., Wulff, J., Peebles, W., Strobel, H., Zhou, B., Torralba, A.: Seeing what a gan cannot generate. In: ICCV. pp. 4502–4511 (2019)
6. Bigdeli, S.A., Zwicker, M., Favaro, P., Jin, M.: Deep mean-shift priors for image restoration. In: NIPS. pp. 763–772 (2017)
7. Brock, A., Donahue, J., Simonyan, K.: Large scale gan training for high fidelity natural image synthesis. In: ICLR (2019)
8. Chen, Y., Pock, T.: Trainable nonlinear reaction diffusion: A flexible framework for fast and effective image restoration. *TPAMI* **39**(6), 1256–1272 (2016)
9. Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In: CVPR (2018)
10. Creswell, A., Bharath, A.A.: Inverting the generator of a generative adversarial network. In: *IEEE transactions on neural networks and learning systems* (2018)
11. Dai, B., Fidler, S., Urtasun, R., Lin, D.: Towards diverse and natural image descriptions via a conditional gan. In: ICCV. pp. 2970–2979 (2017)
12. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR. pp. 248–255 (2009)
13. Donahue, J., Krähenbühl, P., Darrell, T.: Adversarial feature learning. In: ICLR (2017)
14. Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. In: *TPAMI*. vol. 38, pp. 295–307. IEEE (2015)
15. Geman, S., Geman, D.: Stochastic relaxation, gibbs distributions, and the bayesian restoration of images. *TPAMI* (6), 721–741 (1984)
16. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: NIPS. pp. 2672–2680 (2014)
17. Gu, J., Shen, Y., Zhou, B.: Image processing using multi-code gan prior. CVPR (2020)
18. He, K., Sun, J., Tang, X.: Single image haze removal using dark channel prior. *TPAMI* **33**(12), 2341–2353 (2010)
19. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016)
20. Hussein, S.A., Tirer, T., Giryas, R.: Image-adaptive gan based reconstruction. arXiv preprint arXiv:1906.05284 (2019)
21. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML. pp. 448–456 (2015)
22. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: ECCV. pp. 694–711. Springer (2016)



23. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: CVPR. pp. 4401–4410 (2019)
24. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
25. Larsson, G., Maire, M., Shakhnarovich, G.: Learning representations for automatic colorization. In: ECCV. pp. 577–593. Springer (2016)
26. Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z., et al.: Photo-realistic single image super-resolution using a generative adversarial network. In: CVPR. pp. 4681–4690 (2017)
27. Mittal, A., Soundararajan, R., Bovik, A.C.: Making a completely blind image quality analyzer. *IEEE Signal processing letters* **20**(3), 209–212 (2012)
28. Nguyen, A., Yosinski, J., Clune, J.: Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In: CVPR. pp. 427–436 (2015)
29. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
30. Roth, S., Black, M.J.: Fields of experts: A framework for learning image priors. In: CVPR. pp. 860–867 (2005)
31. Rudin, L.I., Osher, S., Fatemi, E.: Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena* **60**(1-4), 259–268 (1992)
32. Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. In: NIPS. pp. 2234–2242 (2016)
33. Samangouei, P., Kabkab, M., Chellappa, R.: Defense-gan: Protecting classifiers against adversarial attacks using generative models. In: ICLR (2018)
34. Shaham, T.R., Dekel, T., Michaeli, T.: Singan: Learning a generative model from a single natural image. In: ICCV. pp. 4570–4580 (2019)
35. Shen, Y., Gu, J., Tang, X., Zhou, B.: Interpreting the latent space of gans for semantic face editing. CVPR (2020)
36. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Deep image prior. In: CVPR. pp. 9446–9454 (2018)
37. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. In: CVPR (2018)
38. Wang, X., Yu, K., Dong, C., Tang, X., Loy, C.C.: Deep network interpolation for continuous imagery effect transition. In: CVPR. pp. 1692–1701 (2019)
39. Xiangli\*, Y., Deng\*, Y., Dai\*, B., Loy, C.C., Lin, D.: Real or not real, that is the question. In: ICLR (2020)
40. Yang, C., Shen, Y., Zhou, B.: Semantic hierarchy emerges in deep generative representations for scene synthesis. arXiv preprint arXiv:1911.09267 (2019)
41. Yeh, R.A., Chen, C., Yian Lim, T., Schwing, A.G., Hasegawa-Johnson, M., Do, M.N.: Semantic image inpainting with deep generative models. In: CVPR. pp. 5485–5493 (2017)
42. Zhang, K., Zuo, W., Gu, S., Zhang, L.: Learning deep cnn denoiser prior for image restoration. In: CVPR. pp. 3929–3938 (2017)
43. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: ECCV. pp. 649–666. Springer (2016)
44. Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., Torralba, A.: Places: A 10 million image database for scene recognition. *TPAMI* **40**, 1452–1464 (2017)
45. Zhu, J.Y., Krähenbühl, P., Shechtman, E., Efros, A.A.: Generative visual manipulation on the natural image manifold. In: ECCV (2016)
46. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017)

47. Zhu, S.C., Mumford, D.: Prior learning and gibbs reaction-diffusion. TPAMI **19**(11), 1236–1250 (1997)

## Appendix

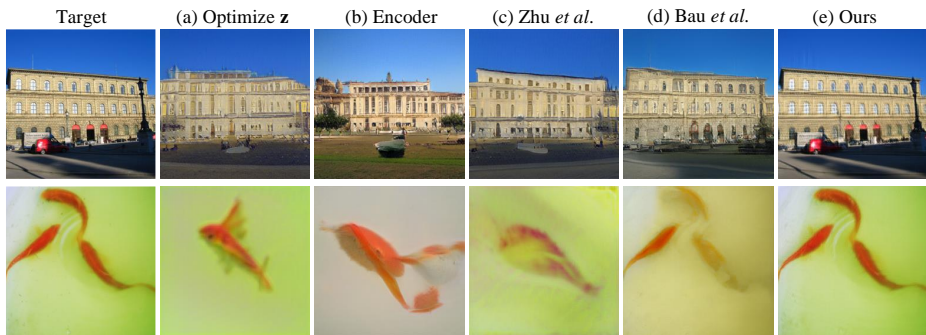
In this appendix, we provide more qualitative results and the implementation details in our experiments. Readers can see restoration and manipulation videos at our github repo.

### A Qualitative Examples

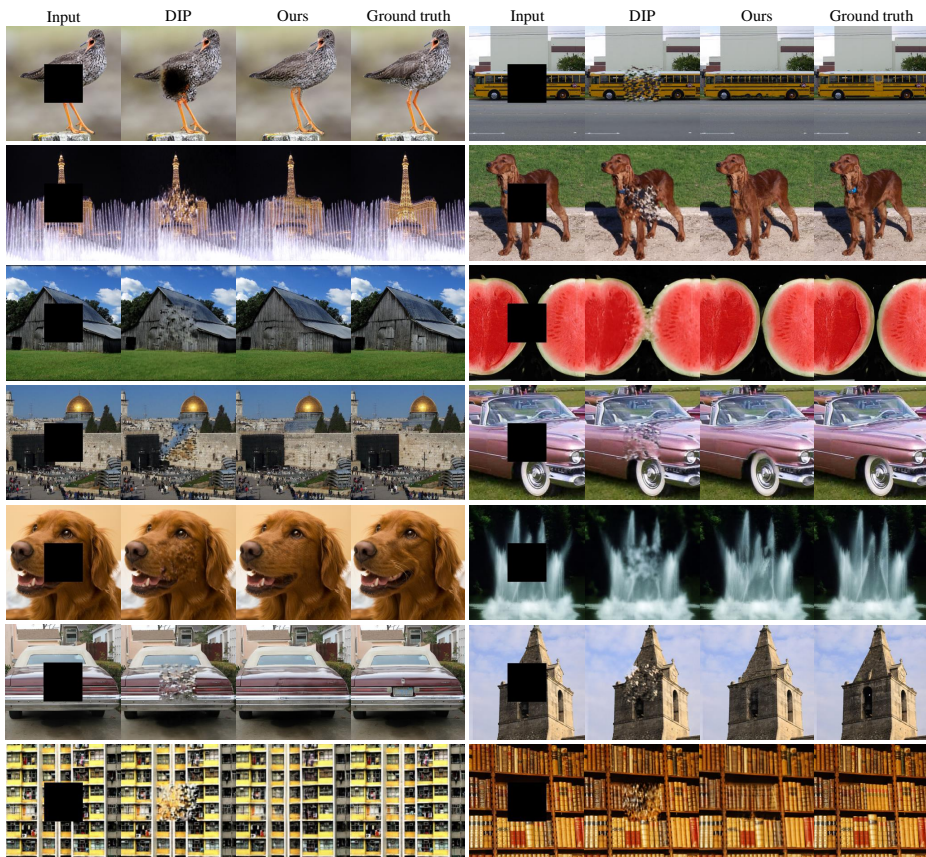
We extend the figures of the main paper with more examples, as shown from Fig. 14 to Fig. 24.



**Fig. 14. Colorization.** This is an extension of Fig.5 in the main paper.

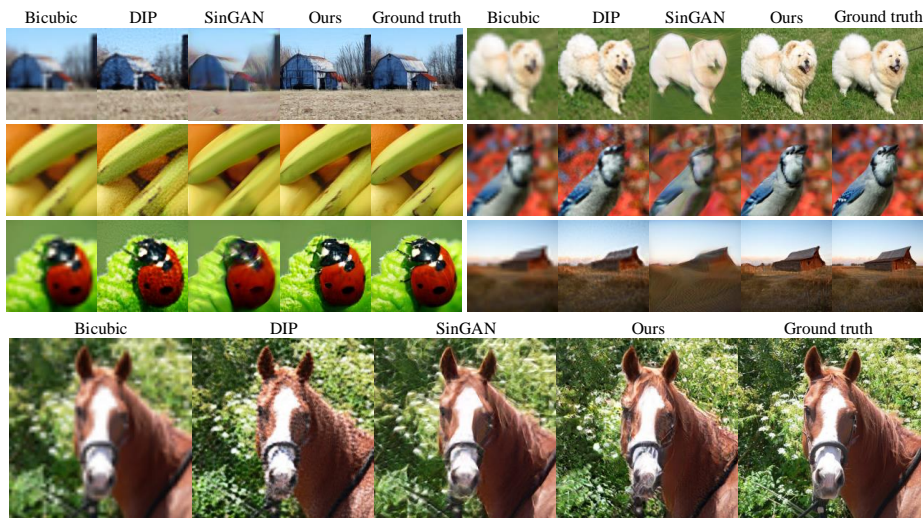


**Fig. 15. Image reconstruction.** We compare our method with other GAN-inversion methods including (a) optimizing latent vector [10,2], (b) learning an encoder [45], (c) a combination of (a)(b) [45], and (d) adding small perturbations to early stages based on (c) [5].



**Fig. 16. Inpainting.** This is an extension of Fig.6 in the main paper. The proposed DGP tends to recover the missing part in harmony with the context. Images of the last row are scratched from the Internet.

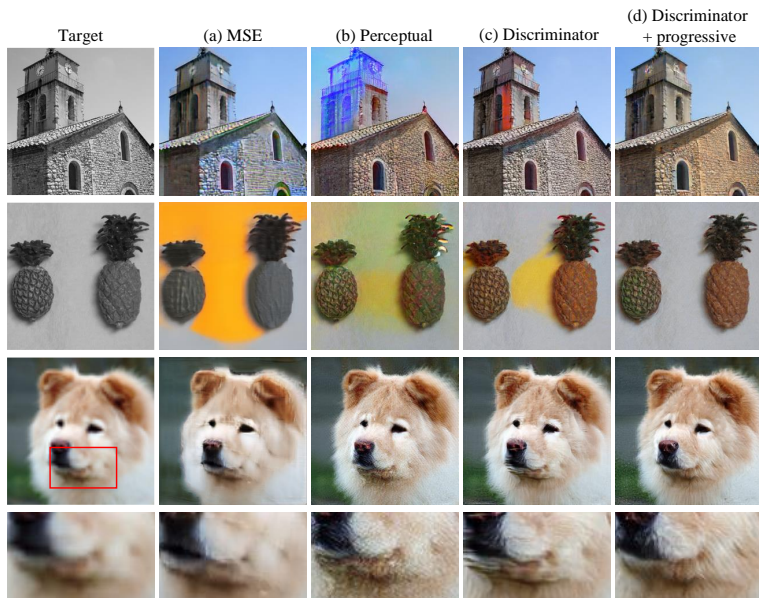




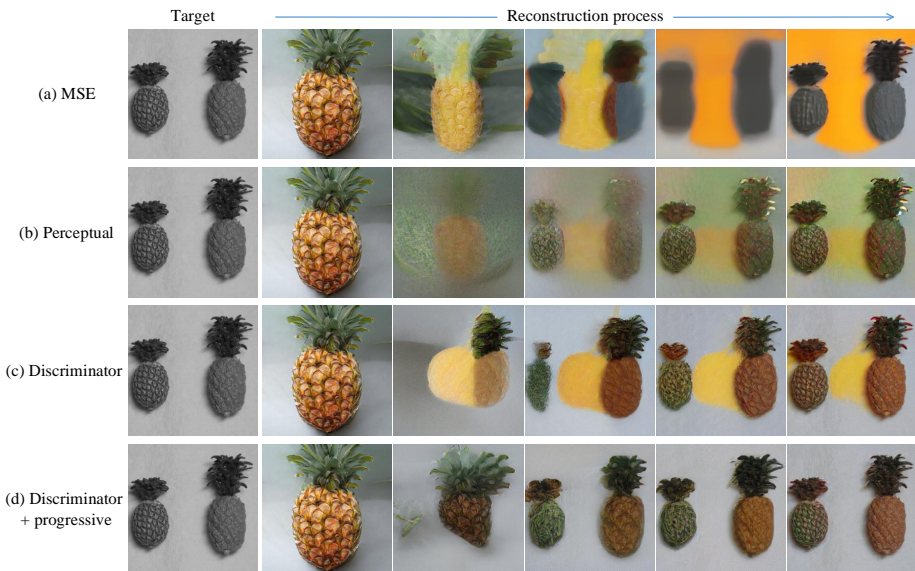
**Fig. 17. Super-resolution ( $\times 4$ ) on  $32 \times 32$  (above) and  $64 \times 64$  (below) size images. This is an extension of Fig.7 in the main paper.**



**Fig. 18. The reconstruction process of DGP in various image restoration tasks.**



**Fig. 19.** Comparison of different loss types and optimization techniques in colorization and super-resolution, including (a) MSE loss, (b) perceptual loss with VGG network [22], (c) discriminator feature matching loss, and (d) combined with progressive reconstruction.



**Fig. 20.** Comparison of different loss types and optimization techniques when fine-tuning the generator to restore the image.



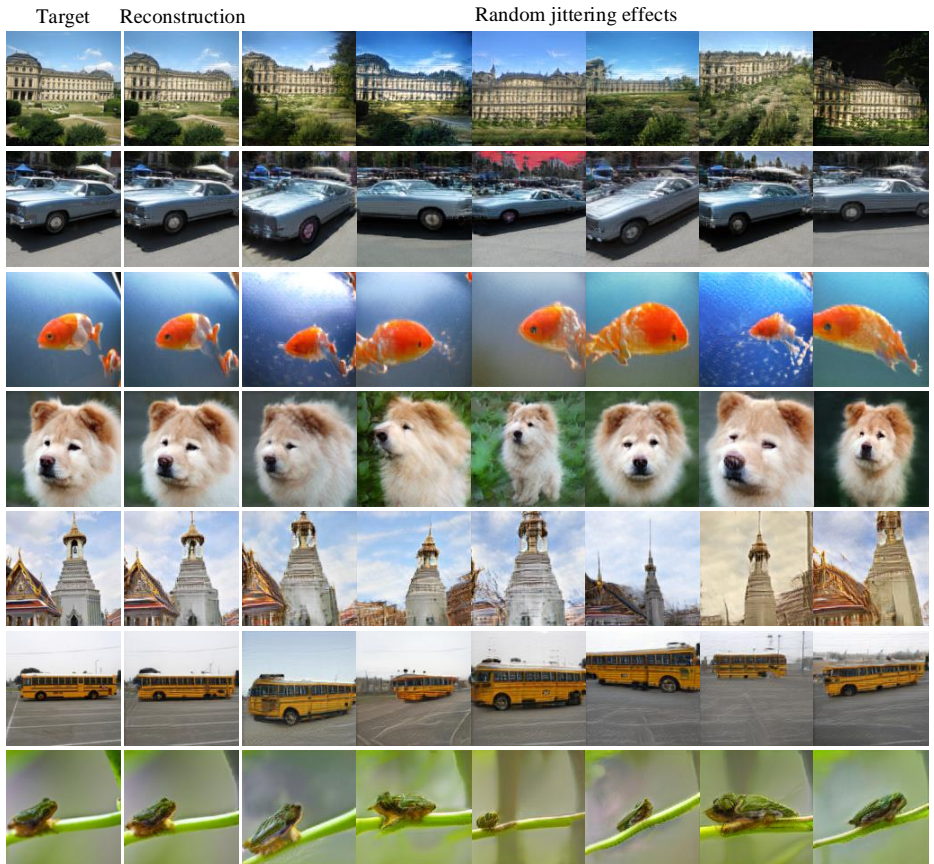


Fig. 21. Random jittering. This is an extension of Fig.11 in the main paper.

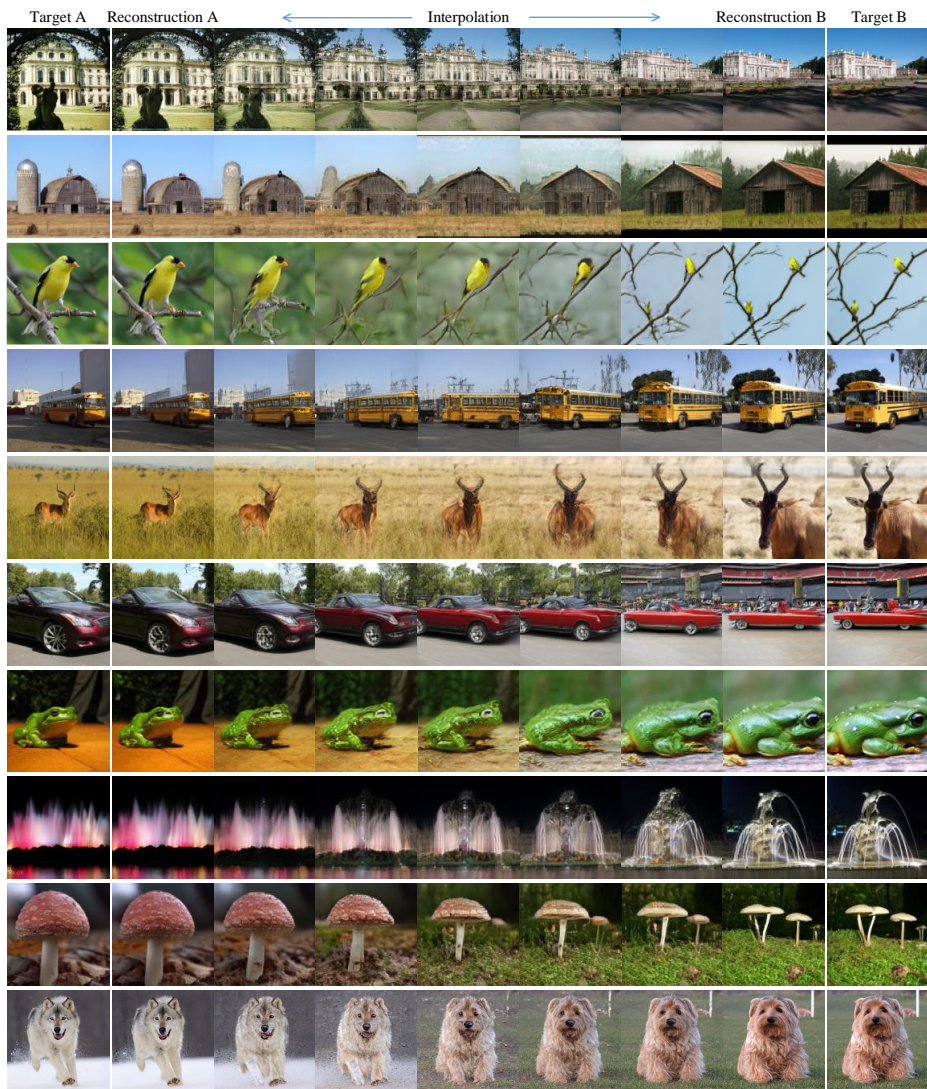
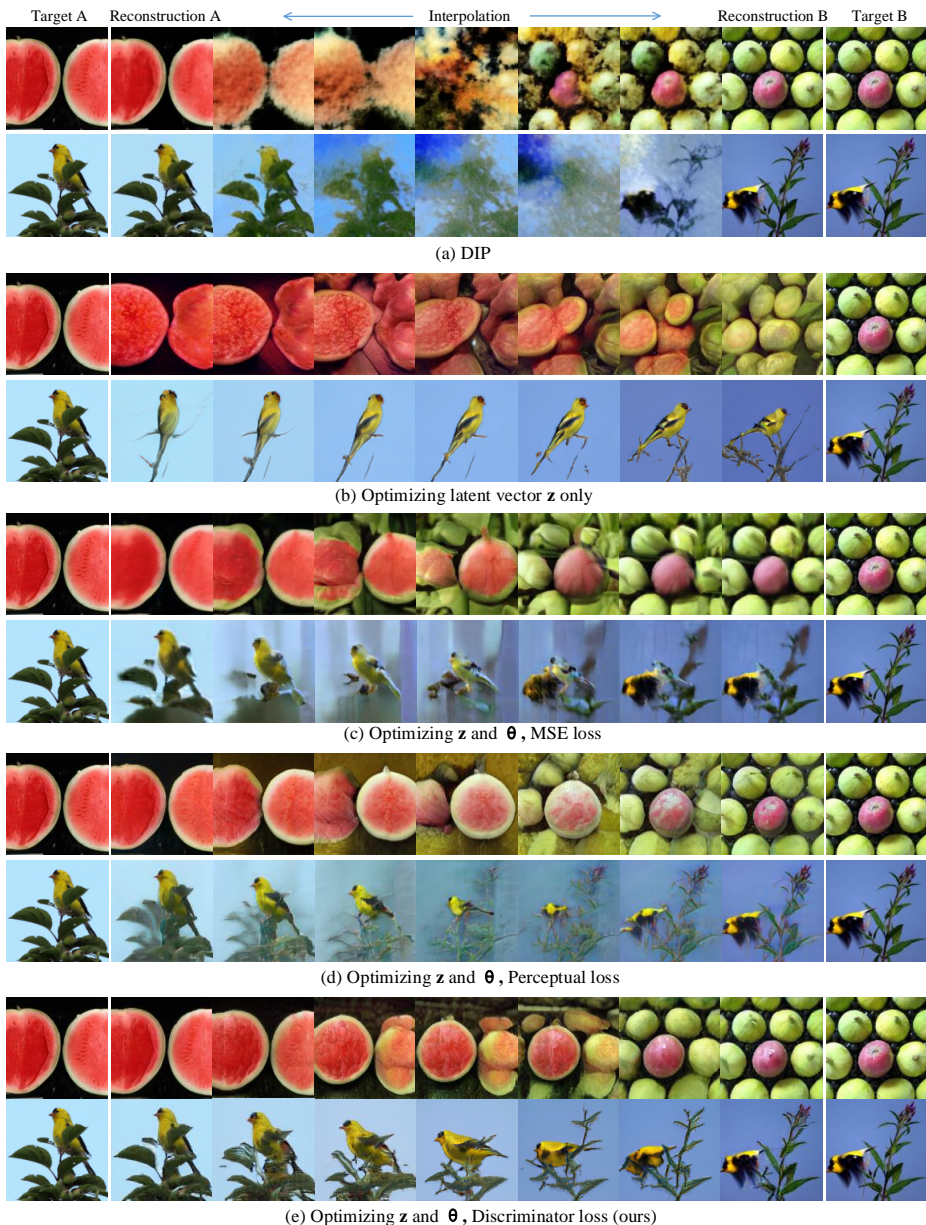


Fig. 22. Image morphing. This is an extension of Fig.12 in the main paper.





**Fig. 23.** Comparison of various methods in image morphing, including (a) using DIP, (b) optimizing the latent vector  $\mathbf{z}$  of the pre-trained GAN, and (c)(d)(e) optimizing both  $\mathbf{z}$  and the generator parameter  $\theta$  with (c) MSE loss, (d) perceptual loss with VGG network [22], and (e) discriminator feature matching loss. (b) fails to produce accurate reconstruction while (a)(c)(d) could not obtain realistic interpolation results. In contrast, our results in (e) are much better.





**Fig. 24. Category transfer.** The red box shows the target, and the blue box shows the reconstruction. Others are category transfer results.

## B Implementation Details

**Architectures.** We adopt the BigGAN[7] architectures of  $128^2$  and  $256^2$  resolutions in our experiments. For the  $128^2$  resolution, we use the best setting of [7], which has a channel multiplier of 96 and a batchsize of 2048. As for the  $256^2$  resolution, the channel multiplier and batchsize are respectively set to 64 and 1920 due to limited GPU resources. We train the GANs on the ImageNet training set, and the  $128^2$  and  $256^2$  versions have Inception scores of 103.5 and 94.5 respectively. Our experiments are conducted based on PyTorch [29].

**Initialization.** In order to ease the optimization goal of Eq.4 in the paper, it is a good practice to start with a latent vector  $\mathbf{z}$  that produces an approximate reconstruction. Therefore, we randomly sample 500 images using the GAN, and select the nearest neighbor of the target image under the discriminator feature metric as the starting point. Since encoder based methods tend to fail for degraded input images, they are not used in this work.

Note that in BigGAN, a class condition is needed as input. Therefore, in order to reconstruct an image, its class condition is required. This image classification problem could be solved by training a corresponding deep network classifier and is not the focus of this work, hence we assume the class label is given except for the adversarial defense task. For adversarial defense and images whose classes are not given, both the latent vector  $\mathbf{z}$  and the class condition are randomly sampled.

**Fine-tuning.** With the above pre-trained BigGAN and initialized latent vector  $\mathbf{z}$ , we fine-tune both the generator and the latent vector to reconstruct a target image. As the batchsize is only 1 during fine-tuning, we use the tracked global statistics (*i.e.*, running mean and running variance) for the batch normalization (BN) [21] layers to prevent inaccurate statistic estimation. The discriminator of BigGAN is composed of a number of residual blocks (6 blocks and 7 blocks for  $128^2$  and  $256^2$  resolution versions respectively). The output features of these blocks are used as the discriminator loss, as described in Eq.(6) of the paper. In order to prevent the latent vector from deviating too much from the prior gaussian distribution, we add an additional L2 loss to the latent vector  $\mathbf{z}$  with a loss weight of 0.02. We adopt the ADAM optimizer [24] in all our experiments. The detailed training settings for various tasks are listed from Table.6 to Table.11, where the parameters in these tables are explained below:

*Blocks num.:* the number of generator blocks to be fine-tuned. For example, for blocks num.=1, only the shallowest block is fine-tuned.

*D loss weight:* the factor multiplied to the discriminator loss.

*MSE loss weight:* the factor multiplied to the MSE loss.

*Iterations:* number of training iterations of each stage.

*G lr:* the learning rate of the generator blocks.

*z lr:* the learning rate of the latent vector  $\mathbf{z}$ .

For inpainting and super-resolution, we use a weighted combination of discriminator loss and MSE loss, as the MSE loss is beneficial for the PSNR metric. We

**Table 6.** The fine-tuning setting of colorization. The explanation of these parameters are in the main text

Stage	1	2	3	4	5
Blocks num.	1	2	3	4	5
D loss weight	1	1	1	1	1
MSE loss weight	0	0	0	0	0
Iterations	200	200	300	400	300
G lr	5e-5	5e-5	5e-5	5e-5	2e-5
z lr	2e-3	1e-3	5e-4	5e-5	2e-5

**Table 8.** The fine-tuning setting of super-resolution. This setting is biased towards MSE loss

Stage	1	2	3	4	5
Blocks num.	1	2	3	4	5
D loss weight	1	1	1	0.5	0.1
MSE loss weight	1	1	1	50	100
Iterations	200	200	200	200	200
G lr	2e-4	2e-4	1e-4	1e-4	1e-5
z lr	1e-3	1e-3	1e-4	1e-4	1e-5

**Table 10.** The fine-tuning setting of adversarial defense. The fine-tuning is stopped if the MSE loss reaches 5e-3

	stage 1	stage 2
Blocks num.	6	6
D loss weight	0	0
MSE loss weight	1	1
Iterations	100	900
G lr	2e-7	1e-4
z lr	5e-2	1e-4

**Table 7.** The fine-tuning setting of inpainting. In this task we also fine-tune the class embedding apart from the generator blocks

Stage	1	2	3	4
Blocks num.	5	5	5	5
D loss weight	1	1	0.1	0.1
MSE loss weight	1	1	100	100
Iterations	400	200	200	200
G lr	2e-4	1e-4	1e-4	1e-5
z lr	1e-3	1e-4	1e-4	1e-5

**Table 9.** The fine-tuning setting of super-resolution. This setting is biased towards discriminator loss

Stage	1	2	3	4	5
Blocks num.	1	2	3	4	5
D loss weight	1	1	1	1	1
MSE loss weight	1	1	1	1	1
Iterations	200	200	200	200	200
G lr	5e-5	5e-5	2e-5	1e-5	1e-5
z lr	2e-3	1e-3	2e-5	1e-5	1e-5

**Table 11.** The fine-tuning setting of manipulation tasks including random jittering, image morphing, and category transfer

	stage 1	stage 2	stage 3
Blocks num.	5	5	5
D loss weight	1	1	1
MSE loss weight	0	0	0
Iterations	125	125	100
G lr	2e-7	2e-5	2e-6
z lr	1e-1	2e-3	2e-6

also seamlessly replace BN with instance normalization (IN) for the setting in Table. 7, Table. 8, and Table. 10, which enables higher learning rate and leads to better PSNR. This is achieved by initialize the scale and shift parameters of IN with the statistics of the output features of BN. Our quantitative results on adversarial defense is based on the 256<sup>2</sup> resolution model, while those for other tasks are based on the 128<sup>2</sup> resolution models.