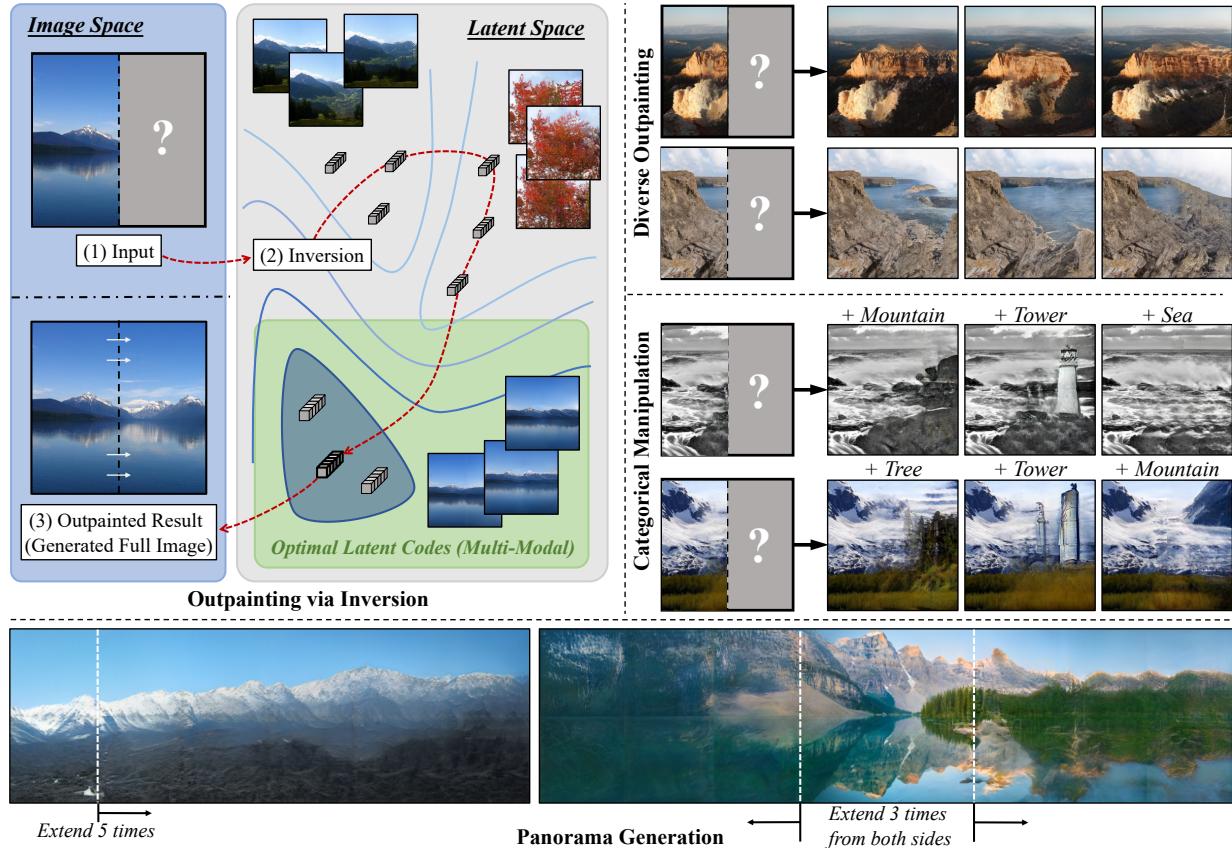


000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  

# Diverse Image Outpainting via GAN Inversion

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
Anonymous ICCV submission092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107  
Paper ID 2914

038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
Figure 1: (top-left) Given an input image and a trained generator, the proposed algorithm searches for latent codes that can generate images containing the input image. We can naturally achieve (top-right) diverse image outpainting, (middle-right) categorical manipulation for outpainting area, and (bottom) generate panorama with rich and complex structure.

## Abstract

Image outpainting seeks for a semantically consistent extension of the input image beyond its available content. Compared to inpainting — filling in missing pixels in a way coherent with the neighboring pixels — outpainting can be achieved in more diverse ways since the problem is less constrained by the surrounding pixels. Existing image outpainting methods pose the problem as a conditional image-to-image translation task, often generating repetitive structures and textures by replicating the content available in the input image. In this work, we formulate the problem from the perspective of inverting generative adversar-

ial networks. Our generator renders micro-patches conditioned on their joint latent code as well as their individual positions in the image. To outpaint an image, we seek for multiple latent codes not only recovering available patches but also synthesizing diverse outpainting by patch-based generation. This leads to richer structure and content in the outpainted regions. Furthermore, our formulation allows for outpainting conditioned on the categorical input, thereby enabling flexible user controls. Extensive experimental results demonstrate the proposed method performs favorably against existing in- and outpainting methods, featuring higher visual quality and diversity.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

## 1. Introduction

Given an input image, we can easily picture how adjacent images might look, had they been captured. For example, given an image of mountains, we can picture the surroundings covered by forests or snow, imagine a lake beneath the hillside, and visualize cliffs near the ocean. This mental skill depends on our prior experience and exposure to diverse scenery. In other words, this is an *image outpainting* task. It can enable various content creation applications such as image editing using extrapolated regions, panorama image generation, and extended immersive experience in virtual reality, to name a few.

Recent advances in image inpainting [22, 26, 38, 39] do not directly address the outpainting problem as the former has more context to deal with — the missing pixels have a larger amount of available surrounding pixels, serving as the boundary conditions and providing crucial guidance for inpainting. In contrast, the outpainting problem can rely only on the context of the available image, with only a scarce number of pixels near the boundary available as the boundary condition. Furthermore, the texture and semantics of the outpainted regions should be coherent with that of the input. Finally, outpainting methods ought to support diversity in the generated content. A similar analogy is between video interpolation and video prediction, where the former deals with existing events [14] while the latter tries to model multiple futures [35].

In the literature, image outpainting is addressed from the image-to-image translation (I2I) perspective [30, 37]. These methods aim to learn a deterministic mapping from the domain of partial images to the domain of complete outpainted images. This formulation is limited in several respects. First, for the I2I methods, the available pixels serve as a strong source of context, thereby facilitating leakage of textures and structures of the input to the output and leading to the repetitive nature of the outpainting (as shown in panorama results in [30]). Second, existing I2I-based methods are deterministic [30, 37], while in reality there exist numerous ways each image can be outpainted. Applying the available multimodal I2I methods [13, 18] to the outpainting problem is non-trivial.

In this work, we tackle the outpainting problem by inverting generative adversarial networks (GANs) [1, 4, 7, 23, 43]. Similar to Lin et al. [20], we extend a StyleGAN2-based [16] generator to perform generation in a coordinate conditional manner and independently generate spatially consistent micro-patches. Each micro-patch shares the global latent code with the rest of micro-patches in the image, while having a unique coordinate label. Outpainting can then be formulated as finding the optimal latent codes for the available input micro-patches, followed by generating the desired regions by providing the proper coordinate conditioning. To search for the latent code, we propose a

GAN inversion process that finds multiple latent codes producing diverse outpainted regions, unlocking diversity in the output. In addition, we propose a categorical generation schema to enable flexible user control. Figure 1 shows examples of multi-modal and categorical outpainting.

We qualitatively and quantitatively evaluate the proposed method on the Place365 [42] dataset, and the Flickr-Scenery dataset which we collected. We leverage Fréchet Inception Distance (FID) [12] and conduct a user study to evaluate the realism of outpainted images. Since the proposed method can achieve multi-modal generation, we measure the diversity using the Learned Perceptual Image Patch Similarity (LPIPS) metric [41]. Finally, we demonstrate the scenario of categorical generation in the outpainting area and the panorama generation.

## 2. Related Work

**Generative Adversarial Networks.** Generative models aim to model and sample from a target distribution. Generative adversarial networks [11], among various generative models, have demonstrated superior performance in generating high-quality samples. The core idea of GANs is a two-player game between a generator aiming to map noise vectors to realistic images and a discriminator attempting to discriminate the generated images from the real ones. GANs facilitate a variety of creation tasks such as image-to-image translation [45], text-to-image generation [31, 40], and video generation [19, 33]. However, most of the models generate new images from scratch given various conditional contexts, and generally lack the ability to perform editing and interactive manipulation on existing images.

**GAN Inversion.** To fully exploit the ability and explore the interpretability of well-trained GANs, GAN inversion has been proposed to find the latent codes that can accurately recover given images for a trained GAN model. There are two main branches of approaches. *Encoder-based methods* [5, 9, 27] adopt an additional encoder to learn the mapping from the image domain to the latent space. *Optimization-based methods* [1, 2, 7, 21, 23] use gradient-based optimization methods (i.e., stochastic gradient decent and ADAM) with reconstruction loss as the objective function to find latent codes that can recover input images. Other variants use encoders to get an initialization for the optimization process [4, 44], or modify the training framework by incorporating $\odot$  invertibility [8, 43]. In this work, we adopt the optimization-based technique to tackle the image outpainting task.

**Image Inpainting.** From the aspect of filling missing pixels in images with generative models, the inpainting problem is conceptually related to the outpainting task. Existing image inpainting methods can be categorized into two groups. The first line of work leverages patch similarity and diffusion to

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

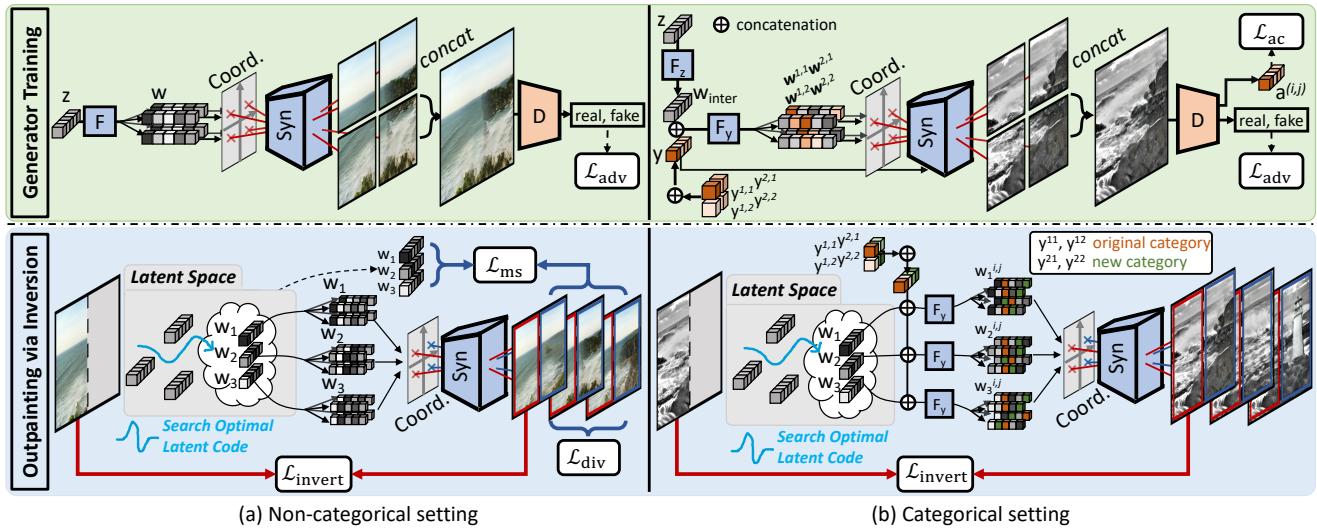


Figure 2: **Overview. Generator training:** The proposed generator adopts StyleGAN2 as the backbone and incorporates coordinate conditioning, where micro-patches are generated conditioned on their positions in the image. **Outpainting via inversion:** We search for latent codes that can recover the given partial images and synthesize diverse samples in the outpainted regions. In addition to the *unconditional setting*, we introduce a *categorical setting* variant, which enables flexible user controls for the categorical manipulation on each outpainting micro-patch.

obtain essential information from known regions [10, 37]. These methods usually work well on textures but fail to learn semantic structures of images. The other line of work adopts a learning-based approach to gain better semantic understanding [22, 26, 38, 39]. Most methods apply an encoder-decoder model with the reconstruction loss and adversarial loss to ensure the filled content is smooth and realistic. In this work, we focus on image outpainting instead of inpainting. Image outpainting is more challenging since it entails *creating* new contents rather than *filling in* partial regions, requiring a substantial understanding of scenes.

**Image Outpainting.** Most image outpainting methods [10, 17, 29, 32] apply patch-based retrieval and matching algorithms to predict possible extrapolation. Recently, several approaches [30, 34, 37] apply GAN models and formulate the problem as an image-to-image translation task. However, the conditional formulation relies heavily on the given available pixels and tends to create repetitive textures and structures. To the best of our knowledge, the proposed method is the first attempt to tackle the image outpainting task from the GAN inversion perspective.

### 3. Diverse Outpainting via Inversion

**Overview.** The goal of image outpainting is to outward-synthesize unknown regions with respect to the given input image. Our pipeline consists of two stages, **generator training** and **outpainting via inversion**. In Section 3.1, we first introduce a generator based on the StyleGAN [15, 16]

and COCO-GAN [20]. It is trained to output micro-patches conditioned on the joint latent and on the coordinate of the patch in the output image. We do not specifically optimize the generator to perform outpainting. During the outpainting stage (Section 3.2) we find the optimal latent code for the available input patches in the latent space of the trained patch-based generator. Outpainting is then performed by combining the desired coordinates and the found optimal latent code and generating a new patch. We further propose a categorical-conditioning scheme to enable controllable outpainting. Finally, a simple blending algorithm to further mitigate the artifacts is introduced in Section 3.3.

#### 3.1. Coordinate-conditioned Generator

In this work, we handle two different settings: (a) *non-categorical generation* that synthesizes images from latent codes, and (b) *categorical generation* that uses categorical labels as additional conditional context, enabling more user-control in the following inversion stage.

**Non-categorical generation.** We use the StyleGAN2 [16] as our backbone architecture. Given a latent  $\mathbf{z}$  from the input latent space  $\mathcal{Z}$ , we obtain an intermediate code  $\mathbf{w} \in \mathcal{W}$  by a non-linear mapping network  $F$ . Similar to in [36], we map  $\mathbf{w}$  to a Gaussianized space  $\mathcal{V}$ . The mapping is achieved via a Leaky ReLU (LRU) with a negative slope of 5, that is,  $\mathbf{v} = \text{LRU}_{5.0}(\mathbf{w})$ . The outpainting quality in the later GAN inversion stage can be substantially improved with the additional Gaussianized space. The necessity of adopting the Gaussianized space is discussed in Section 3.2.

We formulate the image outpainting problem as finding the latent codes that synthesizes images overlapping with the input image. In the inversion process, we seek for a latent code for the whole image while having only a part of the image available. Therefore, instead of generating a full image, the generator synthesizes several *micro-patches*  $\{I_{\text{micro}}^{i,j}\}_{i,j=1,\dots,n}$ , which will be concatenated to form a full image  $I_f$ . Each patch depends on the joint latent code and its coordinates. For an  $n \times n$  micro-patches generation setting, the corresponding coordinates to  $\{I_{\text{micro}}^{i,j}\}_{i,j=1,\dots,n}$  are  $\{c^{i,j}\}_{i,j=1,\dots,n}$ . We set  $c^{1,1} = (-1, -1)$ ,  $c^{n,n} = (1, 1)$ , and the rest, if any, are obtained by linear interpolation. The output image  $I_f$  is generated as follows:

$$\begin{aligned} \mathbf{w} &= F(\mathbf{z}), \\ \mathbf{v} &= \text{LRU}_{5.0}(\mathbf{w}), \\ I_{\text{micro}}^{i,j} &= G(\mathbf{v}, c^{i,j}), \\ I_f &= \text{concat}_{i,j=1,\dots,n}(I_{\text{micro}}^{i,j}). \end{aligned} \quad (1)$$

We train the generator using the Wasserstein-GAN loss [3] with real full-images  $I_r$  and generated full-images  $I_f$ :

$$\mathcal{L}_{\text{adv}} = \mathbb{E}_{I_r} [D(I_r)] - \mathbb{E}_z [D(I_f)]. \quad (2)$$

**Categorical generation.** To enable fine-grained user control in the inversion stage, we propose a categorical generation schema. Given a real image  $I_r$ , we divide it into micro-patches  $\{I_{\text{micro}}^{i,j}\}$ . We then obtain the categorical labels  $\{y^{i,j}\}$  for each micro-patch using the off-the-shelf DeepLabV3 [6] model, and set the  $k$ -th element in the multi-class binary label vector  $y^{i,j}$  to 1 if any of the pixels within  $I_{\text{micro}}^{i,j}$  is recognized as the  $k$ -th class.

To use categorical information as a conditional input, we split the nonlinear mapping network  $F$  into  $\{F_z, F_y\}$ . Here,  $F_z$  operates the same way as  $F$  in the non-categorical setting, while  $F_y$  takes  $\{y^{i,j}\}$  as an additional input and fuses the information with the output of  $F_z$ . The new  $w_{i,j}$  under the categorical setting is computed by

$$\begin{aligned} \mathbf{w}_{\text{inter}} &= F_z(\mathbf{z}), \\ \mathbf{w}^{i,j} &= F_y(\mathbf{w}_{\text{inter}}, y^{i,j}). \end{aligned} \quad (3)$$

Next, similar to the non-categorical generation setting, we first Gaussianize the code with  $\mathbf{v}^{i,j} = \text{LRU}_{5.0}(\mathbf{w}^{i,j})$ , generate micro-patches  $I_{\text{micro}}^{i,j} = G(\mathbf{v}, c_{i,j})$ , then concatenate  $\{I_{\text{micro}}^{i,j}\}_{i,j=1,\dots,n}$  into a full image  $I_f$ . We apply an auxiliary classifier [25] that uses the last intermediate features of  $D$  to perform multi-class classification  $a^{i,j}$  for all  $I_{\text{micro}}^{i,j}$ , which aims at learning a proper conditional distribution of  $I_{\text{micro}}^{i,j}$  regarding the  $y^{i,j}$  input to  $G$ .

$$\begin{aligned} \mathcal{L}_{\text{adv}} &= \mathbb{E}_{I_r} [D(I_r)] - \mathbb{E}_z [D(I_f)], \\ \mathcal{L}_{\text{cls}} &= \text{BCE}(a^{i,j}, y^{i,j}), \end{aligned} \quad (4)$$

where BCE is the binary cross entropy loss function. The full training objective is:

$$\min_D \max_G \mathcal{L}_{\text{adv}} + \min_{G,D} \mathcal{L}_{\text{cls}}.$$

### 3.2. GAN Inversion with Diversity Loss

Given a trained coordinate-conditioned generator  $G$  as discussed in the previous section and an input image  $R$  as a reference, we generate a set of possible outpainted images by composing  $R$  with generated micro-patches  $\{O^m\}$ . For brevity and notation clarity, here we assume  $G$  is trained with a grid of  $2 \times 2$  micro-patches,  $\{R_{\text{micro}}^{1,1}, R_{\text{micro}}^{1,2}, R_{\text{micro}}^{2,1}, R_{\text{micro}}^{2,2}\}$ . Furthermore, for presentation simplicity, we assume  $R$  to be on the left and consists of two left-side micro-patches (i.e.,  $R_{\text{micro}}^{1,1}$  and  $R_{\text{micro}}^{1,2}$ ), while the outpainted area  $\{O^m\}$  on the right, as shown in the lower-half of Figure 2. Note that, in practice,  $G$  is not restricted to  $2 \times 2$ ,  $R$  can be of any resolution, and outpainting can be performed using an arbitrary direction.

Similar to existing optimization-based GAN inversion methods, we seek for the optimal latent code  $w$  that recovers the input image. The basic loss function is:

$$\begin{aligned} R_f &= \text{concat}(G(\mathbf{v}, c^{1,1}), G(\mathbf{v}, c^{1,2})), \\ \mathcal{L}_{\text{mse}} &= \|R - R_f\|_2, \\ \mathcal{L}_{\text{percept}} &= \text{Percept}(R, R_f), \end{aligned} \quad (5)$$

where  $\mathbf{v} = \text{LRU}_{5.0}(\mathbf{w})$  and Percept is the perceptual distance proposed in [41].

The outpainting process requires not only the reconstructed parts to be correct  $R$  (i.e.,  $I_{\text{micro}}^{1,1}$  and  $I_{\text{micro}}^{1,2}$ ), but the outpainted parts (i.e.,  $I_{\text{micro}}^{2,1}$  and  $I_{\text{micro}}^{2,2}$ ) to be realistic and consistent. Note that the continuity and consistency between micro-patches are enforced by joint latent and the coordinate conditioning schema. During the generator training, the latent is sampled from a Gaussian distribution. Therefore, it is crucial to encourage the sought latent code  $w$  to belong to the domain of the training data, and be interpretable by  $G$ , instead of overfitting to the given image with the out-of-domain latent code. As the first step, in Section 3.1, we add an additional *Gaussianized* space [36]  $\mathcal{V}$  after  $\mathcal{W}$ , simplifying the complex and arbitrarily shaped  $\mathcal{W}$  with LRU<sub>5.0</sub>. Next, with the Gaussianized  $\mathcal{V}$ , we can easily derive the mean  $\mu$  and covariance matrix  $\Sigma$  of the distribution  $p(\mathbf{v})$ , where  $\mathbf{v} \in \mathcal{V}$ . We encourage recovered  $\mathbf{v}$  to be in the training distribution by regularizing its prior:

$$\mathcal{L}_{\text{prior}} = (\mathbf{v} - \mu)^\top \Sigma^{-1} (\mathbf{v} - \mu). \quad (6)$$

To enable diverse outpainting, we apply two different objective functions. Assuming we target at generating  $m$  different outpainted results, we first explicitly penalize the inverted latent codes with their pairwise distance:

$$\mathcal{L}_{\text{div}} = - \sum_{i=1}^m \sum_{j=i+1}^m \|\mathbf{w}^i - \mathbf{w}^j\|_1. \quad (7)$$

Then, to further encourage the model to seek for different final latent codes within the latent space, we apply a mode-seeking regularization [24]:

Table 1: **Quantitative comparisons.** The proposed method outperforms related state-of-the-art baselines on both Places365 and Flickr-Scenery datasets in FID and IS metrics, measuring both the visual quality and diversity.

Method	Place365		Flickr-Scenery	
	FID ↓	IS ↑	FID ↓	IS ↑
Boundless [30]	35.02	6.15	61.98	6.98
NS-outpaint [37]	50.68	4.70	61.16	4.76
DeepFillv2 [38, 39]	56.14	5.69	62.47	5.38
Image2StyleGAN++ [2]	25.36	6.71	40.39	7.10
InOut (Ours)	<b>23.57</b>	<u>7.18</u>	<b>30.34</b>	<b>7.16</b>
InOut-C (Ours)	29.24	<b>7.69</b>	<u>33.17</u>	<u>7.15</u>

$$\mathcal{L}_{\text{ms}} = \sum_{i=1}^m \sum_{j=i+1}^m \left( \frac{\|G(\mathbf{w}^i) - G(\mathbf{w}^j)\|_1}{\|\mathbf{w}^i - \mathbf{w}^j\|_1} \right). \quad (8)$$

The full objective of our optimization-based inversion is:

$$\begin{aligned} \arg \min_{\{\mathbf{w}_i\} \in \mathcal{W}} & \lambda_{\text{mse}} \mathcal{L}_{\text{mse}} + \lambda_{\text{percept}} \mathcal{L}_{\text{percept}} + \\ & \lambda_{\text{prior}} \mathcal{L}_{\text{prior}} + \lambda_{\text{div}} \mathcal{L}_{\text{div}} + \lambda_{\text{ms}} \mathcal{L}_{\text{ms}}, \end{aligned} \quad (9)$$

where the hyper-parameters  $\lambda$ 's control the importance of each term. Note that such an inversion paradigm is the same for both non-categorical and categorical settings, except the categorical setting seeks for  $\mathbf{w}_{\text{inter}}$  instead of  $\mathbf{w}$ .

### 3.3. Patch Blending

As described in Section 3.2, the inversion process requires the reconstruction of the given parts and the prediction of the outpainting parts, where the continuity and consistency are enforced in the training stage. However, even with the help of the prior loss  $\mathcal{L}_{\text{prior}}$ , the outpainting occasionally leads to tiny seams between patches after concatenating patches. Due to simple merging of patches, the outpainted images are likely to contain artifacts. As such, we introduce an image blending method to address this issue. In addition to the reference image  $R$  and outpainted area  $O$ , we generate the patches located halfway between  $R$  and  $O$ . Take  $R = I_{\text{micro}}^{1,1}, I_{\text{micro}}^{1,2}$  and  $O = I_{\text{micro}}^{2,1}, I_{\text{micro}}^{2,2}$  as an example, the additional region  $A$  is generated with coordinate  $(0, -1)$  and  $(0, 1)$ . We then linearly blend the overlapped area between  $R$  and  $A$  and the area between  $O$  and  $A$ .

Despite its simplicity, this post-processing step provides sufficient quality for our purpose. In practice, we observe that the generator can accurately interpolate the positions of the extended silhouette of landscapes with respect to the coordinate interpolation. The only rare artifacts of blending that occur in practice are in the categorical setting when a large foreground object (e.g., tower) is rendered from patches with slight ghosting artifacts.

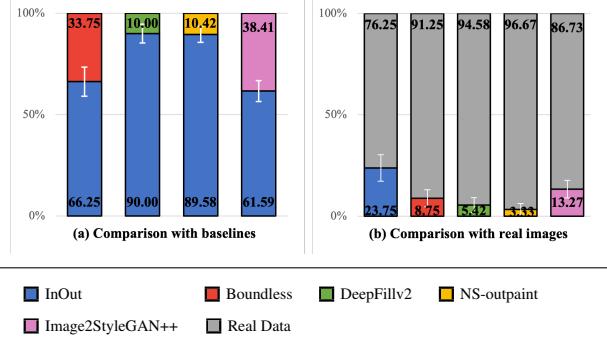


Figure 3: **User Studies.** We conduct user studies to quantify the visual quality in two settings: (a) comparison with baselines, and (b) comparison with real images. We mark the 95% confidence interval with white bars.

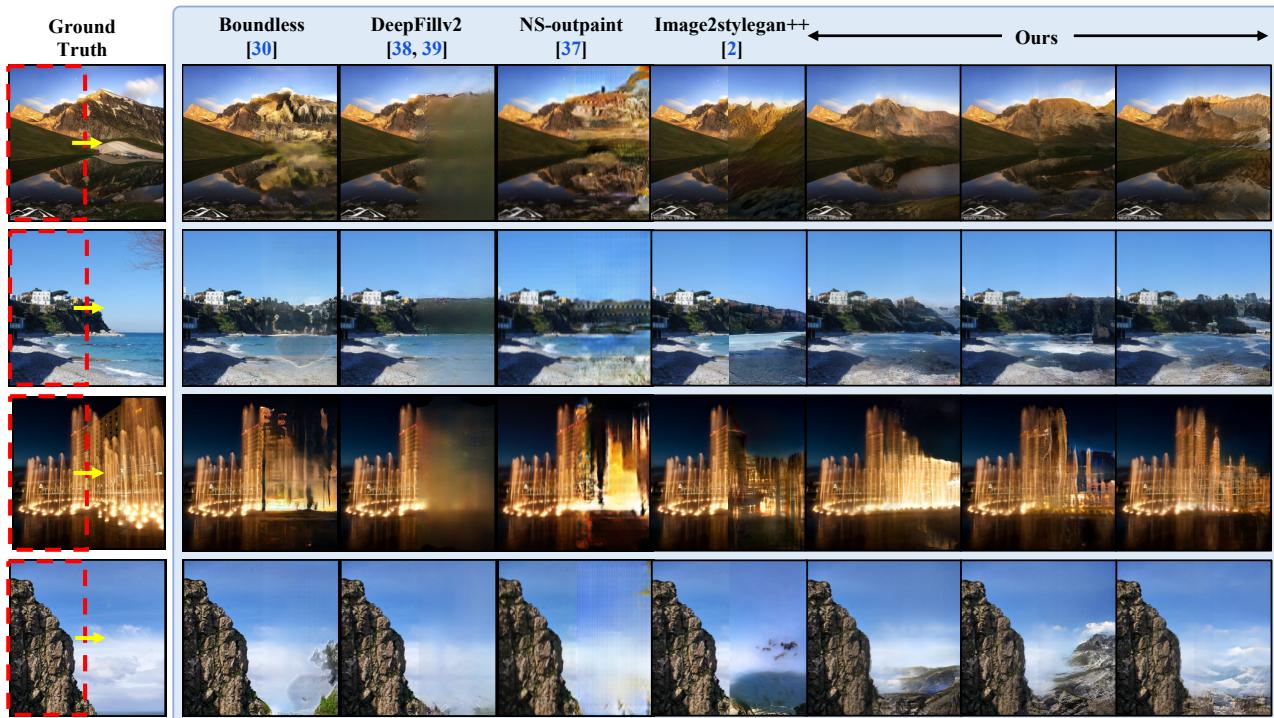
Table 2: **Ablation studies.** We show the necessity of each component using FID and LPIPS for quality and diversity.

# output	m=2		m=3	
	Method	FID ↓	Diversity ↑	FID ↓
InOut w/o $\mathcal{L}_{\text{div}}$ , $\mathcal{L}_{\text{ms}}$	30.12	<b>0.183</b>	30.28	<b>0.176</b>
InOut w/o $\mathcal{L}_{\text{ms}}$	29.85	0.201	29.80	0.206
InOut w/o $\mathcal{L}_{\text{div}}$	29.75	0.204	33.97	0.201
InOut w/o $\mathcal{L}_{\text{prior}}$	<b>36.56</b>	0.216	<b>36.53</b>	0.220
InOut	30.26	0.211	30.18	0.223

## 4. Experimental Results

**Dataset.** We evaluate our method on scenery datasets since they are the most representative and natural use-cases of out-painting. We perform experiments on the Places365 [42] dataset and a collected Flickr-Scenery dataset. More results on images with structured samples (e.g., buildings) could be found in the supplementary material. Similar to [30], we evaluate our method on a subset of the Places365 dataset. We select 25 scenery classes from the Places365 dataset with a subset of 62,500 samples. To further analyze the generalization of our method, we construct a Flickr-Scenery dataset by collecting a large-scale scenery image database of 54,710 images from Flickr. All images are center-cropped and resized to 256×256 pixels. For both datasets, we split the data into 80%, 10%, 10% for training, validation, and testing. All quantitative and qualitative experiments are evaluated on the testing set only. The source code, trained models, and Flickr-Scenery dataset will be made publicly available.

**Hyperparameters.** We use the default setup and parameters from the StyleGAN2 [16], including architecture, losses, optimizer, use of lazy regularizer, and implementation of inversion pipeline. The hyperparameters of our model are set as follows:



**Figure 4: Comparison to related work.** The qualitative comparison against other related methods show that the proposed approach is more stable, synthesizes richer context with more complex structures, and is able to handle some of the difficult complex scenes. (The input regions to all methods are marked with red dashes.)

- **Generator training.** Following the example in Section 3.2, we train our generator with a grid of  $2 \times 2$  micro-patches, and different from COCO-GAN [20], we train the discriminator with full images.
  - **Outpainting via inversion.** The weighting factors in Equation 9 are:  $\lambda_{\text{mse}} = 0.01$ ,  $\lambda_{\text{percept}} = 1$ ,  $\lambda_{\text{prior}} = 0.001$ ,  $\lambda_{\text{div}} = 0.001$ , and  $\lambda_{\text{ms}} = 0.001$ .

**Evaluated Methods.** We carry out quantitative and qualitative experiments with the state-of-the-art image outpainting methods (Boundless [30] and NS-outpaint [37]) and also image-inpainting methods (DeepFillv2 [38, 39] and Image2stylegan++ [2]).

#### **4.1. Quantitative Evaluation**

We evaluate the results in terms of realism and diversity using the Fréchet Inception Distance (FID) [12] and Inception Score (IS) [28]. Note that we *do not* apply the blending scheme introduced in Section 3.3 across the quantitative evaluations, as we aim to demonstrate the strength of the proposed pipeline without additional postprocessing.

As shown in Table 1, all of the proposed inversion-based methods perform favorably against I2I-based methods. The FID and IS results demonstrate that the generated-image distributions from our InOut variants to the real distribution are significantly more similar than the generation distribu-

tions from I2I-based baseline methods. Furthermore, compared to Image2stylegan++ [2], the results show that coordinate conditioning not only enables the categorical manipulation feature of InOut-C, but also naturally improves the generation diversity and quality of InOut-basic.

**User studies.** We conduct user studies to make explicit pairwise qualitative comparisons in two settings: (a) our method against each of the baseline methods, and (b) all methods against real samples. For each round of comparison, we present a pair of two outpainting results generated from the same real sample to the users. The images are either sampled from our method, baselines, or real images. Then, the subjects are asked to select a more realistic and preferred sample out of the image pair. We collect the results from 80 volunteers. Each of them makes 21 rounds of selections, resulting in 1,680 data points.

Figure 3 shows that subjects prefer the outpainting results by our model than those by other evaluated methods. Especially in comparison with real images, we observe a noticeable gap between the proposed InOut and Boundless. This may be attributed to that our method can frequently synthesize complex structures and novel objects (as shown in Figure 4 and Figure 6) that match the varieties and details of real images, while Boundless tends to create overly-smoothed results with raindrop-shaped artifacts.

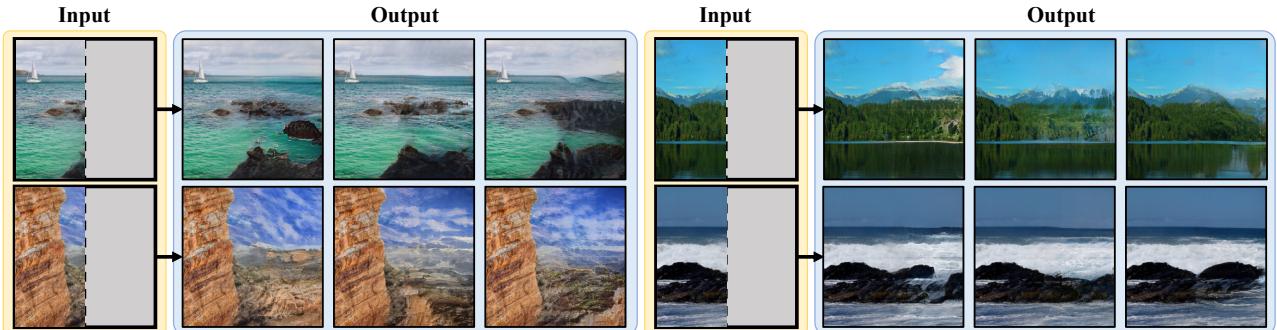
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714659  
660715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730

Figure 5: **Diverse outpainting.** We show that the proposed method can seek various solutions for a given input, achieving a high-variety of outpainting results without sacrificing the generation quality.

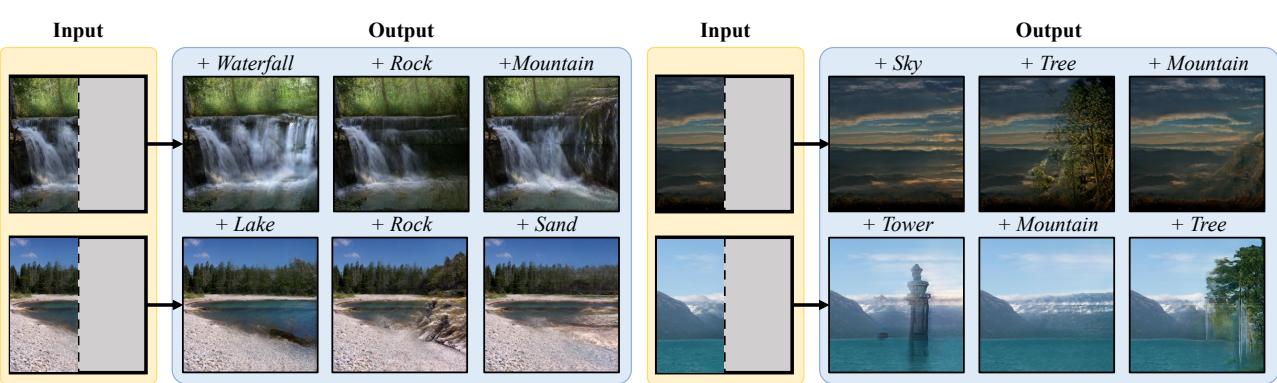
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

Figure 6: **Categorical generation.** We show the effectiveness of categorical manipulation by assigning different categorical labels to the outpainting area of the same real-image input. The results demonstrate that the proposed method can smoothly impose novel objects and calibrate the landscape to accommodate different categorical controls from users.

**Ablation studies.** In Table 2, we analyze and quantify the effectiveness of each component with ablation studies. We introduce a diversity score that measures the perceptual distance [41] among  $m$  outpainting results with respect to a real image. Each diversity score is averaged over 2,048 samples in the testing set. The diversity loss functions  $\mathcal{L}_{\text{ms}}$  and  $\mathcal{L}_{\text{div}}$  improve the diversity without compromising the quality. Note that the worst-case diversity quantity without any diversity loss is unlikely to be zero. This is due to the stochastic nature of gradient descent and randomization techniques introduced in [16] for encouraging the exploration during optimization. Furthermore, we show that the prior loss  $\mathcal{L}_{\text{prior}}$  is essential for securing the visual quality. As we have discussed in Section 3.2, the  $\mathcal{L}_{\text{prior}}$  regularizes the final state of the inverted latent codes to be located within the dense area of the Gaussian prior. Hence, the generated contents within the outpainted area remain realistic, instead of creating artifacts with the unconstrained latent codes that drift far away from the training distribution. As shown in Figure 8, inversion without  $\mathcal{L}_{\text{prior}}$  resulting in either obvious seams between input regions and outpainted areas or replicating input regions to the outpainted areas.

We evaluate the effect of different  $m$  values of  $\mathcal{L}_{\text{ms}}$  and  $\mathcal{L}_{\text{div}}$ . We demonstrate that the diversity losses provide sig-

nificant improvement in diversity score without compromising visual quality by seeking distinctive latent codes.

## 4.2. Qualitative Evaluation

In this section, we demonstrate the visual quality and diversity of the proposed method, and present applications, including categorical generation, panorama generation, and outpainting from different shapes and directions. Please refer to the supplementary materials for more visual results.

**Visual quality.** In Figure 4, we compare the visual quality of outpainting results from the proposed InOut against baselines. The results show that InOut is generally more realistic, coherent, diverse, exhibit more novel structures/objects, yet introduces fewer noticeable artifacts. In contrast, Boundless [30] tends to introduce raindrop-shaped artifacts, DeepFillv2 [38, 39] creates blurry extensions, while NS-outpaint [37] and Image2stylegan++ [2] frequently generates strong artifacts and obvious color differences.

**Diverse outpainting.** In Figure 5, we show the diverse outpainting results when  $m = 3$ . The results show that the proposed diversity loss enables the inversion pipeline to seek for different outpainting solutions. Notice that despite the variety of outpainting solutions, all inverted results remain visually compelling and matches the real-image input.

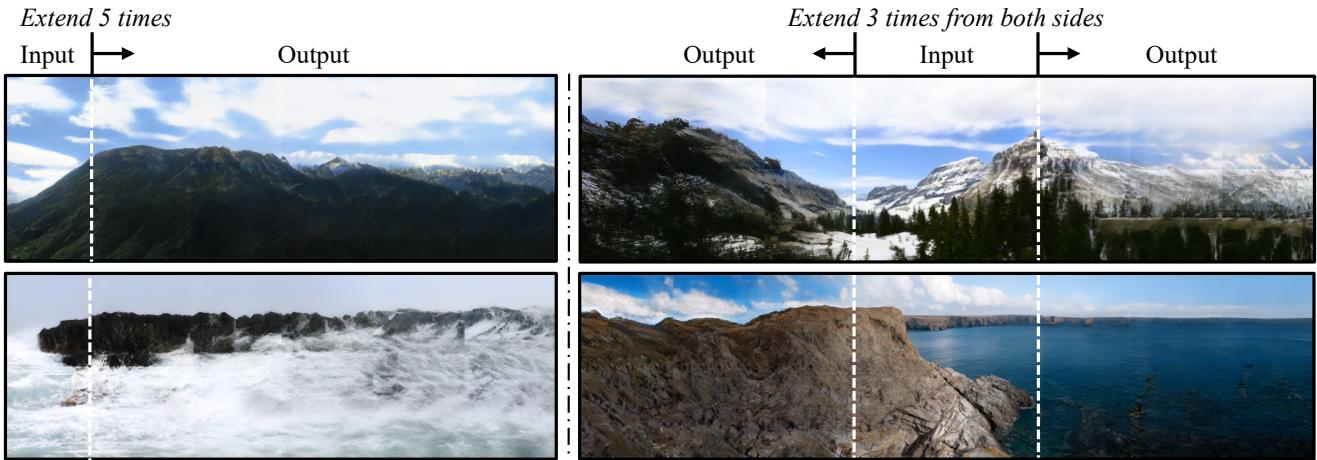


Figure 7: **Panorama generation.** We synthesize panoramic images by performing recursive outpainting. The results are of high quality and high structural complexity without repeating patterns.



Figure 8: **Qualitative on  $\mathcal{L}_{\text{prior}}$ .** Without  $\mathcal{L}_{\text{prior}}$ , the inversion will overfit to the reconstruction loss, resulting in latent codes extremely far away from the training distribution. The outpainted area may result in obvious seams (left) or replication of the input image (middle and right).

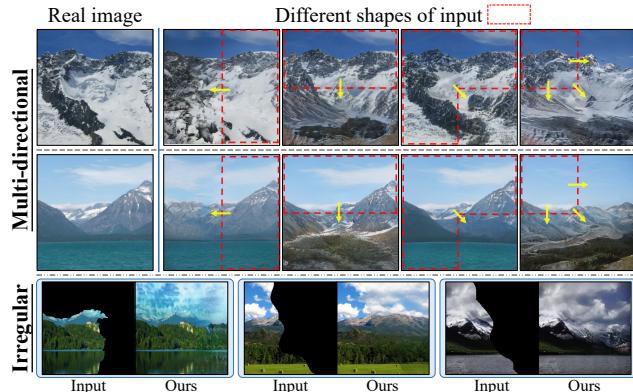


Figure 9: **Multi-directional and irregular boundary outpainting.** Our pipeline can inherently tackle (*top*) different outpainting directions and (*bottom*) irregular input shapes.

**Categorical Generation.** Figure 6 shows the results of categorical manipulation enabled with the InOut-C variant. Users can insert class-specific objects or manipulate the out-painted landscape structure with the categorical conditions of the two micro-patches on the right. The generator is able to automatically complete the background as well as blending the presented objects into the scene.

**Panorama generation.** Our framework naturally supports panorama generation by recursively taking previous out-painted micro-patches as the new inversion target. Figure 7

shows the panoramas generated from our method by outpainting to the left and right. The results show that the recursively outpainted area contains highly diverse structures without repeating patterns.

**Multi-directional and irregular-boundary outpainting.** For brevity, most results presented are generated given two micro-patches on the left-hand side as inputs. Nevertheless, the proposed method can perform outpainting from various directions with different input shapes and even arbitrary input shape. In the top of Figure 9, from left to right, we demonstrate outpainting results generated from the right, from the top, given three micro-patches, and given one micro-patch. In the bottom of Figure 9, we present outpainted results given inputs of irregular boundaries.

## 5. Conclusion

In this work, we tackle the image outpainting task from the GAN inversion perspective. We first train a generator to synthesize micro-patches conditioned on their positions. Based on the trained generator, we propose an inversion process that seeks for multiple latent codes recovering available regions as well as predicting outpainting regions. The proposed framework can generate diverse samples and support categorical specific outpainting, enabling more flexible user controls. Qualitative and quantitative experiments demonstrate the effectiveness of the proposed framework in terms of visual quality and diversity.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

864

## References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *ICCV*, 2019. 2
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *CVPR*, 2020. 2, 5, 6, 7
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. In *ICML*, 2017. 4
- [4] David Bau, Hendrik Strobelt, William Peebles, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. In *SIGGRAPH*, 2019. 2
- [5] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Neural photo editing with introspective adversarial networks. In *ICLR*, 2017. 2
- [6] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 4
- [7] Antonia Creswell and Anil Anthony Bharath. Inverting the generator of a generative adversarial network. *TNNLS*, 2018. 2
- [8] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. In *2017*, 2017. 2
- [9] Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. In *ICLR*, 2017. 2
- [10] Alexei A Efros and Thomas K Leung. Texture synthesis by non-parametric sampling. In *ICCV*, 1999. 3
- [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 2
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, 2017. 2, 6
- [13] Xun Huang, Ming-Yu Liu, Serge Belongie, and Jan Kautz. Multimodal unsupervised image-to-image translation. In *ECCV*, 2018. 2
- [14] Huaiyu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super slomo: High quality estimation of multiple intermediate frames for video interpolation. In *CVPR*, 2018. 2
- [15] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 3
- [16] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *CVPR*, 2020. 2, 3, 5, 7
- [17] Johannes Kopf, Wolf Kienzle, Steven Drucker, and Sing Bing Kang. Quality prediction for image completion. *ACM TOG (Proc. SIGGRAPH)*, 2012. 3

- [18] Hsin-Ying Lee, Hung-Yu Tseng, Jia-Bin Huang, Maneesh Kumar Singh, and Ming-Hsuan Yang. Diverse image-to-image translation via disentangled representations. In *ECCV*, 2018. 2
- [19] Yitong Li, Martin Renqiang Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from text. In *AAAI*, 2018. 2
- [20] Chieh Hubert Lin, Chia-Che Chang, Yu-Sheng Chen, Da-Cheng Juan, Wei Wei, and Hwann-Tzong Chen. Coco-gan: generation by parts via conditional coordinating. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4512–4521, 2019. 2, 3, 6
- [21] Zachary C Lipton and Subarna Tripathi. Precise recovery of latent vectors from generative adversarial networks. In *ICLR Workshop*, 2017. 2
- [22] Guilin Liu, Fitzsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *ECCV*, 2018. 2, 3
- [23] Fangchang Ma, Ulas Ayaz, and Sertac Karaman. Invertibility of convolutional generative networks from partial measurements. In *NeurIPS*, 2018. 2
- [24] Qi Mao, Hsin-Ying Lee, Hung-Yu Tseng, Siwei Ma, and Ming-Hsuan Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *CVPR*, 2019. 4
- [25] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier gans. In *International conference on machine learning*, pages 2642–2651, 2017. 4
- [26] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2, 3
- [27] Guim Perarnau, Joost Van De Weijer, Bogdan Raducanu, and Jose M Alvarez. Invertible conditional gans for image editing. In *NIPS Workshop*, 2016. 2
- [28] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *NIPS*, 2016. 6
- [29] Josef Sivic, Biliana Kaneva, Antonio Torralba, Shai Avidan, and William T Freeman. Creating and exploring a large photorealistic virtual space. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–8. IEEE, 2008. 3
- [30] Piotr Teterwak, Aaron Sarna, Dilip Krishnan, Aaron Maschinot, David Belanger, Ce Liu, and William T Freeman. Boundless: Generative adversarial networks for image extension. In *ICCV*, 2019. 2, 3, 5, 6, 7
- [31] Hung-Yu Tseng, Hsin-Ying Lee, Lu Jiang, Weilong Yang, and Ming-Hsuan Yang. Retrievegan: Image synthesis via differentiable patch retrieval. In *ECCV*, 2020. 2
- [32] Miao Wang, Yu-Kun Lai, Yuan Liang, Ralph R Martin, and Shi-Min Hu. Biggerpicture: data-driven image extrapolation using graph matching. *ACM TOG (Proc. SIGGRAPH)*, 2014. 3
- [33] Tsun-Hsuan Wang, Yen-Chi Cheng, Chieh Hubert Lin, Hwann-Tzong Chen, and Min Sun. Point-to-point video generation. In *ICCV*, 2019. 2

- 972 [34] Yi Wang, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Wide- 1026  
973 context semantic image extrapolation. In *CVPR*, 2019. 3 1027  
974 [35] Nevan Wichters, Ruben Villegas, Dumitru Erhan, and 1028  
975 Honglak Lee. Hierarchical long-term video prediction with- 1029  
976 out supervision. In *ICML*, 2018. 2 1030  
977 [36] Jonas Wulff and Antonio Torralba. Improving inversion and 1031  
978 generation diversity in stylegan using a gaussianized latent 1032  
979 space. *arXiv preprint arXiv:2009.06529*, 2020. 3, 4 1033  
980 [37] Zongxin Yang, Jian Dong, Ping Liu, Yi Yang, and Shuicheng 1034  
981 Yan. Very long natural scenery image prediction by outpaint- 1035  
982 ing. In *ICCV*, 2019. 2, 3, 5, 6, 7 1036  
983 [38] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and 1037  
984 Thomas S Huang. Generative image inpainting with contex- 1038  
985 tual attention. In *CVPR*, 2018. 2, 3, 5, 6, 7 1039  
986 [39] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and 1040  
987 Thomas S Huang. Free-form image inpainting with gated 1041  
988 convolution. In *ICCV*, 2019. 2, 3, 5, 6, 7 1042  
989 [40] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiao- 1043  
990 gang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stack- 1044  
991 gan: Text to photo-realistic image synthesis with stacked 1045  
992 generative adversarial networks. In *ICCV*, 2017. 2 1046  
993 [41] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, 1047  
994 and Oliver Wang. The unreasonable effectiveness of deep 1048  
995 features as a perceptual metric. In *CVPR*, 2018. 2, 4, 7 1049  
996 [42] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, 1050  
997 and Antonio Torralba. Places: A 10 million image database 1051  
998 for scene recognition. *TPAMI*, 2017. 2, 5 1052  
1000 [43] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In- 1053  
1001 domain gan inversion for real image editing. In *ECCV*, 2020. 1054  
1002 2 1055  
1003 [44] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and 1056  
1004 Alexei A Efros. Generative visual manipulation on the natu- 1057  
1005 ral image manifold. In *ECCV*, 2016. 2 1058  
1006 [45] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A 1059  
1007 Efros. Unpaired image-to-image translation using cycle- 1060  
1008 consistent adversarial networkss. In *ICCV*, 2017. 2 1061  
1009 1062  
1010 1063  
1011 1064  
1012 1065  
1013 1066  
1014 1067  
1015 1068  
1016 1069  
1017 1070  
1018 1071  
1019 1072  
1020 1073  
1021 1074  
1022 1075  
1023 1076  
1024 1077  
1025 1078  
1026 1079