

Multimodal Image Outpainting With Regularized Normalized Diversification

Lingzhi Zhang Jiancong Wang Jianbo Shi
University of Pennsylvania

{zlz, jshi}@seas.upenn.edu, jiancong.wang@pennmedicine.upenn.edu



Figure 1: Given only a small foreground region, our model can learn to outpaint a set of diverse and plausible missing backgrounds in both face image and street scene image.

Abstract

In this paper, we study the problem of generating a set of realistic and diverse backgrounds when given only a small foreground region. We refer to this task as image outpainting. The technical challenge of this task is to synthesize not only plausible but also diverse image outputs. Traditional generative adversarial networks suffer from mode collapse. While recent approaches [32, 28] propose to maximize or preserve the pairwise distance between generated samples with respect to their latent distance, they do not explicitly prevent the diverse samples of different conditional inputs from collapsing. Therefore, we propose a new regularization method to encourage diverse sampling in conditional synthesis. In addition, we propose a feature pyramid discriminator to improve the image quality. Our experimental results show that our model can produce more diverse images without sacrificing visual quality compared to state-of-the-arts approaches in both the CelebA face dataset [29] and the Cityscape scene dataset [2]. Code is available at: <https://github.com/owenzlz/DiverseOutpaint>

1. Introduction

Humans have the ability to hallucinate the possible backgrounds for a given object. For example when one shops for a couch (single foreground object) online, one can imagine how the couch might look inside the living room, one can also imagine how the couch might look in the office (various backgrounds). Is it possible for a machine to do the same?

In this paper, we aim to have the machine learn and synthesize a set of diverse and reasonable affordance backgrounds when given a foreground object, especially for cases where large portions of pixels are missing in an image. We refer to this task as image outpainting.

To outpaint the reasonable background for a foreground, the network has to understand the affordance relationship between the foreground and the background. For example, when given features of a person’s eyes, the machine needs to infer the possible facial expressions and other facial features of a person. Or given a car or a pedestrian pose, a machine needs to infer the street layout, as shown in Fig.(1). While affordance learning [10, 5, 22, 26, 35, 38, 41] aims to learn how objects interact in an environment, our task focuses on inverse affordance, which hallucinates the environment or background for the objects. Why would this task be useful besides generating interesting images? Some potential applications include facial recognition when large regions of the face are occluded, or synthesizing diverse images for product advertisements.

Intuitively, our task is multimodal common sense learning, which means there exists many possible backgrounds when given only a foreground region. The goal of this paper is to generate not only plausible but also diverse and thorough outputs. To make this more clear, we draw comparison between out-painting and the more common inpainting, the task of filling in missing pixels in an image. Within inpainting all background and some of the foreground objects are usually given and the semantic relationship between fore-

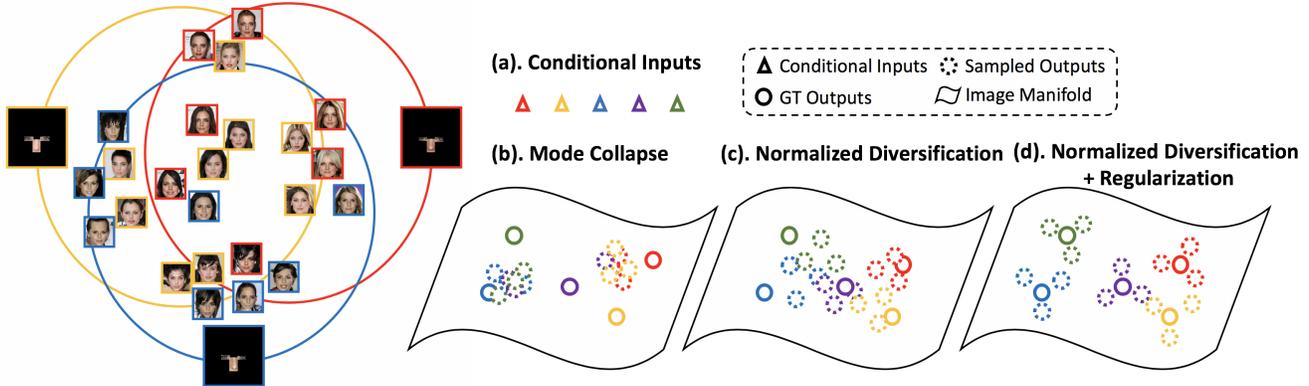


Figure 2: **Motivation** of our diversity regularization. Given conditional inputs (a), generative models could sample many outputs for each conditional input but collapse to a few modes (b). The current solution to this problem (c), normalized diversification [28], could preserve pairwise distance between sampled outputs for each conditional input, but it does not guarantee that the sampled outputs of different inputs could collapse together. Our solution (d) could not only preserve the pairwise distance of samples for each input but also prevent the samples of different inputs from collapsing together.

ground and background objects are pre-determined by the large portion of available pixels. This precludes the generator from hallucinating multiple possible semantic relationship between foreground/background objects. Outpainting, on the contrary, requires to fill in large portion of missing background and there is more degree of freedom in the common-sense semantic relationship that needs to be filled in by the generator. The outpainting therefore imposes higher diversity requirement on the generation framework.

We summarize the contributions of this work as follows. First, we formulated a new image outpainting task and provided a multimodal image synthesis solution. Second, we proposed a new diversity regularization technique to encourage diverse sampling without sacrificing image quality in this conditional synthesis task. In addition, we proposed a novel feature pyramid discriminator to check multi-scale information of outpainted images to improve visual quality. Overall, our proposed method can achieve more diversity and similar or better quality compared to the state-of-the-arts multimodal generative methods in both CelebA [29] face dataset and Cityscape [2] street dataset.

2. Related Work

2.1. Deep Generative Models

Deep generative models have produced exciting results. One type of generative models is Generative Adversarial Network (GAN) [9]. It consists of a generator network (G) and a discriminator network (D). During training, G tries to generate data as similar as the real data while D tries to differentiate the data from the real data. Once the adversarial training reaches an equilibrium, G is able to generate data that is indistinguishable from the real data distribution. Some applications of GAN will be elaborated in section 2.2.

Another popular generative model is variational auto-

encoder (VAE), which embeds high-dimensional data into a low-dimensional Gaussian distribution, samples a latent code and decodes it to the output space. The VAE framework is often used in multimodal prediction tasks, where the latent distribution models the uncertainty in the output space. For example, it has been used in multimodal image-to-image translation [54, 32, 21], predicting uncertain future motions [40, 6, 46], hallucinating diverse human affordance [42, 24] and so on.

We will discuss the related works that address the mode collapse issue in generative modeling. Mode collapse refers to the degenerate case where the generator produces limited or even single output mode. BourGAN [44] proposes to model the modes as a geometric structure of data distribution in a metric space, and uses mixture of Gaussians to construct latent space in order to map to different modes without collapse in unconditional generation. In conditional generation, mode seeking GAN (MSGAN) [32] proposes to maximize the ratio of two sampled images over their corresponding latent variables as a simple and intuitive diversity regularization. Lastly, normalized diversification [28] proposes to enforce the model to preserve the normalized pairwise distance between the sparse samples from a latent distribution to the corresponding high-dimensional output space. On top of the normalized diversification, we proposed a simple but effective diversity regularization to further encourage diversity in conditional image generation. The experimental results show that our method can generate more diverse images without sacrificing image quality compared to the state-of-the-arts approaches.

2.2. Conditional Image Synthesis

Deep generative models have been applied to many conditional image synthesis tasks. In super resolution

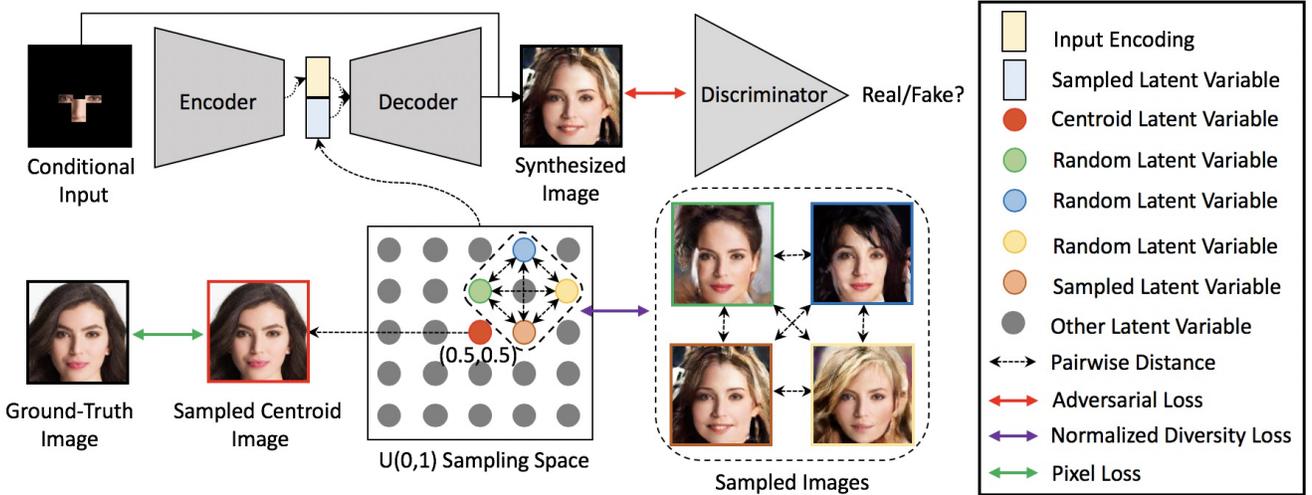


Figure 3: **Model Architecture.** The top part shows the architecture of the network, and the bottom part illustrates how the normalized diversification and diversity regularization are implemented.

[4, 18, 16, 49, 50, 51, 43], deep models learn the image texture prior to upsample a low-resolution image into high-resolution version. In style transfer [7, 16, 13, 30, 25], images can be transformed into an arbitrary style while its content being maintained by simultaneously minimizing the content and style loss w.r.t content and style images in feature space. In text-to-image synthesis [37, 17, 23], models can synthesize image layout and texture based on the input text.

Image inpainting, the task of filling parts of missing pixels in an image, is the most similar to our task among conditional image synthesis tasks. Early works [20, 45] train a deep convolutional network for denoising or inpainting small regions in the image. [36] proposes to learning useful features using image inpainting with adversarial training. [14] introduces the global and local discriminators to check the global and local consistency. [47] iteratively transverses the image manifold to find the closest encoding with respect to the input occluded image and uses it to decode the completed image. More recently, [27] introduces partial convolution, which is weighted to focus more on the valid regions rather than the hole regions. [48] first produces a coarse prediction of the missing region in the first stage, and then refines the texture-level details using an attention mechanism by searching for a set of background regions with the highest similarity with the coarse prediction. [34] proposes to inpaint an image by hallucinating the edge connection in the first stage, and then uses the connected edge map together with occluded image as inputs to produce the final completed image in the second stage.

Different from the above works, our goal is to produce diverse outputs conditional on a small foreground region.

Thus, we only compare our method to the methods that can generate multimodal image solutions.

3. Methods

In the image outpainting task, we aim to synthesize a set of plausible and diverse images when given a single foreground input. The previous approaches mostly leverage VAE [19] to encode a distribution of possible solutions and GAN [9] to synthesize realistic image. However, these approaches suffer from mode collapse. To build on top of normalized diversification[28], we proposed a new regularization technique to further encourage image diversity in this conditional image synthesis and a multi-scale discriminator to improve the visual quality.

3.1. Normalized Diversification

In normalized diversification, the generator learns mapping from a uniform latent space to an unknown output space. The key idea is to preserve the normalized pairwise distance of sparse samples between the latent space and the corresponding output space. In details, the Euclidean distance is used as the distance metric, which are shown in Eq.(1) and Eq.(2).

$$d_z(z_i, z_j) = \|z_i - z_j\|_2 \quad (1)$$

$$d_x(G(z)^i, G(z)^j) = \|G(z)^i - G(z)^j\|_2 \quad (2)$$

In above, $d_z(z_i, z_j)$ and $d_x(G(z)^i, G(z)^j)$ denote the pairwise distance of samples in latent space and output space respectively. We denote z as latent variable, G as generator, $G(z)$ as generated output, and i, j as sample indices.

The normalized pairwise distance matrices are further defined as $D_{ij}^z, D_{ij}^x \in \mathbb{R}^{N \times N}$ as follows.

$$D_{ij}^z = \frac{d_z(z_i, z_j)}{\sum_j d_z(z_i, z_j)} \quad (3)$$

$$D_{ij}^x = \frac{d_x(G(z)_i, G(z)_j)}{\sum_j d_x(G(z)_i, G(z)_j)} \quad (4)$$

During training, we treat the denominator in Eq.(3) and Eq.(4) as a constant when back-propagating the gradient to the generator network. This ensures that we optimize the absolute pairwise distance, rather than adjusting denominator to satisfy the loss constraint.

Finally, the normalized diversity constraint is implemented by a hinge loss, where we only penalize the generator when D_{ij}^x is smaller than D_{ij}^z multiplied by a scale factor.

$$\mathcal{L}_{ndiv}(x, z) = \frac{1}{N^2 - N} \sum_{i=1}^N \sum_{i \neq j}^N \max(0, \alpha D_{ij}^z - D_{ij}^x) \quad (5)$$

As shown in the diversity loss Eq.(5), α is the scale hyperparameter, and we do not consider the diagonal elements of the distance matrix, which are all zeros.

Coupled with the normalized diversity loss, the adversarial loss is used to check the generated diverse outputs are realistic compared to the real data distribution. For the discriminator,

$$\mathcal{L}_D = \mathbb{E}_{x \sim P_{data}(x)} [\min(0, 1 - D(x))] + \mathbb{E}_{z \sim P_{data}(z), z \sim P_z(z)} [\min(0, 1 + D(G(z)))] \quad (6)$$

For the generator,

$$\mathcal{L}_G = -\mathbb{E}_{z \sim P_z(z)} [D(G(z))] \quad (7)$$

We refer to \mathcal{L}_G as \mathcal{L}_{adv} in the following discussion.

Within our implementation, we use hinge loss to optimize the generator and the discriminator. Spectral normalization[33] is applied to scale down the weights in discriminator by their largest singular values, which effectively restricts the Lipschitz constant of the discriminator and thus stabilize training.

3.2. Diversity Regularization

In normalized diversification [28], the diversity loss L_{div} enforces the model to actively explore the output space while the adversarial loss L_{adv} constraints the generated outputs to be reasonable. However, in conditional generation, this framework only enforces the diversity for each conditional input, but do not explicitly prevent sampled outputs of different conditional inputs from collapsing to few modes. In particular, our image outpainting task has a large

degree of freedom to synthesize reasonable outputs for a foreground input, and thus the sampled images of a conditional input could be at very diverse locations on the image manifold. Therefore, it is very likely that the sampled images of different conditional inputs could be very visually similar or close on the image manifold.

To alleviate this issue, we propose a simple yet effective diversity regularization in addition to normalized diversification in this conditional synthesis task. The overall aim is to pull the diverse sampled outputs of different conditional inputs away from each other. Our approach is to impose a hard constraint on the generated outputs decoded from the center point of the uniform latent space, and enforce the generated and corresponding ground-truth outputs to be as similar as possible. This hard constraint is implemented by a pixel-wise Euclidean distance.

$$\mathcal{L}_{pixel} = \|G(z^*) - x\|_2 \quad (8)$$

In above, z^* denotes the center point in the uniform latent space and x is the ground-truth image corresponding to the conditional input. A visual demonstration of this insight is shown in Fig.(2).

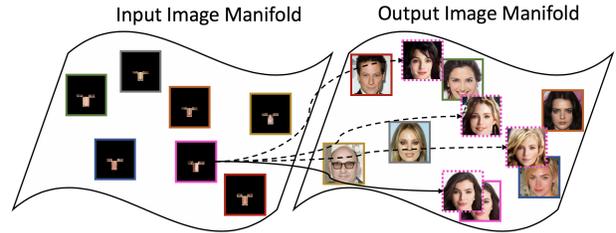


Figure 4: This figure demonstrates the motivation of identity regularization. The solid and dashed arrowed lines indicate the learned mapping with and without this regularization respectively.

With this diversity regularization, the sampling outputs of each conditional input will be ideally center around its corresponding ground truth output, as shown in Fig.(4). Although this regularization loss might not be fully optimized during training, the sampled outputs of different conditional inputs are pulled away from each other and thus alleviate mode collapse issue we mentioned earlier. The improvement of using this regularization is shown both qualitatively and quantitatively in section 4.

3.3. Feature Pyramid Discriminator

Generative Adversarial Network (GAN) [9] are commonly used in image synthesis. Prior to GAN, the synthesized images with only pixel reconstruction loss tends to be blurry. The main advantage of GAN is that the discriminator can provide supervisory signal for the generator to

synthesize realistic texture of images similar to the real data distribution. Indeed, a recent work [8] empirically finds out that CNNs tend to focus on or bias towards visual texture.

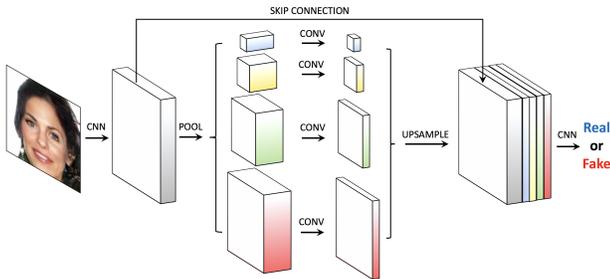


Figure 5: Feature Pyramid Discriminator.

How to design a discriminator that can check both texture realism and structural plausibility? Our insight is to explicitly design a discriminator network that can focus on both low-level textures and high-level structural semantics. Inspired by the pyramid scene parsing network PSPNet [53], we propose to integrate the pyramid pooling module that explicitly computes feature at multiple scales in the discriminator network as above.

This discriminator first extracts features of an image and downscales the features into multiple scales using average pooling. Then, the downscaled features are squeezed to fewer channel dimensions by a layer of convolution. Finally, the downscaled features are concatenated with original feature and are used jointly to compute the real or fake probability of an image. In summary, our feature pyramid discriminator aims to check multi-scale information of an image and is proven to consistently improve image quality across different datasets, as shown in Table.(2).

3.4. Implementation Details

Our generator network consists of an encoder and a decoder with skip connections at each spatial scale. The discriminator is described in section 3.3. Each convolution and deconvolution layers with stride of 2 are followed by a leaky relu layer with a negative slope of 0.2 and an instance normalization layer. The final output image is combined from the foreground input image and the synthesized output image.

At every step of updates, our model jointly optimizes the diversity loss, the adversarial loss as well as the diversity regularization loss. The overall optimization objective is shown Eq.(9).

$$\mathcal{L}_{total} = \lambda_1 * \mathcal{L}_{div} + \lambda_2 * \mathcal{L}_{adv} + \lambda_3 * \mathcal{L}_{reg} \quad (9)$$

The hyperparameters λ indicate the weights of different optimization objectives. During training, the loss functions

are jointly optimized by Adam optimizer with learning rate of $3e-4$, beta 1 of 0.5, and beta 2 of 0.99. We use $\lambda_1 = 0.1$, $\lambda_2 = 1$, $\lambda_3 = 5$ for loss weights.

4. Experiments

4.1. Preliminaries

Datasets. To generate synthetic training data, we sampled 400 conditional input points from a 2D uniform distribution and computed the corresponding outputs with a discrete non-linear transformation function. In the real image experiments, we used both CelebA [29] face dataset and Cityscape [2] street scene dataset. For the face dataset, we center-cropped and scaled the image down to be 128×128 , and cut out everything but two eyes and nose as inputs. The eyes and nose input region are localized by running a pre-trained facial landmark detector Super-FAN [1]. For the street scene dataset, we scaled down the images into 256×128 and then cut them by half into 128×128 as inputs. We used an instance segmentation network Mask R-CNN [11] to crop out the foreground region as inputs.

Evaluations. For the synthetic experiments, we evaluated the results by visualizing the sampled output space and calculating the generated data plausibility and diversity quantitatively. For the real image experiments, we did use FID[12] score to evaluate the image quality and pairwise LPIPS[52] score to quantify image diversity. The larger FID score indicates the better image quality, and the larger LPIPS score indicates the better image diversity. Qualitative evaluations are also provided.

4.2. Baseline Models

To demonstrate the effectiveness of our method, we compared the results with several state-of-the-arts models as strong baselines.

cVAE-GAN. [21] The conditional variational auto-encoder GAN (cVAE-GAN) encodes the input images into a parametric Gaussian distribution and decodes the sampled latent code into the output images, where the model is trained with a reconstruction loss, an adversarial loss and the KL divergence.

BicycleGAN. [54] The BicycleGAN model combines both cVAE-GAN and conditional latent regressor GAN (cLR-GAN) [3]. The key idea is to enforce the connection between latent encoding and output in both directions simultaneously.

MSGAN. [32] Built on top of the conditional GAN [15], the MSGAN model maximizes the ratio of the distance between generated images with respect to the corresponding latent codes in order to encourage the network to explore more minor modes in the data distribution.

NDiv. [28] The NDiv model preserves the normalized pairwise distance between the sparse samples from a latent

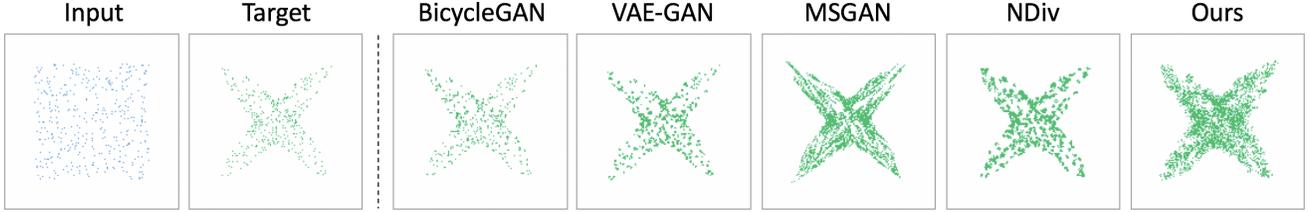


Figure 6: Qualitative results on synthetic data.

Methods	Frechet Distance ↓	Pairwise Distance ↑	Number of Modes ↑
BicycleGAN[54]	1.6477	0.0998	482
VAE-GAN[21]	0.8804	0.5898	962
NDiv[28]	1.2829	0.9917	1416
MSGAN[32]	1.2985	4.0199	2041
Ours	1.0751	4.3332	2481

Table 1: Quantitative results on synthetic data.

distribution and the generated output space, where the latent space is parameterized using a uniform distribution.

All of these approaches aim to generate multimodal solutions in conditional image synthesis. Both cVAE-GAN [21] and BicycleGAN [54] leverage VAE framework for variational inference but they do not explicitly enforce sampled diverse outputs, and thus mode collapse happens during both training and inference. MSGAN [32] proposes to enforce the generated outputs to be as diverse as possible with respect to the corresponding latent codes, but its Gaussian latent distribution puts a strong prior assumption on the output distribution. NDiv [28] does not explicitly prevent the sampled outputs of different conditional inputs from collapsing into a few modes in conditional generation, as we discussed in section 3.2.

4.3. Synthetic Data Experiment

To demonstrate the performance of diverse sampling in conditional generation, we start our experiments with a synthetic dataset. This dataset contains a set of sampled points from a uniform space $R^2 \sim U_2(0, 100)$ as inputs and the corresponding output points in a star-shaped space within the same range. We designed a discrete non-linear function to map the input points to the output points, and train the generative models to model such non-linear mapping. Both training and testing contain 400 sampled data points.

The task is to train a conditional generative model, given an input from the uniform space and a 2-dimension random vector sampled from either a normal distribution (BicycleGAN, VAE-GAN) or from a uniform distribution (NidV, ours), generate a corresponding point in the four-star space. In Fig.(6), the left two plots indicate the testing conditional inputs and ground truth outputs. The rest of the plots on the right are sampled outputs using different methods. During

inference, we sampled ten times for each conditional inputs. Qualitatively, the more distributed the sampled points lie in the output space indicate more diversity the model can produce. As shown in the figure, our model can generate more diverse outputs than other state-of-the-arts methods, since more sampled output points exist.

In addition, we demonstrate the quantitative comparison study in Table.(1). To evaluate the plausibility of the generated outputs, we use the Frechet Distance (FD) to compare the similarity between the generated output distribution and ground truth output distribution. The FD score is computed by averaging across ten batches of sampling outputs with respect to the ground truth. To evaluate the diversity of the generated outputs, we first compute the pairwise distance between the sampled outputs for each conditional input, and then calculate the number of existing output points (modes) in the 2D space. The number of mode is calculated by discretizing the generated output to the closest integer and number of different integers is counted. In conclusion, the quantitative results show that our model can generate most diverse and also plausible outputs compared to the other methods. Although VAE-GAN [21] achieves a slightly better score in terms of plausibility, it only generates limited modes in the output space.

4.4. Real Image Experiment

We conducted image outpainting experiments in both CelebA [29] face dataset and Cityscape [2] street scene dataset. In the testing phase, we used 1000 images from both datasets and sample 10 different output images for each conditional input.

To evaluate the generated image diversity, we first show the image manifold of 100 generated output images, as shown in Fig.(7). The 2D image manifold is obtained by

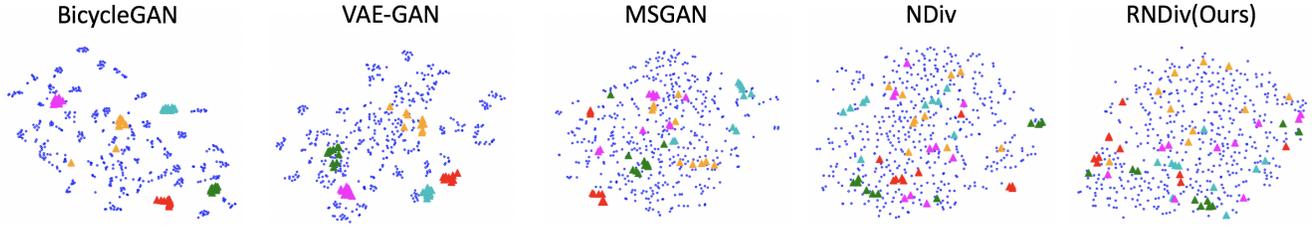


Figure 7: To visualize how the sampled images are located on the image manifold, we sampled 10 outpainted images for 100 testing input images in CelebA dataset[29]. Then, we extracted features for all images using pretrained VGG network[39] and ran t-sne[31] on the features to visualize the manifold in two dimension. The colored points (pink, orange, green, red, cyan) indicate the sampled images for five specific conditional inputs. Within the same color points, the more spread the points indicate more diverse the sampled outputs for the specific conditional input. The blue dots represent the rest of sampled images.

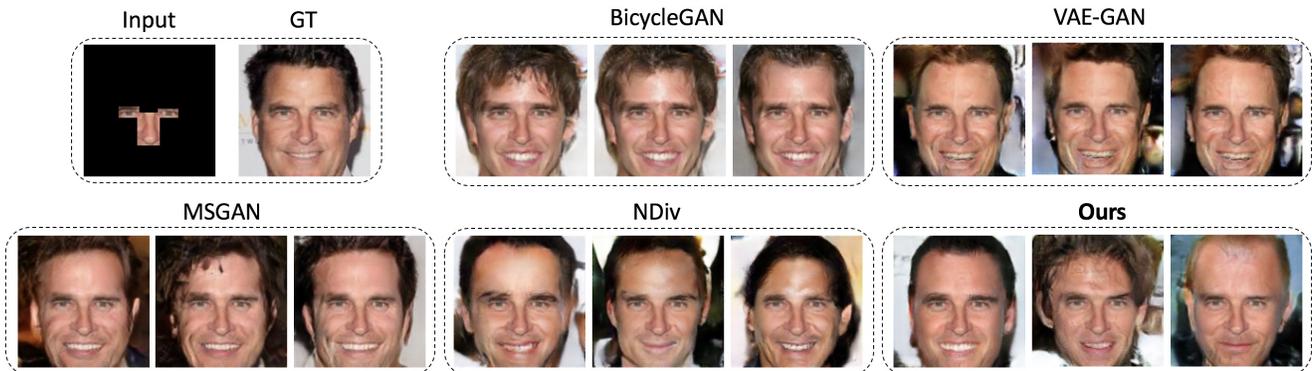


Figure 8: For a specific testing input image, we randomly sampled three possible outpainted images with different methods. We intended to demonstrate the diversity levels that each method can produce. Note that the mouths are the same in BicycleGAN[54], VAE-GAN[21], and MSGAN[32], and hair types are similar in NDiv[28]. In contrast, our method can generate both diverse types of hair styles and mouths.

Methods	CelebA[29]		CityScape[2]	
	Quality (FID) ↓	Diversity (LPIPS) ↑	Quality (FID) ↓	Diversity (LPIPS) ↑
BicycleGAN[54]	64.1328	0.0927	98.8635	0.0993
VAE-GAN[21]	66.3423	0.1754	77.8836	0.2915
MSGAN[32]	56.9978	0.2318	96.6312	0.3096
NDiv[28]	68.8545	0.3198	72.9145	0.4238
Ours (w/o FPD)	62.0442	0.3101	66.1893	0.4351
Ours (full model)	59.4232	0.3274	61.1454	0.4783

Table 2: Quantitative results on real images.

feature extraction using pretrained VGG network [39] and t-SNE [31] dimensionality reduction. In Fig.(7), we use five different color (pink, orange, green, red, cyan) to indicate sampled output images for five specific conditional input. This is intended to show the locations of diverse generated output images on the image manifold for the same conditional input. BicycleGAN [54] and VAE-GAN [21] both

have obvious mode collapse issue, since the sampled images are mostly clustered together into few modes, which leads to the big "holes" in the image manifold. With explicitly diversity losses, both MSGAN [32] and NDiv [28] can generate much more diverse outputs, but some generated images still collapse together, such as the red points. In contrast, our method can generate the most diverse out-

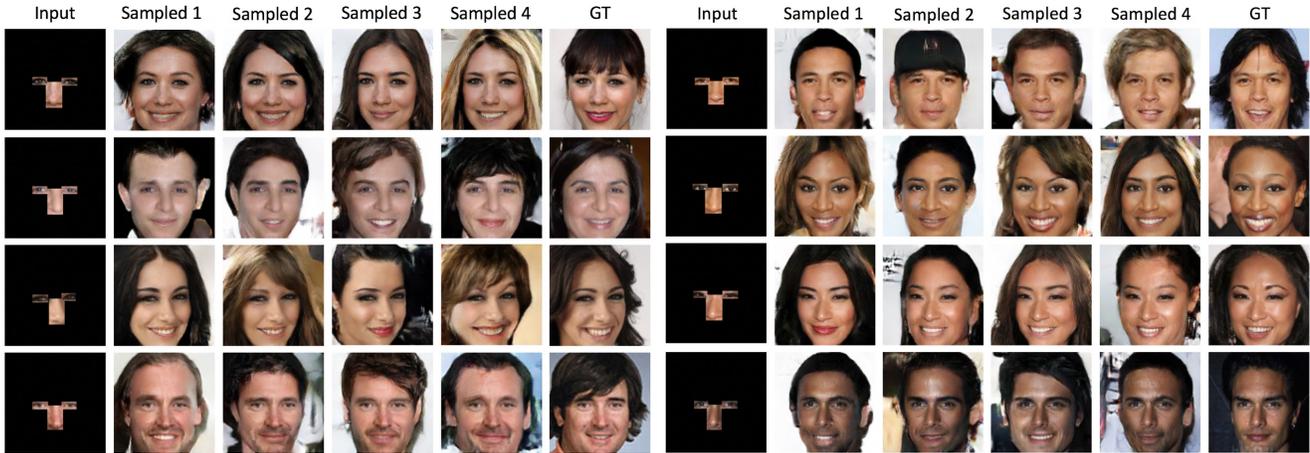


Figure 9: Qualitative results of sampled images in CelebA face dataset [29].

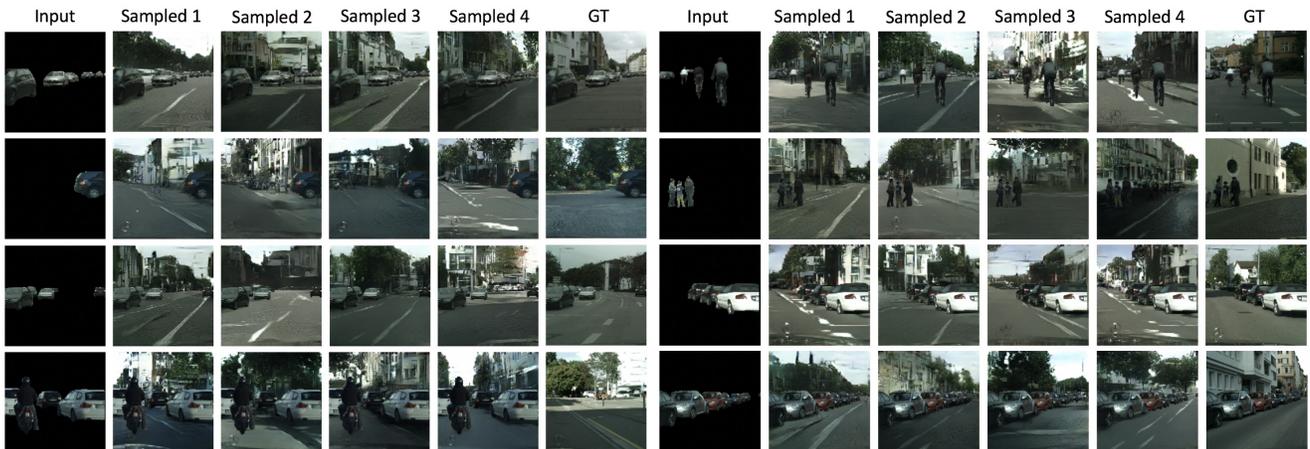


Figure 10: Qualitative results of sampled images in Cityscape street scene dataset [2].

puts compared to these methods. On the image manifold, almost all the randomly sampled outputs stays away from each other and thus results in a more expanded manifold than the others. A set of generated images from different methods are shown in Fig.(8).

In the quantitative evaluation, we use the Frechet Inception Distance (FID) [12] to measure the image quality and the pairwise Learned Perceptual Image Patch Similarity (LPIPS) [52] to measure the image diversity. The pairwise LPIPS score is computed between ten sampled generated images and is averaged across the entire testing set. As shown in Table.(2), our model can generate substantially more diverse images and achieve similar or better image quality compared to the state-of-the-arts methods on both datasets. Note that our proposed feature pyramid discriminator (FPD) improves the image quality and our proposed diversity regularization improves the image diversity consistently in both datasets.

5. Conclusion

In this paper, we formulated the image outpainting task, which aims to synthesize a set of realistic and diverse backgrounds when given only a small existing region. Based on the normalized diversification, we proposed a new regularization technique to further resolve mode collapse issue in this conditional image synthesis task. We also proposed a feature pyramid discriminator to improve the visual quality of generated images by checking image information at multi-scale. Finally, we demonstrated the effectiveness of our method compared to the state-of-the-arts methods in terms of both image diversity and quality in both synthetic dataset and two real-world datasets. In the future, we believe that our work could be applied to occluded face recognition for forensic purpose or common image editing, and could potentially extended to understand object affordance within detection/segmentation tasks.

References

- [1] A. Bulat and G. Tzimiropoulos. Super-fan: Integrated facial landmark localization and super-resolution of real-world low resolution faces in arbitrary poses with gans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 109–117, 2018. 5
- [2] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 1, 2, 5, 6, 7, 8
- [3] J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016. 5
- [4] C. Dong, C. C. Loy, K. He, and X. Tang. Image super-resolution using deep convolutional networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(2):295–307, 2016. 3
- [5] D. F. Fouhey, V. Delaitre, A. Gupta, A. A. Efros, I. Laptev, and J. Sivic. People watching: Human actions as a cue for single view geometry. *International journal of computer vision*, 110(3):259–274, 2014. 1
- [6] R. Gao, B. Xiong, and K. Grauman. Im2flow: Motion hallucination from static images for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5937–5947, 2018. 2
- [7] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2414–2423, 2016. 3
- [8] R. Geirhos, P. Rubisch, C. Michaelis, M. Bethge, F. A. Wichmann, and W. Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. 5
- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 2, 3, 4
- [10] A. Gupta, S. Satkin, A. A. Efros, and M. Hebert. From 3d scene geometry to human workspace. In *CVPR 2011*, pages 1961–1968. IEEE, 2011. 1
- [11] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 5
- [12] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017. 5, 8
- [13] X. Huang and S. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1501–1510, 2017. 3
- [14] S. Iizuka, E. Simo-Serra, and H. Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):107, 2017. 3
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 5
- [16] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pages 694–711. Springer, 2016. 3
- [17] J. Johnson, A. Gupta, and L. Fei-Fei. Image generation from scene graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1219–1228, 2018. 3
- [18] J. Kim, J. Kwon Lee, and K. Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1646–1654, 2016. 3
- [19] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [20] R. Köhler, C. Schuler, B. Schölkopf, and S. Harmeling. Mask-specific inpainting with deep neural networks. In *German Conference on Pattern Recognition*, pages 523–534. Springer, 2014. 3
- [21] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther. Autoencoding beyond pixels using a learned similarity metric. *arXiv preprint arXiv:1512.09300*, 2015. 2, 5, 6, 7
- [22] D. Lee, S. Liu, J. Gu, M.-Y. Liu, M.-H. Yang, and J. Kautz. Context-aware synthesis and placement of object instances. In *Advances in Neural Information Processing Systems*, pages 10393–10403, 2018. 1
- [23] W. Li, P. Zhang, L. Zhang, Q. Huang, X. He, S. Lyu, and J. Gao. Object-driven text-to-image synthesis via adversarial training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12174–12182, 2019. 3
- [24] X. Li, S. Liu, K. Kim, X. Wang, M.-H. Yang, and J. Kautz. Putting humans in a scene: Learning affordance in 3d indoor environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12368–12376, 2019. 2
- [25] Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M.-H. Yang. Universal style transfer via feature transforms. In *Advances in neural information processing systems*, pages 386–396, 2017. 3
- [26] C.-H. Lin, E. Yumer, O. Wang, E. Shechtman, and S. Lucey. St-gan: Spatial transformer generative adversarial networks for image compositing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9455–9464, 2018. 1
- [27] G. Liu, F. A. Reda, K. J. Shih, T.-C. Wang, A. Tao, and B. Catanzaro. Image inpainting for irregular holes using partial convolutions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 85–100, 2018. 3
- [28] S. Liu, X. Zhang, J. Wangni, and J. Shi. Normalized diversification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10306–10315, 2019. 1, 2, 3, 4, 5, 6, 7

- [29] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. 1, 2, 5, 6, 7, 8
- [30] F. Luan, S. Paris, E. Shechtman, and K. Bala. Deep photo style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4990–4998, 2017. 3
- [31] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008. 7
- [32] Q. Mao, H.-Y. Lee, H.-Y. Tseng, S. Ma, and M.-H. Yang. Mode seeking generative adversarial networks for diverse image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1429–1437, 2019. 1, 2, 5, 6, 7
- [33] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018. 4
- [34] K. Nazeri, E. Ng, T. Joseph, F. Qureshi, and M. Ebrahimi. Edgeconnect: Generative image inpainting with adversarial edge learning. *arXiv preprint arXiv:1901.00212*, 2019. 3
- [35] X. Ouyang, Y. Cheng, Y. Jiang, C.-L. Li, and P. Zhou. Pedestrian-synthesis-gan: Generating pedestrian data in real scene and beyond. *arXiv preprint arXiv:1804.02047*, 2018. 1
- [36] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2536–2544, 2016. 3
- [37] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee. Generative adversarial text to image synthesis. *arXiv preprint arXiv:1605.05396*, 2016. 3
- [38] A. Roy and S. Todorovic. A multi-scale cnn for affordance segmentation in rgb images. In *European conference on computer vision*, pages 186–201. Springer, 2016. 1
- [39] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 7
- [40] J. Walker, C. Doersch, A. Gupta, and M. Hebert. An uncertain future: Forecasting from static images using variational autoencoders. In *European Conference on Computer Vision*, pages 835–851. Springer, 2016. 2
- [41] X. Wang, R. Girdhar, and A. Gupta. Binge watching: Scaling affordance learning from sitcoms. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2596–2605, 2017. 1
- [42] X. Wang, R. Girdhar, and A. Gupta. Binge watching: Scaling affordance learning from sitcoms. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2596–2605, 2017. 2
- [43] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 3
- [44] C. Xiao, P. Zhong, and C. Zheng. Bourgan: Generative networks with metric embeddings. In *Advances in Neural Information Processing Systems*, pages 2269–2280, 2018. 2
- [45] L. Xu, J. S. Ren, C. Liu, and J. Jia. Deep convolutional neural network for image deconvolution. In *Advances in Neural Information Processing Systems*, pages 1790–1798, 2014. 3
- [46] T. Xue, J. Wu, K. Bouman, and B. Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Advances in neural information processing systems*, pages 91–99, 2016. 2
- [47] R. A. Yeh, C. Chen, T. Yian Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do. Semantic image inpainting with deep generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5485–5493, 2017. 3
- [48] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang. Generative image inpainting with contextual attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5505–5514, 2018. 3
- [49] X. Yu and F. Porikli. Ultra-resolving face images by discriminative generative networks. In *European Conference on Computer Vision*, pages 318–333. Springer, 2016. 3
- [50] X. Yu and F. Porikli. Face hallucination with tiny unaligned images by transformative discriminative neural networks. In *AAAI*, volume 2, page 3, 2017. 3
- [51] X. Yu and F. Porikli. Hallucinating very low-resolution unaligned and noisy face images by transformative discriminative autoencoders. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3760–3768, 2017. 3
- [52] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. *arXiv preprint*, 2018. 5, 8
- [53] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. 5
- [54] J.-Y. Zhu, R. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman. Toward multimodal image-to-image translation. In *Advances in Neural Information Processing Systems*, pages 465–476, 2017. 2, 5, 6, 7