# Image Fine-grained Inpainting

Zheng Hui, Jie Li, Xinbo Gao, *Senior Member, IEEE* and Xiumei Wang

*Abstract*—Image inpainting techniques have shown promising improvement with the assistance of generative adversarial networks (GANs) recently. However, most of them often suffered from completed results with unreasonable structure or blurriness. To mitigate this problem, in this paper, we present a one-stage model that utilizes dense combinations of dilated convolutions to obtain larger and more effective receptive fields. Benefited from the property of this network, we can more easily recover large regions in an incomplete image. To better train this efficient generator, except for frequently-used VGG feature matching loss, we design a novel self-guided regression loss for concentrating on uncertain areas and enhancing the semantic details. Besides, we devise a geometrical alignment constraint item (feature center coordinates alignment) to compensate for the pixel-based distance between prediction features and ground-truth ones. We also employ a discriminator with local and global branches to ensure local-global contents consistency. To further improve the quality of generated images, discriminator feature matching on the local branch is introduced, which dynamically minimizes the similarity of intermediate features between synthetic and ground-truth patches. Extensive experiments on several public datasets demonstrate that our approach outperforms current state-of-the-art methods. Code is available at https://github.com/Zheng222/DMFN.

*Index Terms*—image fine-grained inpainting, self-guided regression, geometrical alignment.

## I. INTRODUCTION

Image inpainting (a.k.a. image completion) aims to synthesize proper contents in missing regions of an image, which can be used in many applications. For instance, it allows removing unwanted objects in image editing tasks, while filling the contents that are visually realistic and semantically correct. Early approaches to image inpainting are mostly based on patches of low-level features. PatchMatch [1], a typical method, iteratively searches optimal patches to fill in the holes. It can produce plausible results when painting image background or repetitive textures. However, it cannot generate pleasing results for cases where completing regions include complex scenes, faces, and objects, which is due to PatchMatch cannot synthesize new image contents, and missing patches cannot be found in remaining regions for challenging cases.

With the rapid development of deep convolutional neural networks (CNN) and generative adversarial networks (GAN) [2], image inpainting approaches have achieved remarkable success. Pathak *et al.* proposed context-encoder [3], which employs a deep generative model to predict missing parts of the scene from their surroundings using reconstruction and adversarial losses. Yang *et al.* [4] introduced style

The authors are with the Video & Image Processing System (VIPS) Lab, School of Electronic Engineering, Xidian University, No.2, South Taibai Road, Xi'an 710071, China. (e-mail: zheng_hui@aliyun.com, leejie@mail.xidian.edu.cn, xbgao@mail.xidian.edu.cn, wangxm@xidian.edu.cn)

transfer into image inpainting to improve textural quality that propagates the high-frequency textures from the boundary to the hole. Li *et al.* [5] presented semantic parsing in the generation to restrict synthesized semantically valid contents for the missing facial key parts from random noise. To be able to complete large regions, Iizuka *et al.* [6] adopted stacked dilated convolutions in their image completion network to obtain lager spatial support and reached realistic results with the assistance of a globally and locally consistent adversarial training approach. Shortly afterward, Yu *et al.* [7] extended this insight and developed a novel contextual attention layer, which uses the features of known patches as convolutional kernels to compute the correlation between the foreground and background patches. More specifically, they calculated attention score for each pixel and then performed transposed convolution on attention score to reconstruct missing patches with known patches. It might be failing when the relationship between unknown and known patches is not close (*e.g.*masking all of the critical components of a facial image). Wang *et al.* [8] proposed a generative multi-column convolutional neural network (GMCNN) that uses varied receptive fields in branches by adopting different sizes of convolution kernels (*i.e.* $3 \times 3$, $5 \times 5$, and $7 \times 7$) in a parallel manner. This method produces advanced performance but suffers from substantial model parameters (12.562M) caused by large convolution kernels. In terms of image quality (more photo-realistic, fewer artifacts), it is still room for improvement.

The goals pursued by image inpainting are ensuring produced images with global semantic structure and finely detailed textures. Additionally, completed image should be approaching the ground truth as much as possible, especially for building and face images. Previous techniques more focus on solving how to yield holistically reasonable and photo-realistic images. This problem has been mitigated by GAN [2] or its improved version WGAN-GP [9] that is frequently utilized in image inpainting methods [3], [6], [4], [5], [7], [10], [11], [8], [12], [13]. However, concerning fine-grained details, there is still much room to enhance. Besides, these existing methods haven't taken into account the consistency between outputs and targets, *i.e.*, semantic structures should be as much similar as possible for facial images and building images.

To overcome the limitations of the methods as mentioned above, we present a unified generative network for image inpainting, which is denoted as *dense multi-scale fusion network* (DMFN). The *dense multi-scale fusion block* (DMFB), serving as the basic block of DMFN, is composed of four-way dilated convolutions as illustrated in Figure 2. This basic block adopts the combination and fusion of hierarchical features extracted from various convolutions with different dilation rates to obtain better multi-scale features, compared with general dilated convolution (dense v.s. sparse). For generating

Fig. 1. The inpainted results on FFHQ dataset [14] by using our method. The missing areas are shown in white. *It is worth noting that they also recover well in terms of lighting and texture.*

images with the realistic and semantic structure, we design a *self-guided regression loss* that constrains low-level features of the generated content according to the normalized discrepancy map (the difference between the output and target). *Geometrical alignment constraint* is developed for penalizing the coordinate center of estimated image high-level features away from the ground-truth. This loss can further help the processing of image fine-grained inpainting. We improve the discriminator using relativistic average GAN (RaGAN) [15]. It is noteworthy that we use global and local branches in the discriminator as in [6], where one branch focuses on the global image while the other concentrates on the local patch of the missing region. To explicitly constraint the output and ground-truth images, we utilize the hidden layers of the local branch that belongs to the whole discriminator to evaluate their discrepancy through an adversarial training process. With all these improvements, the proposed method can produce high-quality results on multiple datasets, including faces, building, and natural scene images.

Our contributions are summarized as follows:

- Self-guided regression loss corrects semantic structure errors to some extent through re-weighing VGG features guided by discrepancy map, which is novel for image/video completion task.
- We present the geometrical alignment constraint to supplement the shortage of pixel-based VGG features matching loss, which restrains the results with a more reasonable semantic spatial location.
- We propose the dense multiple fusion block (DMFB, enhancing the dilated convolution) to improve the network

representation, which increases the receptive field while maintaining an acceptable parameter size. Our generative image inpainting framework achieves compelling visual results (as illustrated in Figure 1) on challenging datasets.

As shown in Table I, we summarized the difference between the typical method and the proposed approach. We are committed to improving the dilated convolution that frequently used in image completion, and developing more losses to measure the matching degree of features from different viewpoint. More details can be found in Section III.

The rest of this paper is organized as follows. Section II provides a brief review of related inpainting methods. Section III describes the proposed approach and loss functions in detail. In Section IV, we explain the experiments conducted for this work, experimental comparisons with other state-of-the-art methods, and model analysis. In Section V, we conclude the study.

## II. RELATED WORK

A variety of algorithms for image inpainting have been proposed. Traditional diffusion-based methods [17], [18] propagate information from neighboring regions to the holes. They can work well for small and narrow holes, where the texture and color variance are the same. However, these methods fail to recover meaning contents in the large missing regions. Patch-based approaches, such as [19], [20], search for relevant patches from the known regions in an iterative fashion. Simakov *et al.* [21] proposed bidirectional similarity scheme to capture better and summarize non-stationary visual data. However, these methods are computationally expensive due to

TABLE I
IMAGE INPAINTING METHODOLOGY EMPLOYED BY SOME REPRESENTATIVE MODELS.

| Method | GMCNN (NeurIPS'2018) [8] | GC (ICCV'2019) [16] | DMFN (Ours) |
|---|---|---|---|
| Stage | one-stage | two-stage | one-stage |
| Generator | multi-column CNNs | gated CNN | dense multi-scale fusion network |
| Discriminator | WGAN-GP | SN-PatchGAN | RelativisticGAN |
| Losses | reconstruction + adversarial + ID-MRF | $l1$ reconstruction + SN-PatchGAN | $l1$ + self-guided regression + fm_vgg + fm_dis + adversarial + alignment |

calculating the similarity scores of each output-target pair. To relieve this problem, PatchMatch [1] is proposed, which speeds it up by designing a faster similar patch searching algorithm. Ding *et al.* [22] proposed a Gaussian-weighted nonlocal texture similarity measure to obtain multiple candidate patches for each target patch.

Recently, deep learning and GAN-based algorithms have been a remarkable paradigm for image inpainting. Context Encoders (CE) [3] embed the $128 \times 128$ image with a $64 \times 64$ center hole as a low dimensional feature vector and then decode it to a $64 \times 64$ image. Iizuka *et al.* [6] proposed a high-performance completion network with both global and local discriminators that is critical in obtaining semantically and locally consistent image inpainting results. Also, the authors employ the dilated convolution layers to increase receptive fields of the output neurons. Yang *et al.* [4] use intermediate features extracted by pre-trained VGG network [23] to find hole's most similar patch outside the hole. This approach performs multi-scale neural patch synthesis in a coarse-to-fine manner, which noticeably takes a long time to fill a large image during the inference stage. For face completion, Li *et al.* [5] trained a deep generative model with a combination of reconstruction loss, global and local adversarial losses, and a semantic parsing loss specialized for face images. Contextual Attention (CA) [7] adopted two-stage network architecture where the first step produces a crude result, and the second refinement network using attention mechanism takes the coarse prediction as inputs and improves fine details. Liu *et al.* [24] introduced partial convolution that employs computational operations only on valid pixels and presented an auto-update binary mask to determinate whether the current pixels are valid. Substituting convolutional layers with partial convolutions can help a UNet-like architecture [25] achieve the state-of-the-art inpainting results. Yan *et al.* [11] introduced a special shift-connection to the U-Net architecture for enhancing the sharp structures and fine-detailed textures in the filled holes. This method was mainly developed on building and natural landscape images. Similar to [4], [7], Song *et al.* [10] decoupled the completion process into two stages: coarse inference and fine textures translation. Nazeri *et al.* [26] also proposed a two-stage network that comprises of an edge generator and an image completion network. Similar to this method, Li *et al.* [27] progressively incorporated edge information into the feature to output more structured image. Xiong *et al.* [28] inferred the contours of the objects in the image, then used the completed contours as a guidance to complete the image. Different from frequently-used two-stage processing [29], Sagong *et al.* [30] proposed parallel path for semantic inpainting to reduce the computational costs.
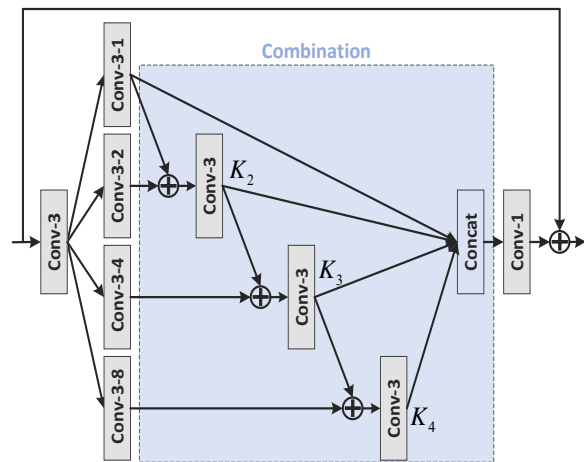
## III. PROPOSED METHOD



Fig. 2. The architecture of the proposed dense multi-scale fusion block (DMFB). Here, "Conv-3-8" indicates $3 \times 3$ convolution layer with the dilation rate of 8 and $\oplus$ is element-wise summation. Instance normalization (IN) and ReLU activation layers followed by the first convolution, second column convolutions and concatenation layer are omitted for brevity. The last convolutional layer only connects an IN layer. The number of output channels for each convolution is set to 64 except for the last $1 \times 1$ convolution (256 channels) in DMFB.

Our proposed inpainting system is trained in an end-to-end way. Given an input image with hole $\mathbf{I}_{in}$, its corresponding binary mask $\mathbf{M}$ (value 0 for known pixels and 1 denotes unknown ones), the output $\mathbf{I}_{out}$ predicted by the network, and the ground-truth image $\mathbf{I}_{gt}$. We take the input image and mask as inputs, *i.e.*, $[\mathbf{I}_{in}, \mathbf{M}]$. We now elaborate on our network as follows.

### A. Network structure

As depicted in Figure 3, our framework consists of a generator, and a discriminator with two branches. The generator produces plausible painted results, and the discriminator conducts adversarial training.

For image inpainting task, the size of the receptive fields should be sufficiently large. The dilated convolution is popularly adopted in the previous works [6], [7] to accomplish this purpose. This way increases the area that can use as input without increasing the number of learnable weights. However, the kernel of dilated convolution is sparse, which skips many pixels during applying to compute. Large convolution kernel (*e.g.*$7 \times 7$) is applied in [8] to implement this intention. However, this solution introduces heavy model parameters. To enlarge the receptive fields and ensure dense convolution kernels simultaneously, we propose our dense multi-scale fusion
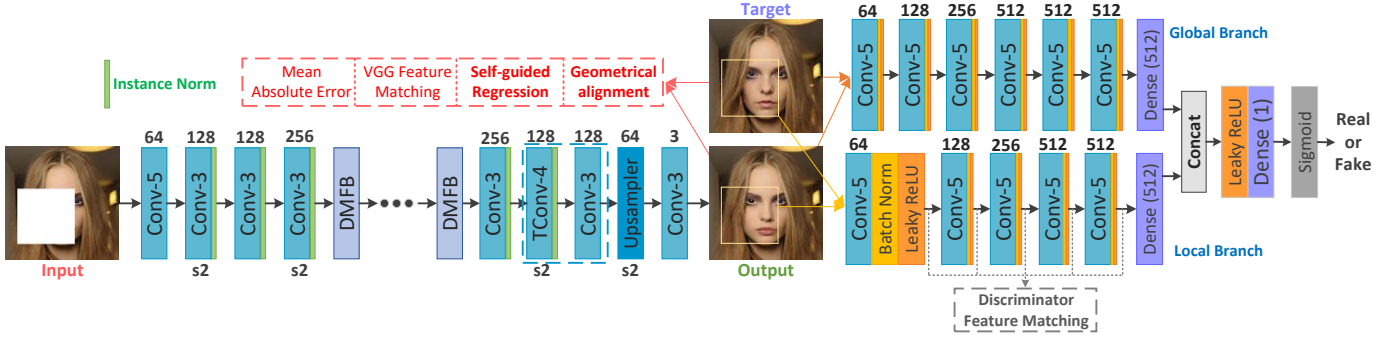
Fig. 3. The framework of our method. The activation layer followed by each "convolution + norm" or convolution layer in the generator is omitted for conciseness. The activation function adopts ReLU except for the last convolution (Tanh) in the generator. Blue dotted box indicates our upsampler module (TConv-4 is $4 \times 4$ transposed convolution) and "$s2$" denotes the stride of 2.

block (DMFB, see in Figure 2) inspired by [31]. Specifically, the first convolution on the left in DMFB reduces the channels of input features to $64$ for decreasing the parameters, and then these processed features are sent to four branches to extract multi-scale features, denoted as $\mathbf{x}_i$ ($i = 1, 2, 3, 4$), by using dilated convolutions with different dilation factors. Except for $\mathbf{x}_1$, each $\mathbf{x}_i$ has a corresponding $3 \times 3$ convolution, denoted by $K_i(\cdot)$. Through a cumulative addition fashion, we can get dense multi-scale features from the combination of various sparse multi-scale features. We denote by $\mathbf{y}_i$ the output of $K_i(\cdot)$. The combination part can be formulated as

$$\mathbf{y}_i = \begin{cases} \mathbf{x}_i, & i = 1; \\ K_i(\mathbf{x}_{i-1} + \mathbf{x}_i), & i = 2; \\ K_i(\mathbf{y}_{i-1} + \mathbf{x}_i), & 2 < i \leq 4. \end{cases} \quad (1)$$

The following step is the fusion of concatenated features simply using a $1 \times 1$ convolution. In a word, this basic block especially enhances the general dilated convolution and has fewer parameters than large kernels.

Different from previous generative inpainting networks [7], [8] that apply WGAN-GP [9] for adversarial training, we propose to use RaGAN [15] to pursue more photo-realistic generated images [32]. This discriminator also considers the consistency of global and local images.
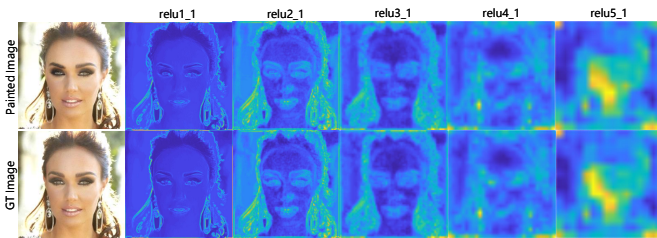
*B. Loss functions*



Fig. 4. Visualization of average VGG feature maps.

*1) Self-guided regression loss:* Here, we address the semantic structure preservation issue. We scheme to take self-guided regression constraint to correct the image semantic level estimation. Briefly, we compute the discrepancy map between

generated contents and corresponding ground truth to navigate the similarity measure of the feature map hierarchy from the pre-trained VGG19 [23] network. At first, we investigate the characteristic of VGG feature maps. Given an input image $\mathbf{I}_A$, it is first fed forward through the VGG19 to yield a five-level feature map pyramid, where their spatial resolution reduces low progressively. Specifically, the $l$-th ($l = 1, 2, 3, 4, 5$) level is set to the feature tensor produced by relu$l\_1$ layer of VGG19. These feature tensors are denoted by $F_A^l$. We give an illustration of average feature maps $F_{A\_avg}^l = \frac{1}{M} \sum_{m=1}^{M} F_{A\_m}^l$ in Figure 4, which suggests that the deeper layers of a pre-trained network represent higher-level semantic information, while lower-level features more focus on textural or structural details, such as edges, corners, and other simple conjunctions. In this paper, we would intend to improve the detail fidelity of the completed image, especially for building and face images.
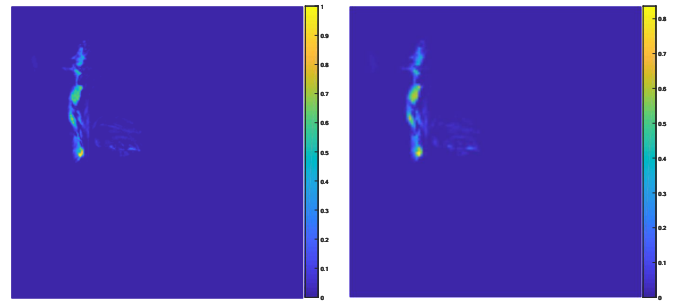


Fig. 5. Visualization of guidance maps. (Left) Guidance map $\mathbf{M}_{guidance}^1$ for "relu1_1" layer. (Right) Guidance map $\mathbf{M}_{guidance}^2$ for "relu2_1" layer. These are corresponding to Figure 4.

To this end, through the error map between the output image produced by the generator and ground truth, we get the guidance map to distinguish between areas of challenging and manageable. Therefore, we propose to use the following equation to gain the average error map:

$$\mathbf{M}_{error} = \frac{1}{3} \sum_{c \in \mathcal{C}} (\mathbf{I}_{out,c} - \mathbf{I}_{gt,c})^2, \quad (2)$$

where $\mathcal{C}$ are the three color channels, $\mathbf{I}_{out,c}$ denotes $c$-th channel of the output image. Then, the normalized guidance

mask can be calculated by

$$\mathbf{M}_{guidance,p} = \frac{\mathbf{M}_{error,p} - \min\left(\mathbf{M}_{error}\right)}{\max\left(\mathbf{M}_{error}\right) - \min\left(\mathbf{M}_{error}\right)}, \quad (3)$$

where $\mathbf{M}_{error,p}$ is the error map value at position $p$. Note that our guidance mask with continuous values between 0 and 1, which is soft instead of binary. $\mathbf{M}_{guidence}^l$ corresponds $l$-th level feature maps and it can be expressed by

$$\mathbf{M}_{guidance}^{l+1} = AP\left(\mathbf{M}_{guidance}^l\right), \quad (4)$$

where $AP$ denotes *average pooling* with kernel size of 2 and stride of 2. Here, $\mathbf{M}_{guidance}^1 = \mathbf{M}_{guidance}$ (Equation 3). In this way, the value range of $\mathbf{M}_{guidance}^l$ is still between 0 and 1. In view of the fact that lower-level feature map contains more detailed information, we choose feature tensors from "relu1_1" and "relu2_1" layers to describe image semantic structures. Thus, our self-guided regression loss is defined as

$$\mathcal{L}_{self-guided} = \sum_{l=1}^{2} w^l \frac{\left\|\mathbf{M}_{guidance}^l \odot \left(\Psi_{\mathbf{I}_{gt}}^l - \Psi_{\mathbf{I}_{output}}^l\right)\right\|_1}{N_{\Psi_{\mathbf{I}_{gt}}^l}}, \quad (5)$$

where $\Psi_{\mathbf{I}_*}^l$ is the activation map of the relu$l$_1 layer given original input $\mathbf{I}_*$, $N_{\Psi_{\mathbf{I}_{gt}}^l}$ is the number of elements in $\Psi_{\mathbf{I}_{gt}}^l$, $\odot$ is the element-wise product operator, and $w^l = \frac{1e3}{\left(C_{\Psi_{\mathbf{I}_{gt}}^l}\right)^2}$ followed by [33]. Here, $C$ is the channel size of feature map $\Psi_{\mathbf{I}_{gt}}^l$.

An obvious benefit for this regularization is to suppress regions with higher uncertainty (as shown in Figure 5). $\mathbf{M}_{guidance}$ can be viewed as a spatial attention map, which preferably optimizes areas that are difficult to handle. Our self-guided regression loss is performed lower-level semantic space instead of pixel space. The merit of this way would appear in the perceptual image synthesis with pleasant structural information.

*2) Geometrical alignment constraint:* In the typical solutions, the metric evaluation in higher-level feature space is only achieved using pixel-based loss, *e.g.*, L1 or L2. It doesn't take the alignment of each high-level feature map semantic hub into account. To better measure the distance between high-level features belong to prediction and ground-truth, we impose the geometrical alignment constraint on the response maps of "relu4_1" layer. This term can help the generator create a plausible image that aligned with the target image in position. Specifically, this term encourages the output feature center to be spatially close to the target feature center. The geometrical center for the $k$-th feature map along axis $u$ is calculated as

$$c_u^k = \sum_{u,v} u \cdot \left(\mathbf{R}\left(k,u,v\right) \Big/ \sum_{u,v} \mathbf{R}\left(k,u,v\right)\right), \quad (6)$$

where response maps $\mathbf{R} = \mathrm{VGG}\left(\mathbf{I};\theta_{vgg}\right) \in \mathbb{R}^{K \times H \times W}$. $\mathbf{R}\left(k,u,v\right) \Big/ \sum_{u,v} \mathbf{R}\left(k,u,v\right)$ represents a spatial probability distribution function. $c_u^k$ denotes coordinate expectation along axis $u$. Then, we pass both the completed image $\mathbf{I}_{output}$ and ground-truth image $\mathbf{I}_{gt}$ through the VGG network and obtain the corresponding response maps $\mathbf{R}'$ and $\mathbf{R}$. Given these response maps, we compute the centers $\left\langle c_u^{k'}, c_v^{k'}\right\rangle$ and $\left\langle c_u^k, c_v^k\right\rangle$ using Equation 6. Then, we formulate the geometrical alignment constraint as

$$\mathcal{L}_{align} = \sum_k \left\|\left\langle c_u^{k'}, c_v^{k'}\right\rangle - \left\langle c_u^k, c_v^k\right\rangle\right\|_2^2. \quad (7)$$

*3) Feature matching losses:* The VGG feature matching loss $\mathcal{L}_{fm\_vgg}$ compares the activation maps in the intermediate layers of well-trained VGG19 [23] model, which can be written as

$$\mathcal{L}_{fm\_vgg} = \sum_{l=1}^{5} w^l \frac{\left\|\Psi_{\mathbf{I}_{gt}}^l - \Psi_{\mathbf{I}_{output}}^l\right\|_1}{N_{\Psi_{\mathbf{I}_{gt}}^l}}, \quad (8)$$

where $N_{\Psi_{\mathbf{I}_{gt}}^l}$ is the number of elements in $\Psi_{\mathbf{I}_{gt}}^l$. Inspired by [34], we also introduce local branch in discriminator feature matching loss $\mathcal{L}_{fm\_dis}$, which is reasonable to assume that the output image are consistent with the ground-truth images under any measurements (*i.e.*, any high-dimensional spaces). This feature matching loss is defined as

$$\mathcal{L}_{fm\_dis} = \sum_{l=1}^{5} w^l \frac{\left\|D_{local}^l\left(\mathbf{I}_{gt}\right) - D_{local}^l\left(\mathbf{I}_{output}\right)\right\|_1}{N_{D_{local}^l\left(\mathbf{I}_{gt}\right)}}, \quad (9)$$

where $D_{local}^l\left(\mathbf{I}_*\right)$ is the activation in the $l$-th selected layer of the discriminator given input $\mathbf{I}_*$ (see in Figure 3). Note that the hidden layers of the discriminator are trainable, which is slightly different from the well-trained VGG19 network trained on the ImageNet dataset. It can adaptively update based on specific training data. This complementary feature matching can dynamically extract features that may be not mined in VGG model.

*4) Adversarial loss:* For improving the visual quality of inpainted results, we use relativistic average discriminator [15] as in ESRGAN [32], which is the recent state-of-the-art perceptual image super-resolution algorithm. For the generator, the adversarial loss is defined as

$$\mathcal{L}_{adv} = -\mathbb{E}_{x_r}\left[\log\left(1 - D_{Ra}\left(x_r, x_f\right)\right)\right]$$
$$- \mathbb{E}_{x_f}\left[\log\left(D_{Ra}\left(x_f, x_r\right)\right)\right], \quad (10)$$

where $D_{Ra}\left(x_r, x_f\right) = sigmoid\left(C\left(x_r\right) - \mathbb{E}_{x_f}\left[C\left(x_f\right)\right]\right)$ and $C\left(\cdot\right)$ indicates the discriminator network without the last *sigmoid* function. Here, real/fake data pairs $\left(x_r, x_f\right)$ are sampled from ground-truth and output images.

*5) Final objective:* With self-guided regression loss, geometrical alignment constraint, VGG feature matching loss, discriminator feature matching loss, adversarial loss, and mean absolute error (MAE) loss, our overall loss function is defined as

$$\mathcal{L}_{total} = \mathcal{L}_{mae} + \lambda\left(\mathcal{L}_{self-guided} + \mathcal{L}_{fm\_vgg}\right)$$
$$+ \eta\mathcal{L}_{fm\_dis} + \mu\mathcal{L}_{adv} + \gamma\mathcal{L}_{align}, \quad (11)$$

where $\lambda$, $\eta$, $\mu$, and $\gamma$ are used to balance the effects between the losses mentioned above.

## IV. EXPERIMENTS

We evaluate the proposed inpainting model on Paris Street View [3], Places2 [35], CelebA-HQ [36], and a new challenging facial dataset FFHQ [14].
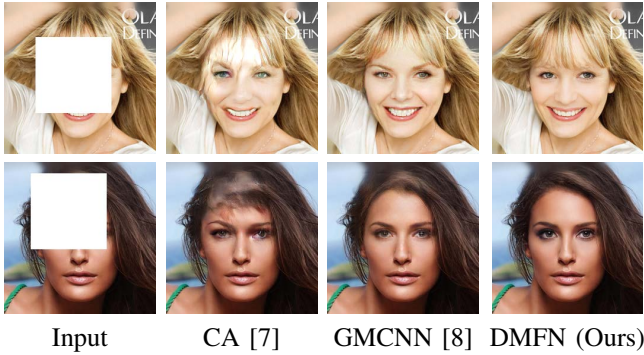
| Input | CA [7] | GMCNN [8] | DMFN (Ours) |

Fig. 6. Visual comparisons on CelebA-HQ. *Best viewed with zoom-in.*

## A. Experimental settings

For our experiments, we empirically set $\lambda = 25$, $\eta = 5$, $\mu = 0.003$ and $\gamma = 1$ in Equation 11. The training procedure is optimized using Adam optimizer [37] with $\beta_1 = 0.5$ and $\beta_2 = 0.9$. We set the learning rate to $2e - 4$. The batch size is 16. We apply PyTorch framework to implement our model and train them using NVIDIA TITAN Xp GPU (12GB memory).

For training, given a raw image $\mathbf{I}_{gt}$, a binary image mask $\mathbf{M}$ (value 0 for known pixels and 1 denotes unknown ones) at a random position. In this way, the input image $\mathbf{I}_{in}$ is obtained from the raw image as $\mathbf{I}_{in} = \mathbf{I}_{gt} \odot (\mathbf{1} - \mathbf{M})$. Our inpainting generator takes $[\mathbf{I}_{in}, \mathbf{M}]$ as input, and produces prediction $\mathbf{I}_{pred}$. The final output image is $\mathbf{I}_{out} = \mathbf{I}_{in} + \mathbf{I}_{pred} \odot \mathbf{M}$. All input and output are linearly scaled to $[-1, 1]$. We train our network on the training set and evaluate it on the validation set (Places2, CelebA-HQ, and FFHQ) or testing set (Paris street view and CelebA). For training, we use images of resolution $256 \times 256$ with the largest hole size $128 \times 128$ as in [7], [8]. For Paris street view ($936 \times 537$), we randomly crop patches with resolution $537 \times 537$ and then scale down them to $256 \times 256$ for training. Similarly, for Places2 ($512 \times *$), $512 \times 512$ sub-images are cropped at a random location. These images are scaled down to $256 \times 256$ for our model. For CelebA-HQ and FFHQ face datasets ($1024 \times 1024$), images are directly scaled to 256. We use the irregular mask dataset provided by [24]. For irregular masks, the random regular regions are cropped and sent to the local discriminator. All results generated by our model are not post-processed.

## B. Qualitative comparisons

As shown in Figures 7, 6, and 8, compared with other state-of-the-art methods, our model gives a noticeable visual improvement on textures and structures. For instance, our network generates plausible image structures in Figure 7, which mainly stems from the dense multi-scale fusion architecture and well-designed losses. The realistic textures are hallucinated via feature matching and adversarial training. For Figure 6, we show that our results with more realistic details and fewer artifacts than the compared approaches. Besides, we give partial results of our method and PICNet [12] on Places2 dataset in Figure 8. The proposed DMFN creates more reasonable, natural, and photo-realistic images. Additionally, we also show some example results (masks at random position)

of our model trained on FFHQ in Figure 9. In Figure 10, our method performs more stable and fine for large-area irregular masks than compared algorithms. More compelling results can be found in the *supplementary material.*

## C. Quantitative comparisons

Following [7], [8], we measure the quality of our results using peak signal-to-noise ratio (PSNR) and structural similarity (SSIM). Learned perceptual image patch similarity (LPIPS) [39] is a new metric that can better evaluate the perceptual similarity between two images. Because the purpose of image inpainting is to pursue visual effects, we adopt LPIPS as the main qualitative assessment. The lower the values of LPIPS, the better. In Places2, 100 validation images from "canyon" scene category are chosen for evaluation. As shown in Table II, our method produces acceptable results compared with CA [7], GMCNN [8], PICNet [12], and PENNet [13] in terms of all evaluation measurements.

We also conducted user studies as illustrated in Figure 12. The scheme is based on blind randomized A/B/C tests deployed on Google Forms platform as in [8]. Each survey includes 40 single-choice questions. Each question involves three options (completed images that are generated from the same corrupted input by three different methods). There are 20 participants invited to accomplish this survey. They are asked to select to the most realistic item in each question. The option order is shuffled each time. Finally, our method outperforms compared approaches by a large margin.

## D. Ablation study

*1) Effectiveness of DMFB:* To validate the representation ability of our DMFB, we replace its middle part (4 dilated convolutions and combination operation) to a $3 \times 3$ dilated convolution (256 channels) with dilation rate of 2 or 8 ("rate=2" or "rate=8", see in Table III). Additionally, to verify the strength of $K_i(\cdot)$ in combination operation, we perform the DMFB without $K_i(\cdot)$ that denoted as "w/o $K_i(\cdot)$" in Table III. Combined with Table III and Figure 13, we can clearly see that our model with DMFB (**Parms**: $471, 808$) predicts reasonable and less artifact images than ordinary dilated convolutions (**Parms**: $803, 392$). Specifically, in the second row of Figure 13, both "w/o combination" and "w/o $K_i(\cdot)$" have different degrees of the partial absence of lampposts. The visual effects of them can be obviously ranked as "DMFB" > "w/o $K_i(\cdot)$" > "w/o combination". Meanwhile, the results of "rate=2" and "rate=8" suggest the importance of spatial support as discussed in [6]. It also demonstrates large and dense receptive field is beneficial to completing images with large holes.

*2) Self-guided regression and geometrical alignment constraint:* To investigate the effect of proposed self-guided regression loss and geometrical alignment constraint, we train a complete DMFN on CelebA-HQ dataset without the corresponding loss. Thanks to the effectiveness of the generator and the loss functions, the baseline model (w/o self-guided regression loss & geometrical alignment constraint) can already achieve almost satisfactory results. Because of the subtle problems existing in the baseline model, we propose two loss

Fig. 7. Visual comparisons on Paris street view.

TABLE II
QUANTITATIVE RESULTS (CENTER REGULAR MASK) ON FOUR TESTING DATASETS.

| Method | Paris street view (100) | Places2 (100) | CelebA-HQ (2,000) | FFHQ (10,000) |
| --- | --- | --- | --- | --- |
| | LPIPS / PSNR / SSIM | LPIPS / PSNR / SSIM | LPIPS / PSNR / SSIM | LPIPS / PSNR / SSIM |
| CA [7] | N/A | 0.1524 / 21.32 / 0.8010 | 0.0724 / 24.13 / 0.8661 | N/A |
| GMCNN [8] | 0.1243 / 24.38 / 0.8444 | 0.1829 / 19.51 / 0.7817 | 0.0509 / 25.88 / 0.8879 | N/A |
| PICNet [12] | 0.1263 / 23.79 / 0.8314 | 0.1622 / 20.70 / 0.7931 | N/A | N/A |
| PENNet [13] | N/A | 0.2384 / 21.93 / 0.7586 | 0.0676 / 25.50 / 0.8813 | N/A |
| DMFN (Ours) | **0.1018 / 25.00 / 0.8563** | **0.1188 / 22.36 / 0.8194** | **0.0460 / 26.50 / 0.8932** | **0.0457 / 26.49 / 0.8985** |



Fig. 8. Visual comparisons on Places2.

TABLE III
QUANTITATIVE RESULTS OF DIFFERENT STRUCTURES ON PARIS STREET VIEW DATASET (CENTER REGULAR MASK).

| Model | rate=2 | rate=8 | w/o combination | w/o $K_i(\cdot)$ | DMFB |
| --- | --- | --- | --- | --- | --- |
| Params | 803,392 | 803,392 | 361,024 | 361,024 | 471,808 |
| LPIPS↓ | 0.1059 | 0.1067 | 0.1083 | 0.1026 | **0.1018** |
| PSNR↑ | 24.93 | 24.91 | 24.24 | 24.93 | **25.00** |
| SSIM↑ | 0.8530 | 0.8549 | 0.8505 | 0.8561 | **0.8563** |

functions to further refine the results, the so-called image fine-grained inpainting. For instance, in the first row of Figure 14, "w/o alignment" shows that the eyelid lines at the left eye (in



Fig. 9. Visual results on FFHQ dataset.

red box) are dislocation, the eyebrows are slightly scattered. "w/o self-guided" yields correct double eyelids, but eyebrows are still unnatural. "with all" shows the best performance. Although the qualitative performance is not much improved, these new loss functions have a corrective effect on a few problematic results produced by the baseline model. And we give the quantitative results in Table IV, which validates the
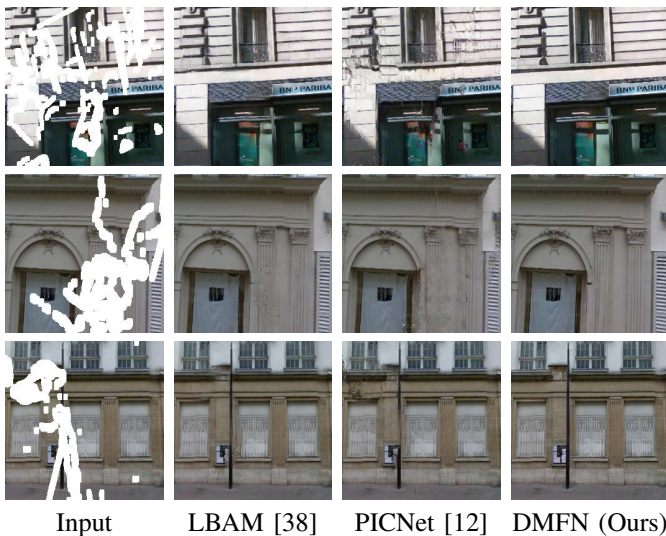
Input  LBAM [38]  PICNet [12]  DMFN (Ours)

Fig. 10. Inpainted images with irregular masks on Paris StreetView.



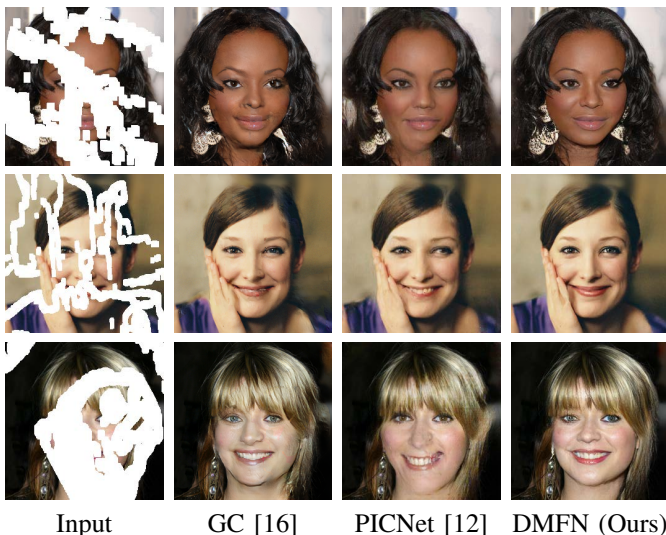Input  GC [16]  PICNet [12]  DMFN (Ours)

Fig. 11. Inpainted images with irregular masks on CelebA-HQ.

effectiveness of various proposed losses.

### E. Discussions

*1) Intention of self-guided regression loss:* The prior works (*e.g.*, CA [7] and GMCNN [8]) assign less weight at masked region centers to formulate the variant of L1 loss. CA only use spatial discounted L1 loss in the coarse network (first stage). GMCNN first train their model with only confidence-driven L1 loss. Without the assistant of GAN, these first stages only aim to obtain a coarse results. Different from them, the self-guided

TABLE IV
INVESTIGATION OF SELF-GUIDED REGRESSION LOSS AND GEOMETRICAL
ALIGNMENT CONSTRAINT ON CELEBA-HQ (RANDOM REGULAR MASK).

| Metric | w/o self-guided | w/o align | w/o dis_fm | with all |
|---|---|---|---|---|
| LPIPS↓ | 0.0537 | 0.0534 | 0.0542 | **0.0530** |
| PSNR↑ | 25.73 | 25.63 | 25.65 | **25.83** |
| SSIM↑ | 0.8892 | 0.8884 | 0.8870 | **0.8892** |



Fig. 12. Results of user study.



Input  rate=2  rate=8  w/o combination  w/o $K_i(\cdot)$  DMFB (Ours)
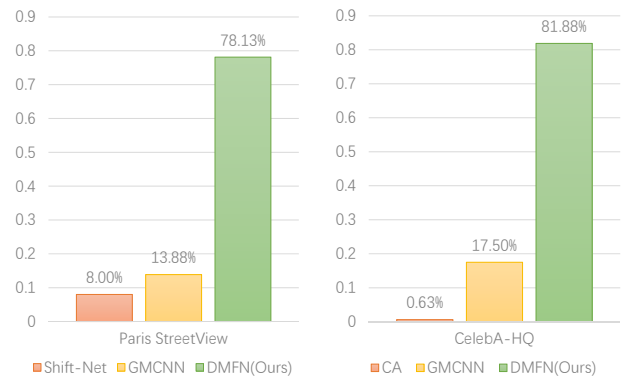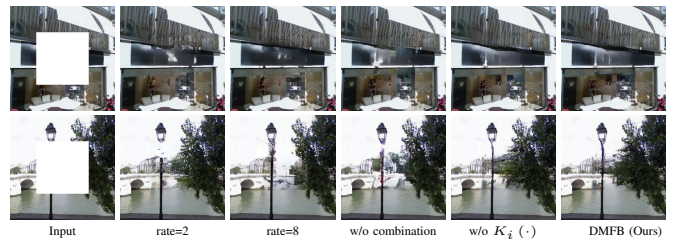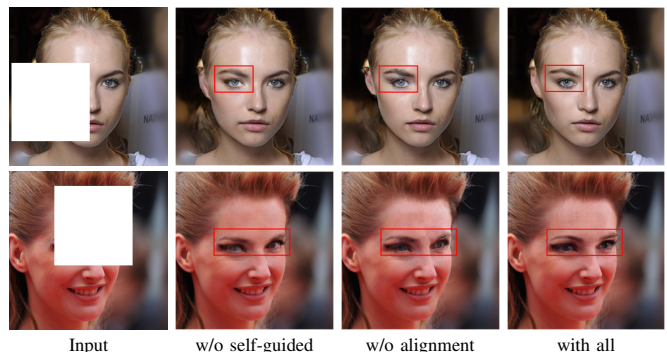
Fig. 13. Visual comparison of different structures. *Best viewed with zoom-in.*

regression loss apply to VGG features focuses on learning the *hard* areas measured by the current guidance map. And our framework is one-stage trained with all losses at the same time.

*2) Analysis of self-guided regression loss:* In this section, we conduct the ablation study of using different distance metrics in the average error map. Table V compares instantiations including *Gaussian*, *dot product*, and *L2* when used in self-guided regression loss. The simple L2 performs the best LPIPS performance, and Gaussian achieves the best PSNR and SSIM performance. Considering that the purpose of image inpainting is the pursuit of plausible visual effect, We choose the simple and efficient L2 distance to measure the average error map. Figure 15 shows the visual comparisons among these metrics, which indicates L2 can recover the better structural information.

*3) Investigation of geometrical alignment constraint:* As illustrated in Figure 16, we visualize the first $64$ feature maps of each selected VGG layer. The response maps of



Input  w/o self-guided  w/o alignment  with all

Fig. 14. Visual comparison of results using different losses. *Best viewed with zoom-in.*

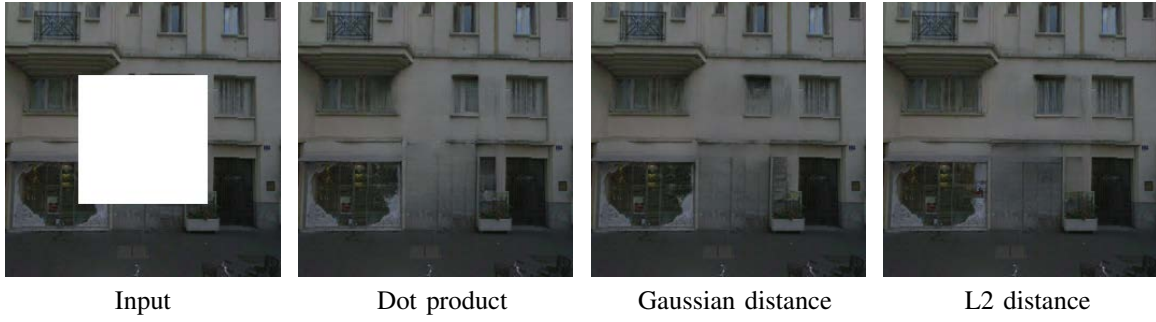| Input | Dot product | Gaussian distance | L2 distance |

Fig. 15. Visual comparisons on Paris street view. *Best viewed with zoom-in.*
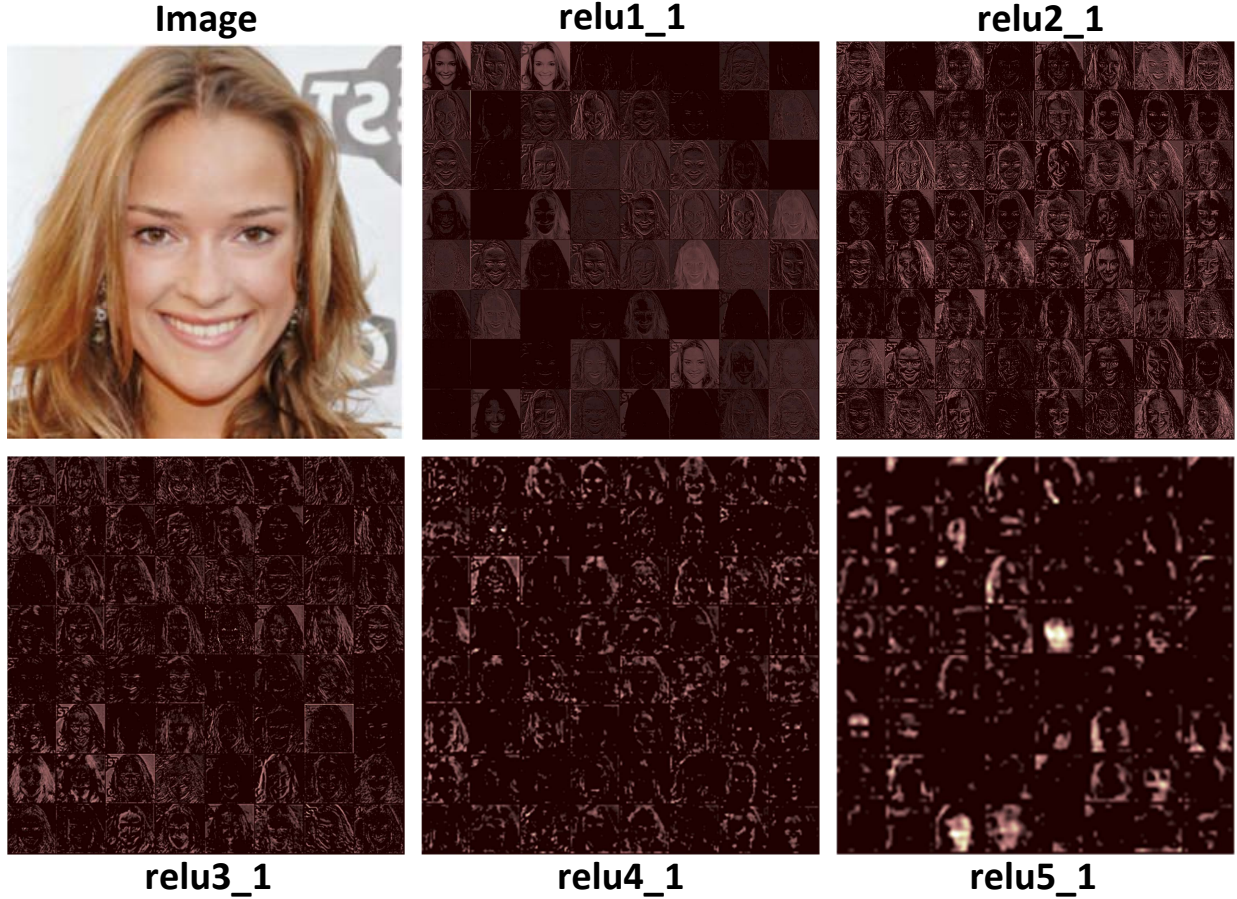


Fig. 16. Visualization of VGG feature maps (*the first* 64 *pieces*).

'relu1_1', 'relu2_1', and 'relu3_1' layers almost have the completed face contour, which is unsuited to aligning the part components using geometrical alignment constraint. For the spatial resolution of each response map generated by 'relu5_1' layer is only $16 \times 16$, it will result in a small coordinate range. Thus, we choose the output response maps of 'relu4_1' layer to compute our geometrical alignment constraint, which guides the coordinate expectation registration.

## V. CONCLUSION

In this paper, we proposed a dense multi-scale fusion network with self-guided regression loss and geometrical alignment constraint for image fine-grained inpainting, which highly improves the quality of produced images. Specifically,

dense multi-scale fusion block is developed to extracted better features. With the assistance of self-guided regression loss, the restoration of semantic structures becomes easier. Additionally, geometrical alignment constraint is inductive to the coordinate registration between generated image and ground-truth, which promotes the reasonableness of painted results.

TABLE V

THE COMPARISON OF SELF-GUIDED REGRESSION LOSS WITH VARIOUS DISTANCE METRICS ON PARIS STREETVIEW. HERE, $X_i \in \mathbb{R}^{C \times 1}$ REPRESENTS THE VECTOR OF THE IMAGE $X \in \mathbb{R}^{C \times H \times W}$ AT POSITION $i$.

| Distance metric | $\phi(X_i, Y_i)$ | PNSR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|---|
| Gaussian distance | $\exp\left\{-\|X_i - Y_i\|_2^2 / 2\sigma^2\right\}$ | **25.06** | **0.8596** | 0.1027 |
| Dot product | $X_i Y_i^T$ | 24.77 | 0.8588 | 0.1035 |
| L2 distance | $\|X_i - Y_i\|_2^2$ | 25.00 | 0.8563 | **0.1018** |

## REFERENCES

[1] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patch-match: a randomized correspondence algorithm for structural image editing," *ACM Transactions on Graphics (TOG)*, vol. 28, no. 3, pp. 24:1–24:11, 2009.

[2] 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 (NeurIPS)*, 2014, pp. 2672–2680.

[3] D. Pathak, P. Krähenbühl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPP)*, 2016, pp. 2536–2544.

[4] C. Yang, X. Lu, Z. Lin, E. Shechtman, O. Wang, and H. Li, "High-resolution image inpainting using multi-scale neural patch synthesis," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6721–6729.

[5] Y. Li, S. Liu, J. Yang, and M.-H. Yang, "Generative face completion," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3911–3919.

[6] S. Iizuka, E. Simo-Serra, and H. Ishikawa, "Globally and locally consistent image completion," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 107:1–107:14, 2017.

[7] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative image inpainting with contextual attention," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5505–5514.

[8] Y. Wang, X. Tao, X. Qi, X. Shen, and J. Jia, "Image inpainting via generative multi-column convolutional neural networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 331–340.

[9] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein gans," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5767–5777.

[10] Y. Song, C. Yang, Z. Lin, X. Liu, Q. Huang, H. Li, and C.-C. Jay Kuo, "Contextual-based image inpainting: Infer, match, and translate," in *European Conference on Computer Vision (ECCV)*, 2018, pp. 3–19.

[11] Z. Yan, X. Li, M. Li, W. Zuo, and S. Shan, "Shift-net: Image inpainting via deep feature rearrangement," in *European Conference on Computer Vision (ECCV)*, 2018, pp. 1–17.

[12] C. Zheng, T.-J. Cham, and J. Cai, "Pluralistic image completion," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1438–1447.

[13] Y. Zeng, J. Fu, H. Chao, and B. Guo, "Learning pyramid-context encoder network for high-quality image inpainting," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1486–1494.

[14] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4401–4410.

[15] A. Jolicoeur-Martineau, "The relativistic discriminator: a key element missing from standard gan," in *International Conference for Learning Representations (ICLR)*, 2019.

[16] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. Huang, "Free-form image inpainting with gated convolution," in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 4471–4480.

[17] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester, "Image in-painting," in *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, 2000, pp. 417–424.

[18] C. Ballester, M. Bertalmio, V. Caselles, G. Sapiro, and J. Verdera, "Filling-in by joint interpolation of vector fields and gray levels," *IEEE Transactions on Image Processing (TIP)*, vol. 10, no. 8, pp. 1200–1211, 2001.

[19] A. A. Efros and W. T. Freeman, "Image quilting for texture synthesis and transfer," in *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques*, 2001, pp. 341–346.

[20] V. Kwatra, I. Essa, A. Bobick, and N. Kwatra, "Texture optimization for example-based synthesis," *ACM Transactions on Graphics (TOG)*, vol. 24, no. 3, pp. 795–802, 2005.

[21] D. Simakov, Y. Caspi, E. Shechtman, and M. Irani, "Summarizing visual data using bidirectional similarity," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008, pp. 1–8.

[22] D. Ding, S. Ram, and J. J. Rodríguez, "Image inpainting using nonlocal texture matching and nonlinear filtering," *IEEE Transactions on Image Processing (TIP)*, vol. 28, no. 4, pp. 1705–1719, 2019.

[23] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference for Learning Representations (ICLR)*, 2015.

[24] 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 *European Conference on Computer Vision (ECCV)*, 2018, pp. 85–100.

[25] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1125–1134.

[26] K. Nazeri, E. Ng, T. Joseph, F. Z. Qureshi, and M. Ebrahimi, "Edgeconnect: Structure guided image inpainting using edge prediction," in *IEEE International Conference on Computer Vision Workshop (ICCVW)*, 2019.

[27] J. Li, F. He, L. Zhang, B. Du, and D. Tao, "Progressive reconstruction of visual structure for image inpainting," in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 5962–5971.

[28] W. Xiong, Z. Lin, J. Yang, X. Lu, C. Barnes, and J. Luo, "Foreground-aware image inpainting," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5840–5848.

[29] Y. Ren, X. Yu, R. Zhang, T. H. Li, S. Liu, and G. Li, "Structureflow: Image inpainting via structure-aware appearance flow," in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 181–190.

[30] M. cheol Sagong, Y. goo Shin, S. wook Kim, S. Park, and S. jea Ko, "Pepsi: Fast image inpainting with parallel decoding network," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 11 360–11 368.

[31] Z. Hui, J. Li, X. Gao, and X. Wang, "Progressive perception-oriented network for single image super-resolution," *Information Sciences*, vol. 546, pp. 769–786, 2021.

[32] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "Esrgan: Enhanced super-resolution generative adversarial networks," in *European Conference on Computer Vision Workshop (ECCVW)*, 2018, pp. 63–79.

[33] Y. Zhou, Z. Zhu, X. Bai, D. Lischinski, D. Cohen-Or, and H. Huang, "Non-stationary texture synthesis by adversarial expansion," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 4, pp. 49:1–49:13, 2018.

[34] C. Wang, C. Xu, C. Wang, and D. Tao, "Perceptual adversarial networks for image-to-image transformation," *IEEE Transactions on Image Processing (TIP)*, vol. 27, no. 8, pp. 4066–4079, 2018.

[35] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 40, no. 6, pp. 1452–1464, 2017.

[36] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," in *International Conference for Learning Representations (ICLR)*, 2018.

[37] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *International Conference for Learning Representations (ICLR)*, 2015.

[38] C. Xie, S. Liu, C. Li, M.-M. Cheng, W. Zuo, X. Liu, S. Wen, and E. Ding, "Image inpainting with learnable bidirectional attention maps," in *IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 8858–8867.

[39] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in

*IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 586–595.