# Image Inpainting with Contextual Reconstruction Loss

Yu Zeng[1]    Zhe Lin[2]    Huchuan Lu[3]    Vishal M. Patel[4]

zengxianyu18@qq.com, zlin@adobe.com, lhchuan@dlut.edu.cn, vpatel36@jhu.edu

[1]Tencent Lightspeed & Quantum Studios    [2]Adobe Research
[3]Dalian University of Technology    [4]Johns Hopkins University

## Abstract

*Convolutional neural networks (CNNs) have been observed to be inefficient in propagating information across distant spatial positions in images. Recent studies in image inpainting attempt to overcome this issue by explicitly searching reference regions throughout the entire image to fill the features from reference regions in the missing regions. This operation can be implemented as contextual attention layer (CA layer) [24], which has been widely used in many deep learning-based methods. However, it brings significant computational overhead as it computes the pair-wise similarity of feature patches at every spatial position. Also, it often fails to find proper reference regions due to the lack of supervision in terms of the correspondence between missing regions and known regions. We propose a novel contextual reconstruction loss (CR loss) to solve these problems. First, a criterion of searching reference region is designed based on minimizing reconstruction and adversarial losses corresponding to the searched reference and the ground-truth image. Second, unlike previous approaches which integrate the computationally heavy patch searching and replacement operation in the inpainting model, CR loss encourages a vanilla CNN to simulate this behavior during training, thus no extra computations are required during inference. Experimental results demonstrate that the proposed inpainting model with the CR loss compares favourably against the state-of-the-arts in terms of quantitative and visual performance. Code is available at https://github.com/zengxianyu/crfill.*

## 1 Introduction

Image inpainting is a task of filling missing regions in images. It is an important problem in computer vision and can be used in many applications, *e.g.* image restoration, compositing, manipulation, re-targeting, and image-based rendering [2, 10, 15]. Traditional methods such as [4, 9, 2] borrow example patches from known regions or external dataset to paste into the missing regions. They cannot hallu-
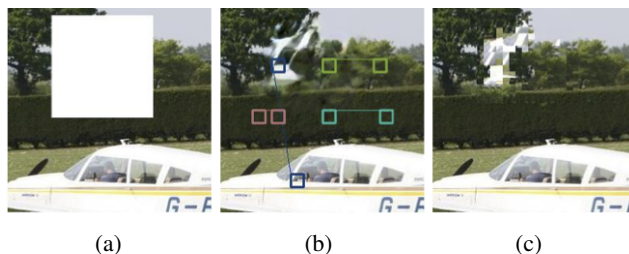


Figure 1: Artifacts in the inpainting results can often find their origin from unsuitable reference regions. (a) Input. (b) Inpainting results of Deepfillv2 and the reference patches selected by the CA layer. (c) Images produced by moving each reference patch to the corresponding location in the missing region.

cinate novel image contents for challenging cases involving complex, non-repetitive structures. Recent research efforts have shifted the attention to data-driven deep CNN-based approaches [16, 6, 24, 12, 25]. The generative adversarial framework [5] is often adopted due to its success in generating visually realistic images. Significant amounts of training data can be obtained simply by corrupting images deliberately and using the original images before corruption as the ground-truths. By training a deep model on such a large dataset, one can generate plausible inpainting results for complex scenes. However, it has been observed that CNNs often fail to effectively model long-term correlations between distant contextual information and the missing regions which leads to distorted structures and blurry textures inconsistent with surrounding known regions [24].

Recently, contextual attention layer (CA layer) [24] was proposed to alleviate this problem. It is inserted between convolution layers to replace feature patches in the missing regions with a weighted sum of the contextual patches. It allows the network to explicitly borrow known patches as references to better predict the missing region. Since the reference patches are searched throughout the entire image, long-term correlations can be captured. Also, since each patch is borrowed as a whole, its local structure is kept, resulting in less-distorted structures and better textures. CA layer and its variants have been widely adopted in many

deep learning-based inpainting methods [25, 14, 26, 27, 20]. However, there are some important issues that prevent CA layer from being implemented in practical systems. First, it requires to compute the similarity of every pair of patches in the input image, which is computationally expensive especially for high-resolution images. The computational complexity of CA layer is $O(N^2)$, where $N$ is the number of pixels in the input image and is the real bottleneck with growing resolution. Second, there is no supervision on the searching of the reference patches for CA layer. Sometimes it chooses inappropriate reference patches, which leads to artifacts in the output image. Consider the image shown in Fig. 1 (a) with missing pixels. Fig. 1 (b) shows an inpainting result of Deepfillv2 [25] and the reference patches selected by CA layer. We can see that some reference patches are not appropriate. For clarity, we move each reference patch to the corresponding location in the missing region and produce Fig. 1 (c). By comparing Fig. 1 (b) and Fig. 1 (c), we can see how the artifacts have resulted from the incorrect reference regions. For instance, the white stripes on the tree are generated by the patches from the car.

In this paper, we propose an effective yet more efficient alternative to CA layer to address these issues. First, we design a criterion of searching reference patches in known regions: the reference patches themselves should have the smallest inpainting loss to the ground-truth at the corresponding positions. In other words, moving the reference patches to their corresponding positions in the missing region should result in a plausible image. Second, we propose a contextual reconstruction loss (CR loss) to guide the inpainting network to utilize the reference patches subjected to this criterion. Unlike CA layer that explicitly puts the computational load by searching and replacing patches in the inpainting model during inference stage, the CR loss matches the generated patches with the reference patches only during training. More specifically, during training, the pair-wise similarity of generator feature patches are computed, and then an auxiliary image is produced by reconstructing the generator output (i. e. the inpainting result) using the weighted sum of known contextual patches with softmax of this pair-wise similarity as the weight. This resembles the way we have produced Fig. 1 as described above but is made differentiable by exploiting softmax and an autoencoder-like auxiliary network. By minimizing the inpainting loss of the auxiliary image to the ground-truth, the generator features will be encouraged to closer to features of the image patches of the smallest inpainting loss.

Since CR loss is required only during training, there is no overhead brought to the inpainting model during inference. As a result, the obtained model is very efficient. The contributions of this paper are summarized as follow:

- A novel CR loss is proposed for training a deep network for image inpainting. CR loss is general and can

be applied to different CNN-based inpainting models without changing their architectures.
- Similar to CA layer, CR loss enables a network explicitly to borrow known regions as references to fill in a hole. Different from CA layer, the CR loss brings no extra computational burden to the inpainting model and the searching of reference known regions is subject to a well-defined criterion.
- Experiments show the effectiveness of the CR loss and the performance of our inpainting method compares favorably against the state-of-the-art methods.

## 2  Related work

Earlier inpainting methods rely on the principle of borrowing known regions to fill missing regions. Diffusion-based methods [1, 3] propagate neighboring content to the missing regions and often result in significant artifacts when filling large holes or when texture is of large variation. Patch-based methods [4, 9, 2] search for the most similar patches from known regions to complete missing regions. They can produce high-quality results for textures and repeating patterns. However, due to the lack of high-level structural understanding and inability of generating new content, the results may not be semantically reasonable.

Encouraged by the success of deep CNNs in image restoration tasks, recent research efforts have shifted their attention to deep learning-based methods. To produce sharper results, these methods typically adopt adversarial training inspired by GANs [5]. Pathak *et al.* [16] first attempted to use a CNN for hole filling. Li *et al.* [11] propose a deep generative model for face completion. Iizuka *et al.* [6] use two discriminators for adversarial training to make the inpainted content both locally and globally consistent. By learning from a large corpus of data, deep CNN-based methods can understand image structure and hence can handle more difficult cases. However, inpainting general scenes still remains a difficult task for deep CNN-based methods. CNN-based methods often suffer from the problem of generating distorted structures and artifacts. This is possibly due to their ineffectiveness in exploiting structured image features from distant known regions [24].

Inspired by classical inpainting methods that borrow known patches to fill missing regions, analogous operations have been integrated into deep learning models. Yu *et al.* [24] propose a contextual attention layer (CA layer) which replaces the generated features in the missing region with linear combinations of feature patches from known region using similarity as weight. Zeng *et al.* [26] propose to use region affinity from a high-level feature map to guide the patch replacement operation in the previous low-level feature map. Yang *et al.* [22] propose a multi-scale neural patch synthesis approach based on joint optimization of image content and texture constraints. Zeng *et al.* [27] use a re-

lated neural patch-vote approach to upsample but avoids the slow optimization by using a modified contextual attention layer. Yan *et al.* [21] shift the encoder features of the known region to the missing region in the mirrored layer of the decoder serving as an estimation of the missing parts. Song *et al.* [17] propose a patch-swap layer that replaces each feature patch in the missing regions with the most similar patch on the known regions. Beside these methods, related variants of patch replacement operations have also been used in other image inpainting methods [14, 13, 25, 20].

Contextual attention in image inpainting models are related to self-attention [18] and non-local layers [19]. The response of self-attention and non-local layers for each element in a feature map is a weighted sum using the affinity of elements as weight. This is very similar to contextual attention or patch replacement operations. The main difference is that contextual attention measures similarity and construct response at patch level, while self-attention and non-local layers operate on a per-element basis.

All of the above mentioned methods including CA layer and related patch replacement operations for image inpainting as well as self-attention and non-local layers for other tasks all require to compute the similarity of every pair of elements of a feature map during inference, which is not very efficient especially with growing resolution. Different from them, our proposed method uses a novel contextual reconstruction loss to encourage a vanilla CNN to simulate these operations. To the best of our knowledge, this is the first attempt to implicitly model the patch replacement operation in CNNs guided by a loss.

## 3 Generative inpainting network

We adopt generative adversarial network-based approach for image inpainting, which has a generator and a discriminator. The objective of the discriminator is to discriminate between real images (without any missing pixels) and images inpainted by the generator.

**Discriminator.** We use a PatchGAN discriminator [7] with spectral normalization following [25]. Inpainted images or the ground-truth complete images are passed to the discriminator, and it outputs a score map, where each element is a score corresponding to a local region of the input covered by its receptive field. The loss for the discriminator is:

$$\mathcal{L}_D = \mathbb{E}_{X \sim p_{data}(X)} \big[ \mathrm{ReLU}(\mathbb{1} - D(X)) \big] + \\ \mathbb{E}_{U \sim p_U(U)} \big[ \mathrm{ReLU}(\mathbb{1} + D(G(U) \circ M + U)) \big], \quad (1)$$

where $D$ denotes the discriminator, $X$ represents the real image (ground-truth), $U$ represents the incomplete image with the pixels in the missing regions set to zero, $M$ represents a binary mask corresponding of the missing region where $M_{xy} = 1$ indicates that pixel at $x, y$ is missing and $M_{xy} = 0$ indicates that pixel at $x, y$ is valid/known, $G(\cdot)$ represents the generator and $\circ$ denotes element-wise multiplication. The inpainting result $G(U) \circ M + U$ is composed

by putting the generated content $G(U)$ in the missing region and keeping the original content of $U$ in the known region.

**Coarse-to-fine generator.** Fig. 2 shows the overall architecture of our generator network. It is a coarse-to-fine architecture, similar to the one in DeepFillv2 [25] but the main difference is that the CA layer is removed and CR loss is applied instead. The coarse network and the refinement network are convolutional encoder-decoder type networks. Dilated convolution layers are inserted to enlarge the receptive fields. We use gated convolution [25] in all convolution and dilated convolution layers. The coarse network takes an incomplete image where missing pixels are set to zero and a binary mask indicating the missing region as input and generates an initial prediction. Then the refinement network takes this initial prediction as input and outputs the final inpainting result.
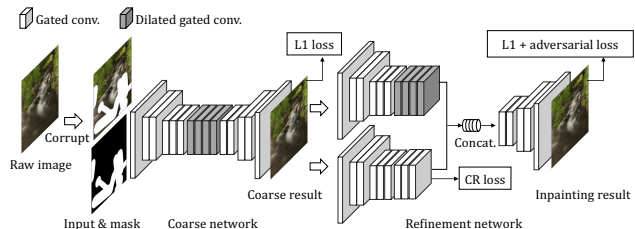


Figure 2: Overall architecture of the generator network.

**Training the generator.** We expect the coarse network to complete the global structure and then details will be filled in by the refinement network. So only the L1 loss is used to train the coarse network. On the other hand, a combination of L1 loss, adversarial loss and the proposed CR loss is used to train the refinement network. Let $Y = G(U)$ represent the refinement network output. Then the loss for the refinement network is defined as follow:

$$\mathcal{L}_G = \mathbb{E}_{U,X \sim p(U,X)}[L(Y) + \lambda L_{CR}], \quad (2)$$

where $L(Y)$ is the sum of L1 loss and adversarial loss

$$L(Y) = \mathrm{ReLU}(\mathbb{1} - D(Y \circ M + U)) + \beta \|Y - X\|_1 . \quad (3)$$

Here, $L_{CR}$ denotes CR loss, which will be elaborated in Section 4. $\lambda$ is set equal to 0.5. $\beta$ is set equal to 1.5.

## 4 Contextual reconstruction loss

**Revisit contextual attention layer.** CA layer borrows patches from the known region to fill in the missing region. It can be inserted between convolution layers and applied on the feature maps. First, a feature map is split into foreground patches and background patches which correspond to the missing region and known region in the input image, respectively. Then the similarity of each pair of foreground and background patches are measured and every foreground patch is replaced by the linear combination of background patches using softmax of the similarity as weight. There are two problems with contextual attention layer. First, in approximation, each foreground patch is replaced with its

nearest neighbor among the known patches, not necessarily a proper reference for the corresponding location at the missing region. Second, it requires to compute the similarity of each pair of foreground and background patches at the inference stage, which brings heavy computational overhead, especially for high-resolution images.

**Contextual reconstruction loss.** To overcome the above issues, we propose the contextual reconstruction loss. Fig. 3 compares the usage of CA layer and CR loss in a CNN. Unlike CA layer that is inserted in the network, CR loss does not directly involve in the generation of inpainted image and affects the network only during training for learning better features. The overall training system is shown



Figure 3: Comparison of the usage of CA layer and CR loss.

in Fig. 5, where the inpainting model takes an incomplete image as well as the binary mask of the missing regions as input and outputs the inpainting result, which has been defined in Sec. 3. The training system with CR loss consists of a similarity encoder and an auxiliary encoder-decoder type network and re-uses the inpainting loss defined in Sec. 3. The similarity encoder takes the generator feature as input and encodes the similarity among image regions. The auxiliary encoder-decoder network produce an auxiliary image in which the known regions are unchanged while the missing regions are filled with similar known regions based on the similarity provided by the similarity encoder. CR loss of the generator feature is defined to be the inpainting loss ( *i. e.* L1 and adversarial loss) of the auxiliary image. By minimizing CR loss, the generator features are encouraged to be close to the known image features of the smallest inpainting loss. In what follows, we elaborate on the formal definition and explanation for this loss.

For the convenience of the description of the proposed method, here we define the global and patch-wise perspective of a CNN. If we view an image $U$ passed through a CNN as the combination of square patches $u_1, u_2, ...$, a convolution layer feature map $F(U)$ of this image $U$ can be seen as the combination of local function of image patches $f(u_1), f(u_2), ...$. Hereafter an upper case letter, *e. g.* $F(U)$, will represent a whole image or feature map with the corresponding lower case one for a patch in it *e. g.* $f(u_i)$. Each image patch can correspond to feature patches of various sizes in different layers of a CNN. For example, for the layers shown in Fig. 4, an $N \times N$ image patch $u$ centered at $(x, y)$ corresponds to an $\frac{N}{2} \times \frac{N}{2}$ feature patch $f_2(u)$ centered at $(\frac{x}{2}, \frac{y}{2})$ in the second layer and an $\frac{N}{4} \times \frac{N}{4}$ feature patch $f_3(u)$ centered at $(\frac{x}{4}, \frac{y}{4})$ in the third layer. They can be seen as different representations of the same patch. There-
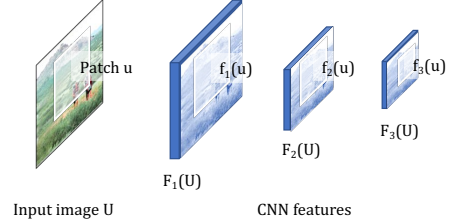


Figure 4: If view an input image as combination of patches, corresponding feature patch in a layer $l$ can be seen as a representation $f_l(u)$ of an image patch $u$.
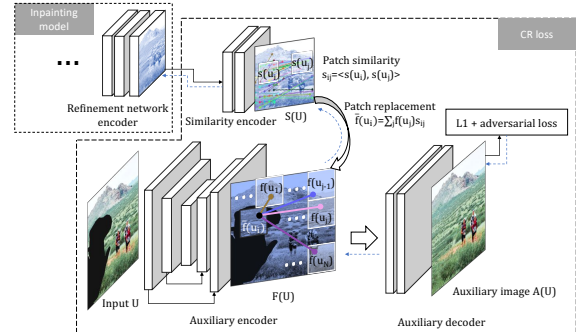


Figure 5: Training system with the proposed contextual reconstruction loss. The dashed blue lines indicates the gradient back-propagation flow.

fore, we can take the similarity of a pair of feature patches $f_l(u_i), f_l(u_j)$ at a layer $l$ and use it as a similarity $s_{ij}$ of the corresponding image patch or feature patches at arbitrary other layers. We take the pair-wise cosine similarity of feature patches $< s(u_i), s(u_j) >$ at the last layer of the similarity encoder as $s_{ij}$:

$$s_{ij} = \frac{s(u_i)^\mathsf{T} s(u_j)}{\|s(u_i)\| \cdot \|s(u_j)\|}, \quad (4)$$

where the feature patches $s(u_i), s(u_j)$ are viewed as vectors when taken the inner product. It can be implemented as processing the feature map with patches extracted from itself as convolution filters as described in [24].

Before the feature of the auxiliary encoder is passed through the auxiliary decoder, each feature patch is replaced with a weighted sum of patches in the known region with the softmax of the similarity provided by the similarity encoder as weight. Let $f(u_i)$ be the auxiliary encoder feature of image patch $u_i$, $\bar{f}(u_i)$ the feature of $u_i$ after patch replacement is obtained as follows,

$$\bar{f}(u_i) = \sum_{j \in \mathcal{V}} \text{softmax}(\alpha s_{ij}) f(u_j), \quad (5)$$

where $\mathcal{V}$ represents the index set of patches in the known region. $\alpha$ is set equal to 10. It can be implemented as transposed convolution which processes the similarity maps with patches extracted from the auxiliary encoder features as filters after dropout with the similarity to all patches in the missing region.

Then the feature map after the patch replacement is

4

transformed into an auxiliary image by the auxiliary decoder. Let $\bar{F}(U)$ be the feature map consisting of feature patches $\bar{f}(u_i)$ given in Eqn. 5, the auxiliary image $A(U)$ of the input image $U$ is obtained as follow,

$$A(U) = H(\bar{F}(U)), \qquad (6)$$

where $H()$ represents the auxiliary decoder connected to the auxiliary encoder. Based on this discussion, the contextual reconstruction loss $L_{CR}$ is defined as the inpainting loss $L()$ (Eqn. 3) of the auxiliary image:

$$L_{CR} = L(A(U)). \qquad (7)$$

Assumes that the auxiliary decoder inverts the auxiliary decoder for known regions. Then the auxiliary encoder-decoder approximately conducts the process of producing Fig. 1 (c) described in Sec. 1: for each patch $u_i$ in the missing regions, the most similar patch $u_{i*}$ from the known regions, *i. e.* $i^* = \arg\max_j s_{ij}$, is moved to the position of $u_i$. This results in a simplified view of the training system shown in Fig. 6, which is an analogue of learning to solve a jigsaw puzzle. To minimize the loss of the auxiliary image, proper known patches should be chosen and moved to the right place. Furthermore, as the reference known patches are selected according to the similarity of generator features, it requires the generator feature in the missing region to be closest to the proper known region. Since the searching of known patches is throughout the whole image, distant relationship can be captured. A more detailed explanation will be given in the following section.
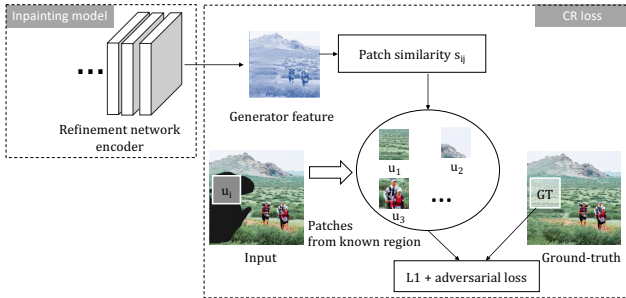


Figure 6: A simplified view of the training system.

**Explanation.** For a patch in the known region, *i. e.* $u_j, j \in \mathcal{V}$, since $\bar{f}(u_j) \approx f(u_j)$ as softmax$(\alpha s_{jj}) \approx 1$, the auxiliary image patch in the known region will be $a(u_j) = h(f(u_j)), j \in \mathcal{V}$. Note that the ground-truth of $a(u_j)$ is $u_j$ itself for $j \in \mathcal{V}$, the auxiliary encoder-decoder $H(F(U))$ will learn lazily to copy the input to the output so it is easy to make the auxiliary decoder invert the auxiliary encoder: $h(f(u_j)) = u_j$ for $j \in \mathcal{V}$. To make it easier, we put skip connection from shallow layers to the deeper layers of the auxiliary encoder, as indicated in Fig. 5.

Assume that softmax in Eqn. 5 to be "hard", then the feature patch after replacement can be approximated as:

$$\bar{f}(u_i) = f(u_{i*}) \text{ where } i^* = \arg\max_{j \in \mathcal{V}} s_{ij}. \qquad (8)$$

Then a patch in the auxiliary image is:

$$a(u_i) = h(f(u_{i*})) = u_{i*} \text{ where } i^* = \arg\max_{j \in \mathcal{V}} s_{ij}. \qquad (9)$$

Resembling the translation between the global and patch-wise perspective of a CNN introduced earlier, the inpainting loss for the whole image also can be distributed to each patch. Thereby the inpainting loss of the auxiliary image $L(A(U))$ can be seen as a sum over local loss of patches: $L(A(U)) = \sum_i l_i(a(u_i))$ (We differentiate the local loss of different patches as the ground-truth for each patch is different, which means $L(A(U))$ is in fact $L(A(U); X) = \sum_i l(a(u_i); x_i)$ and it is denoted as $l_i(a(u_i))$ for short). We have argued earlier that the loss on known regions can be minimized with the auxiliary encoder-decoder copying the input to the output. So what count for $L(A(U))$ are the patches in the missing region:

$$L(A(U)) = \sum_{i \in \mathcal{V}^C} l_i(u_{i*}) \text{ where } i^* = \arg\max_{j \in \mathcal{V}} s_{ij}. \qquad (10)$$

This is a sum over local inpainting loss of known image patches with the largest similarity to the generated patches in the missing region. Note that the number of known patches in an image is limited, so for each local inpainting loss $l_i()$ there must be one patch $u_{i^0}$ that has the smallest loss among all the known patches, *i. e.* $i^0 = \arg\min_{j \in \mathcal{V}} l_i(u_j)$. The minimum of CR loss is achieved when $i^* = i^0$. In other words, when the generator features in missing regions are closest to the features in known regions which have the smallest inpainting loss.

In summary, to minimize the CR loss, two conditions are required. First, appropriate reference patches are selected from the known region. Second, feature patches in missing regions are closest to the corresponding reference patches. Compared with CA layer which searches the known patches in a heuristic way with no guarantee on the reference patches it reaches, CR loss can find more reasonable reference patches. This will be supported by quantitative experimental results in Sec. 5.

## 5 Experimental Results

We implement our method with Python and PyTorch. Detailed network architectures and code can be found in the supplementary material. We train the models using Adam [8] optimizer with the learning rate of 0.0001. To show the effectiveness of the CR loss, we train the network shown in Fig. 2 with and without the CR loss on the Places2 training set. We add images from the salient object segmentation dataset [20] following [27, 20]. Square masks, irregular masks [25] and object-shaped masks [27] are randomly switched at every mini-batch to be used for creating missing regions. Training samples are cropped to $256 \times 256$ and missing regions are put at random positions. The obtained models are denoted as Baseline and Baseline+CR, respectively. To compare the proposed CR loss with CA layer, we also train a network with the same architecture as

DeepFillv2 [25] but remove the CA layer and apply CR loss instead. The training protocol for this network is made to be the same as the official implementation of DeepFillv2 [25], where Places2 training set with square and irregular masks are used for training; $\beta$ in Eqn. 3 is set equal to 1 as in the official implementation of DeepFillv2. The obtained model is denoted as DeepFillv2-CA+CR. All evaluations are conducted on the same platform: an Ubuntu machine with a 3GHz Intel i7-9700F CPU, 32GB memory, 256GB swap space, and a GPU NVIDIA RTX2080Super with 8GB GPU memory.

### 5.1 Comparison with state-of-the-arts

We compare our methods with the following state-of-the-art methods: PENNet [26], Deepfillv2 [25], Rethink [13], Hi-Fill [23]. We use their official implementation and the models trained on Places2 training set provided by the authors of corresponding papers. We use L1 error, PSNR and SSIM to measure the performance quantitatively. For comparison on visual quality, we show the inpainting results and the results of a subject evaluation conducted by human raters.

**Quantitative evaluation.** Table 1 and Table 2 show quantitative comparisons of our method with state-of-the-art methods on Places2 validation set containing 36500 images and ImageNet validation set containing 50000 images. All images are cropped to $256 \times 256$ with missing regions of different shapes ($128 \times 128$ squares, object and irregular shapes) at random positions. From these table, we can see that our method has a smaller L1 error, larger PSNR and SSIM than the existing methods in most cases, which indicates that the inpainting results corresponding to our method are closer to the ground-truth.

**Visual quality.** Fig. 7 shows the visual comparison of our method and existing methods. We can see from the figure that the missing region recovered by our method is more visually coherent with surrounding known regions. This figure implies the effectiveness of CR loss in exploiting structured information in known regions. Furthermore, unlike a previous method PENNet which has high performance in quantitative evaluation but tends to generate blurry images, the results of our method are more visually realistic.

**User study.** We invite 9 human raters to conduct a user study on 45 images randomly sampled from the Places2 validation set. Every $1/3$ of them are corrupted with missing regions of different shapes, *i. e.* square, irregular and object shaped, at random positions. Each time the incomplete image and the results by different methods are presented to the raters in random order. The raters are asked to select one best result. The number of user preference for all methods are shown in the last column of Table 1. From the user study results, we can see that the results of our method are preferred by the human raters most frequently.

**High-resolution inpainting.** The efficiency of the proposed method makes it potentially capable of processing high-resolution input. We make a simple adjustment to the network to adapt to high-resolution input: running the coarse network at $1/2$ the input size and interpolate its output back to the original resolution before passed through the refinement network. Fig. 8 shows high-resolution inpainting results by our adjusted inpainting model and HiFill [23] which was proposed for inpainting at high-resolution. More results and comparison at high-resolution can be found in the supplementary material[1].

### 5.2 Ablation study

**Inpainting results.** Fig. 10 shows a visual comparison of the results obtained by the baseline model (Baseline), the baseline model trained with the CR loss (Baseline+CR), and the results of DeepFillv2 with CA layer. It can be seen that both CA layer and CR loss enable the inpainting model to capture a long term relationship among image regions and the distant known regions while while the baseline model tends to be short-sighted which propagates the most nearby features. For instance, in the second example, both the model with CA layer and CR loss correctly capture the grid pattern while the baseline model just extends the lines.

The first two rows of Table 4 show the quantitative comparison of the baseline model (Baseline) and the baseline model with the CR loss (Baseline+CR). The last two rows show the comparison of the original DeepFillv2 and the model obtained by removing CA layer and applying the CR loss (DeepFillv2-CA+CR) with network architectures and training protocols unchanged. From this table, we can see that the CR loss can improve the quantitative performance of different models, and DeepFillv2 with CR loss yields better results than the official implementation with CA layer.

**Reference patch.** Fig. 9 shows the reference regions used by the CR Loss and the reference region used by CA layer. We can see that in these examples the network trained with the CR loss finds more reasonable reference regions than the network with CA layer. By comparing Fig. 9 (c) and (e) we can see that the reference patches selected by the CR loss better fit the surrounding known regions. To quantify the difference, we measure the average L1 error, PSNR and SSIM of the images constructed by moving each reference patch to the corresponding position in the missing region like piecing together a jigsaw. The quantitative comparison is shown in Table 5, from which we can see that the images pieced by the CR loss have a smaller L1 error, larger PSNR and SSIM, which indicates that the CR loss can find the reference regions closer to the ground-truth.

**Efficiency.** Since CR loss gets rid of the heavy pair-wise similarity computation in the inference stage, the obtained model is computationally efficient. Table 5 shows time complexity and the running time of the networks of the same architecture with the CR loss and CA layer at different

---

[1] https://github.com/zengxianyu/crfill

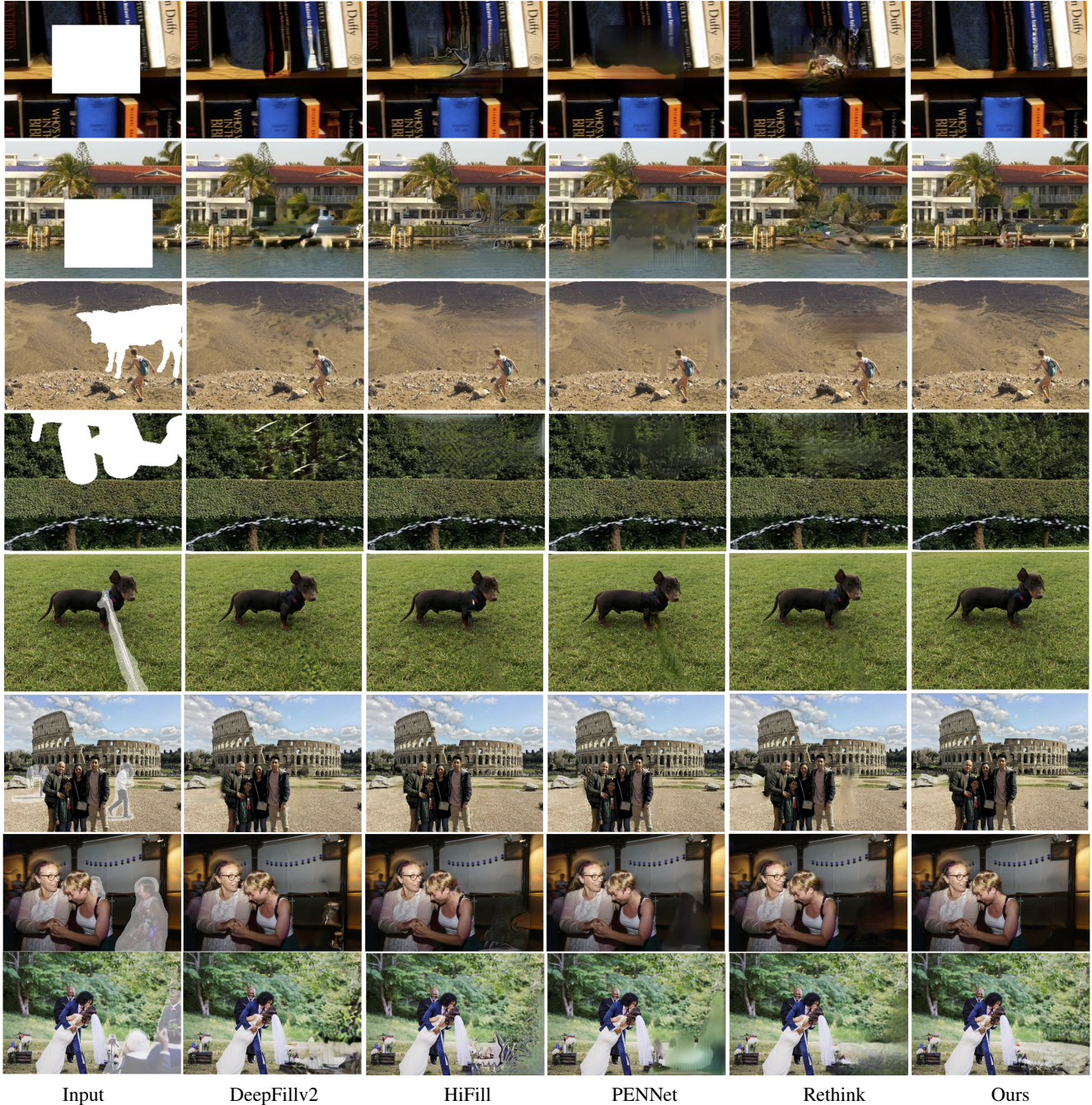| Input | DeepFillv2 | HiFill | PENNet | Rethink | Ours |

Figure 7: Visual comparison of inpainting results of our method and state-of-the-arts. Zoom-in to see the details. Images are compressed due to file size limitation. More results and comparison at high-resolution can be found in the supplementary material.

Table 1: Quantitative evaluation results on the Places2 validation set. The best scores are in bold.

| Method | Square holes | | | Irregular holes | | | Object holes | | | User preference |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | L1 error | PSNR | SSIM | L1 error | PSNR | SSIM | L1 error | PSNR | SSIM | |
| DeepFillv2 [25]ICCV'19 | 0.0290 | 22.72 | 0.8429 | 0.0178 | 27.03 | 0.8947 | 0.0225 | 27.95 | 0.8887 | 88 |
| Rethink [13]ECCV'20 | 0.0307 | 22.30 | 0.8305 | 0.0194 | 26.13 | 0.8796 | 0.0237 | 27.07 | 0.8829 | 13 |
| PENNet [26]ICCV'19 | 0.0284 | 23.28 | 0.8402 | 0.0227 | 25.35 | 0.8699 | 0.0251 | 26.28 | 0.8780 | 6 |
| HiFill [23]CVPR'20 | 0.0329 | 21.35 | 0.8103 | 0.0210 | 25.48 | 0.8712 | 0.0267 | 26.14 | 0.8647 | 27 |
| Ours | **0.0263** | **23.36** | **0.8462** | **0.0157** | **27.79** | **0.9002** | **0.0202** | **28.79** | **0.8929** | **271** |

| Input | HiFill | Ours | Input | HiFill | Ours |

Figure 8: High-resolution results ($1200 \times 1600, 1152 \times 1536$) compared with HiFill. Zoom in to see the details. Images are compressed due to file size limitation.

Table 2: Quantitative evaluation results on the ImageNet validation set. The best scores are in bold.

| Method | Square holes | | | Irregular holes | | | Object holes | | |
|---|---|---|---|---|---|---|---|---|---|
| | L1 error | PSNR | SSIM | L1 error | PSNR | SSIM | L1 error | PSNR | SSIM |
| DeepFillv2 [25]ICCV'19 | 0.0335 | 21.18 | 0.8200 | 0.0205 | 25.55 | 0.8765 | 0.0250 | 26.51 | 0.8739 |
| Rethink [13]ECCV'20 | 0.0347 | 21.13 | 0.8137 | 0.0216 | 25.15 | 0.8682 | 0.0259 | 26.11 | 0.8723 |
| PENNet [26]ICCV'19 | 0.0316 | **22.09** | 0.8257 | 0.0243 | 24.53 | 0.8615 | 0.0269 | 25.48 | 0.8698 |
| HiFill [23]CVPR'20 | 0.0329 | 21.35 | 0.8103 | 0.0210 | 25.48 | 0.8712 | 0.0287 | 25.24 | 0.8570 |
| Ours | **0.0302** | 21.91 | **0.8269** | **0.0180** | **26.43** | **0.8855** | **0.0225** | **27.45** | **0.8808** |



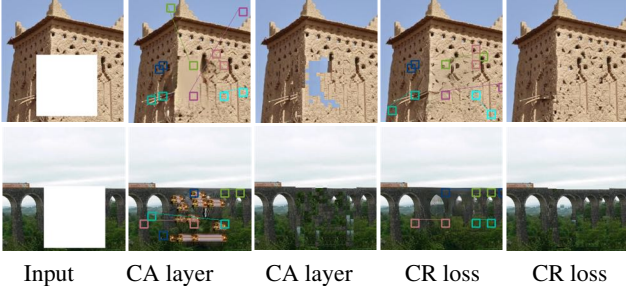| Input | CA layer | CA layer | CR loss | CR loss |

Figure 9: Reference regions selected by the CR loss and CA layer and images pieced by moving each reference patch to the corresponding position in the missing region.

Table 3: L1 error, PSNR, SSIM of the image pieced with reference patches selected by the CR loss and CA layer on the Places2 validation set with square holes.

| Method | Square holes | | | Irregular holes | | |
|---|---|---|---|---|---|---|
| | L1 error | PSNR | SSIM | L1 error | PSNR | SSIM |
| DeepFillv2 (with CA) | .0331 | 21.37 | .8104 | .0232 | 24.39 | .8562 |
| DeepFillv2-CA+CR | .0315 | 21.77 | .8135 | .0205 | 25.27 | .8629 |
| Baseline+CR | .0310 | 21.90 | .8152 | .0202 | 25.39 | .8643 |

Table 4: Comparison with baseline and CA layer on the Places2 validation set.

| Method | Square holes | | | Irregular holes | | |
|---|---|---|---|---|---|---|
| | L1 error | PSNR | SSIM | L1 error | PSNR | SSIM |
| Baseline | .0275 | 22.92 | .8386 | .0171 | 27.03 | .8920 |
| Baseline+CR | .0263 | 23.36 | .8462 | .0157 | 27.79 | .9002 |
| DeepFillv2 (with CA) | .0290 | 22.72 | .8429 | .0178 | 27.03 | .8947 |
| DeepFillv2-CA+CR | .0262 | 23.37 | 0.8447 | .0158 | 27.77 | .8984 |

resolutions. The efficiency advantages become increasingly evident when the resolution increases.
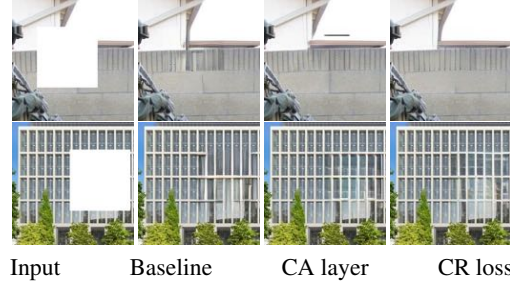


| Input | Baseline | CA layer | CR loss |

Figure 10: Effect of CA layer and CR loss. Both encourage to capture distant relation among image regions.

Table 5: Comparison of time complexity and the running time measured at different resolution. We are not able to evaluate contextual attention layer at $2080 \times 2048$ on GPU due to limited GPU memory size and report the comparison on CPU instead. $N$ is the number of pixels in the input image. $T_s^{GPU}$ and $T_s^{CPU}$ represent the running time (in seconds) on GPU and CPU at $s \times s$ resolution, respectively.

| | Complexity | $T_{512}^{GPU}$ | $T_{1024}^{GPU}$ | $T_{2048}^{GPU}$ | $T_{2048}^{CPU}$ |
|---|---|---|---|---|---|
| CA layer | $O(N^2)$ | 0.063 | 0.371 | Fail | 266 |
| CR loss | $O(N)$ | 0.047 | 0.179 | 0.720 | 28.9 |

## 6 Conclusion

In this paper we proposed a novel contextual reconstruction loss (CR loss) for image inpainting to encourage a CNN to exploit feature patches from known regions as references when inpaint a missing region. CR loss improves the ability of CNNs to capture distance relation among image features and improves both visual quality and quantitative performance of the inpainting results. It brings no extra computation or parameters to the model and can be applied to almost any network architecture. Experiments show that our inpainting model with the proposed CR loss compares favourably against state-of-the-arts in terms of quantitative

measurements and visual quality.

# References

[1] Coloma Ballester, Marcelo Bertalmio, Vicent Caselles, Guillermo Sapiro, and Joan Verdera. Filling-in by joint interpolation of vector fields and gray levels. *IEEE Transactions on Image Processing*, 10(8):1200–1211, 2001. 2

[2] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. 28(3):24, 2009. 1, 2

[3] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *The 27th Annual Conference on Computer Graphics and Interactive Techniques*, 2000. 2

[4] Alexei A Efros and William T Freeman. Image quilting for texture synthesis and transfer. In *The 28th Annual Conference on Computer Graphics and Interactive Techniques*, pages 341–346. ACM, 2001. 1, 2

[5] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014. 1, 2

[6] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):107, 2017. 1, 2

[7] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1125–1134, 2017. 3

[8] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[9] Vivek Kwatra, Irfan Essa, Aaron Bobick, and Nipun Kwatra. Texture optimization for example-based synthesis. 24(3):795–802, 2005. 1, 2

[10] Anat Levin, Assaf Zomet, Shmuel Peleg, and Yair Weiss. Seamless image stitching in the gradient domain. In *European Conference on Computer Vision*, pages 377–389. Springer, 2004. 1

[11] Yijun Li, Sifei Liu, Jimei Yang, and Ming-Hsuan Yang. Generative face completion. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2

[12] Guilin Liu, Fitsum A Reda, Kevin J Shih, Ting-Chun Wang, Andrew Tao, and Bryan Catanzaro. Image inpainting for irregular holes using partial convolutions. In *European Conference on Computer Vision*, 2018. 1

[13] Hongyu Liu, Bin Jiang, Yibing Song, Wei Huang, and Chao Yang. Rethinking image inpainting via a mutual encoder-decoder with feature equalizations. In *European Conference on Computer Vision*, 2020. 3, 6, 7, 8

[14] Hongyu Liu, Bin Jiang, Yi Xiao, and Chao Yang. Coherent semantic attention for image inpainting. In *IEEE International Conference on Computer Vision*, 2019. 2, 3

[15] Eunbyung Park, Jimei Yang, Ersin Yumer, Duygu Ceylan, and Alexander C Berg. Transformation-grounded image generation network for novel 3d view synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3500–3509, 2017. 1

[16] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 1, 2

[17] Yuhang Song, Chao Yang, Zhe Lin, Xiaofeng Liu, Qin Huang, Hao Li, and C-C Jay Kuo. Contextual-based image inpainting: Infer, match, and translate. In *European Conference on Computer Vision (ECCV)*, 2018. 3

[18] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, 2017. 3

[19] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 3

[20] Wei Xiong, Jiahui Yu, Zhe Lin, Jimei Yang, Xin Lu, Connelly Barnes, and Jiebo Luo. Foreground-aware image inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2, 3, 5

[21] Zhaoyi Yan, Xiaoming Li, Mu Li, Wangmeng Zuo, and Shiguang Shan. Shift-net: Image inpainting via deep feature rearrangement. In *European conference on computer vision (ECCV)*, 2018. 3

[22] Chao Yang, Xin Lu, Zhe Lin, Eli Shechtman, Oliver Wang, and Hao Li. High-resolution image inpainting using multi-scale neural patch synthesis. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 2

[23] Zili Yi, Qiang Tang, Shekoofeh Azizi, Daesik Jang, and Zhan Xu. Contextual residual aggregation for ultra high-resolution image inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 6, 7, 8

[24] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Generative image inpainting with contextual attention. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5505–5514, 2018. 1, 2, 4

[25] Jiahui Yu, Zhe Lin, Jimei Yang, Xiaohui Shen, Xin Lu, and Thomas S Huang. Free-form image inpainting with gated convolution. In *IEEE International Conference on Computer Vision*, 2019. 1, 2, 3, 5, 6, 7, 8

[26] Yanhong Zeng, Jianlong Fu, Hongyang Chao, and Baining Guo. Learning pyramid-context encoder network for high-quality image inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1486–1494, 2019. 2, 6, 7, 8

[27] Yu Zeng, Zhe Lin, Jimei Yang, Jianming Zhang, Eli Shechtman, and Huchuan Lu. High-resolution image inpainting with iterative confidence feedback and guided upsampling. In *European Conference on Computer Vision*, 2020. 2, 5